

*Pacific
Journal of
Mathematics*

Volume 257 No. 1

May 2012

PACIFIC JOURNAL OF MATHEMATICS

<http://pacificmath.org>

Founded in 1951 by
E. F. Beckenbach (1906–1982) and F. Wolf (1904–1989)

EDITORS

V. S. Varadarajan (Managing Editor)
Department of Mathematics
University of California
Los Angeles, CA 90095-1555
pacific@math.ucla.edu

Vyjayanthi Chari
Department of Mathematics
University of California
Riverside, CA 92521-0135
chari@math.ucr.edu

Darren Long
Department of Mathematics
University of California
Santa Barbara, CA 93106-3080
long@math.ucsb.edu

Sorin Popa
Department of Mathematics
University of California
Los Angeles, CA 90095-1555
popa@math.ucla.edu

Robert Finn
Department of Mathematics
Stanford University
Stanford, CA 94305-2125
finn@math.stanford.edu

Jiang-Hua Lu
Department of Mathematics
The University of Hong Kong
Pokfulam Rd., Hong Kong
jhlu@maths.hku.hk

Jie Qing
Department of Mathematics
University of California
Santa Cruz, CA 95064
qing@cats.ucsc.edu

Kefeng Liu
Department of Mathematics
University of California
Los Angeles, CA 90095-1555
liu@math.ucla.edu

Alexander Merkurjev
Department of Mathematics
University of California
Los Angeles, CA 90095-1555
merkurev@math.ucla.edu

Jonathan Rogawski
Department of Mathematics
University of California
Los Angeles, CA 90095-1555
jonr@math.ucla.edu

PRODUCTION

pacific@math.berkeley.edu

Silvio Levy, Scientific Editor

Matthew Cargo, Senior Production Editor

SUPPORTING INSTITUTIONS

ACADEMIA SINICA, TAIPEI
CALIFORNIA INST. OF TECHNOLOGY
INST. DE MATEMÁTICA PURA E APLICADA
KEIO UNIVERSITY
MATH. SCIENCES RESEARCH INSTITUTE
NEW MEXICO STATE UNIV.
OREGON STATE UNIV.

STANFORD UNIVERSITY
UNIV. OF BRITISH COLUMBIA
UNIV. OF CALIFORNIA, BERKELEY
UNIV. OF CALIFORNIA, DAVIS
UNIV. OF CALIFORNIA, LOS ANGELES
UNIV. OF CALIFORNIA, RIVERSIDE
UNIV. OF CALIFORNIA, SAN DIEGO
UNIV. OF CALIF., SANTA BARBARA

UNIV. OF CALIF., SANTA CRUZ
UNIV. OF MONTANA
UNIV. OF OREGON
UNIV. OF SOUTHERN CALIFORNIA
UNIV. OF UTAH
UNIV. OF WASHINGTON
WASHINGTON STATE UNIVERSITY

These supporting institutions contribute to the cost of publication of this Journal, but they are not owners or publishers and have no responsibility for its contents or policies.

See inside back cover or pacificmath.org for submission instructions.

The subscription price for 2012 is US \$420/year for the electronic version, and \$485/year for print and electronic. Subscriptions, requests for back issues from the last three years and changes of subscribers address should be sent to Pacific Journal of Mathematics, P.O. Box 4163, Berkeley, CA 94704-0163, U.S.A. Prior back issues are obtainable from Periodicals Service Company, 11 Main Street, Germantown, NY 12526-5635. The Pacific Journal of Mathematics is indexed by Mathematical Reviews, Zentralblatt MATH, PASCAL CNRS Index, Referativnyi Zhurnal, Current Mathematical Publications and the Science Citation Index.

The Pacific Journal of Mathematics (ISSN 0030-8730) at the University of California, c/o Department of Mathematics, 969 Evans Hall, Berkeley, CA 94720-3840, is published monthly except July and August. Periodical rate postage paid at Berkeley, CA 94704, and additional mailing offices. POSTMASTER: send address changes to Pacific Journal of Mathematics, P.O. Box 4163, Berkeley, CA 94704-0163.

PJM peer review and production are managed by EditFLOW™ from Mathematical Sciences Publishers.

PUBLISHED BY PACIFIC JOURNAL OF MATHEMATICS

at the University of California, Berkeley 94720-3840

A NON-PROFIT CORPORATION

Typeset in L^AT_EX

Copyright ©2012 by Pacific Journal of Mathematics

ENERGY AND VOLUME OF VECTOR FIELDS ON SPHERICAL DOMAINS

FABIANO G. B. BRITO, ANDRÉ O. GOMES AND GIOVANNI S. NUNES

We present a “boundary version” for theorems about minimality of volume and energy functionals on a spherical domain of an odd-dimensional Euclidean sphere.

1. Introduction

Let (M, g) be a closed, n -dimensional Riemannian manifold and T^1M the unit tangent bundle of M considered as a closed Riemannian manifold with the Sasaki metric. Let $X : M \rightarrow T^1M$ be a unit vector field defined on M , regarded as a smooth section of the unit tangent bundle T^1M . The volume of X was defined in [Gluck and Ziller 1986] by $\text{vol } X := \text{vol } X(M)$, where $\text{vol } X(M)$ is the volume of the submanifold $X(M) \subset T^1M$. Using an orthonormal local frame $\{e_1, e_2, \dots, e_{n-1}, e_n = X\}$, the volume of the unit vector field X is given by

$$\text{vol } X = \int_M \left(1 + \sum_{a=1}^n \|\nabla_{e_a} X\|^2 + \sum_{a < b} \|\nabla_{e_a} X \wedge \nabla_{e_b} X\|^2 + \dots + \sum_{a_1 < \dots < a_{n-1}} \|\nabla_{e_{a_1}} X \wedge \dots \wedge \nabla_{e_{a_{n-1}}} X\|^2 \right)^{1/2} v_M(g)$$

and the energy of the vector field X is given by

$$\mathcal{E}(X) = \frac{n}{2} \text{vol } M + \frac{1}{2} \int_M \sum_{a=1}^n \|\nabla_{e_a} X\|^2 v_M(g).$$

The Hopf vector fields on \mathbb{S}^{2k+1} are unit vector fields tangent to the classical Hopf fibration $\mathbb{S}^1 \hookrightarrow \mathbb{S}^{2k+1}$. The following theorems gives a characterization of Hopf flows as absolute minima of volume and energy functionals:

Theorem 1 [Gluck and Ziller 1986]. *The unit vector fields of minimum volume on the sphere \mathbb{S}^3 are precisely the Hopf vector fields and no others.*

MSC2010: 53C20.

Keywords: energy of vector fields, volume of vector fields, Hopf flow.

Theorem 2 [Brito 2000]. *The unit vector fields of minimum energy on the sphere \mathbb{S}^3 are precisely the Hopf vector fields and no others.*

We prove in this paper the following boundary version for these theorems:

Theorem 3. *Let U be an open set of the $(2k + 1)$ -dimensional unit sphere \mathbb{S}^{2k+1} and let $K \subset U$ be a connected $(2k + 1)$ -submanifold with boundary of the sphere \mathbb{S}^{2k+1} . Let \vec{v} be a unit vector field on U which coincides with a Hopf flow H along the boundary of K . Then*

$$\mathcal{E}(\vec{v}) \geq \left(\frac{2k+1}{2} + \frac{k}{2k-1} \right) \text{vol } K \quad \text{and} \quad \text{vol } \vec{v} \geq \frac{4^k}{\binom{2k}{k}} \text{vol } K.$$

(Other results for higher dimensions may be found in [Brito et al. 2004; Borrelli and Gil-Medrano 2006; Chacón et al. 2001].)

2. Preliminaries

Let $U \subset \mathbb{S}^{2k+1}$ be an open set of the unit sphere and let $K \subset U$ be a connected $(2k + 1)$ -submanifold with boundary of \mathbb{S}^{2k+1} . Let H be a Hopf vector field on \mathbb{S}^{2k+1} and let \vec{v} be a unit vector field defined on U . We also consider the map $\varphi_t^{\vec{v}} : U \rightarrow \mathbb{S}^{2k+1}(\sqrt{1+t^2})$ given by $\varphi_t^{\vec{v}}(x) = x + t\vec{v}(x)$. This map was introduced in [Asimov 1978; Brito et al. 1981; Milnor 1978].

Lemma 4. *For $t > 0$ sufficiently small, the map $\varphi_t^{\vec{v}}$ is a diffeomorphism.*

Proof. A simple application of the identity perturbation method. \square

From now on, we assume that $t > 0$ is small enough so that the map $\varphi_t^{\vec{v}}$ is a diffeomorphism. In order to find the Jacobian matrix of $\varphi_t^{\vec{v}}$, we define the unit vector field \vec{u} on $\varphi_t^{\vec{v}}(U) \subset \mathbb{S}^{2k+1}(\sqrt{1+t^2})$ by

$$\vec{u}(x) := \frac{1}{\sqrt{1+t^2}} \vec{v}(x) - \frac{t}{\sqrt{1+t^2}} x.$$

Using an adapted orthonormal frame $\{e_1, \dots, e_{2k}, \vec{v}\}$ on a neighborhood V of U , we obtain an adapted orthonormal frame on $\varphi_t^{\vec{v}}(V)$ given by $\{\bar{e}_1, \dots, \bar{e}_{2k}, \vec{u}\}$, where $\bar{e}_i = e_i$ for all $i \in \{1, \dots, 2k\}$.

In this manner, we can write

$$\begin{aligned} d\varphi_t^{\vec{v}}(e_1) &= \langle d\varphi_t^{\vec{v}}(e_1), e_1 \rangle e_1 + \dots + \langle d\varphi_t^{\vec{v}}(e_1), e_{2k} \rangle e_{2k} + \langle d\varphi_t^{\vec{v}}(e_1), \vec{u} \rangle \vec{u}, \\ d\varphi_t^{\vec{v}}(e_2) &= \langle d\varphi_t^{\vec{v}}(e_2), e_1 \rangle e_1 + \dots + \langle d\varphi_t^{\vec{v}}(e_2), e_{2k} \rangle e_{2k} + \langle d\varphi_t^{\vec{v}}(e_2), \vec{u} \rangle \vec{u}, \\ &\vdots \\ d\varphi_t^{\vec{v}}(e_{2k}) &= \langle d\varphi_t^{\vec{v}}(e_{2k}), e_1 \rangle e_1 + \dots + \langle d\varphi_t^{\vec{v}}(e_{2k}), e_{2k} \rangle e_{2k} + \langle d\varphi_t^{\vec{v}}(e_{2k}), \vec{u} \rangle \vec{u}, \\ d\varphi_t^{\vec{v}}(\vec{v}) &= \langle d\varphi_t^{\vec{v}}(\vec{v}), e_1 \rangle e_1 + \dots + \langle d\varphi_t^{\vec{v}}(\vec{v}), e_{2k} \rangle e_{2k} + \langle d\varphi_t^{\vec{v}}(\vec{v}), \vec{u} \rangle \vec{u}. \end{aligned}$$

Now, by Gauss's equation for the trivial immersion $\mathbb{S}^{2k+1} \hookrightarrow \mathbb{R}^{2k+2}$, we have

$$\tilde{\nabla}_Y \vec{v} = d\vec{v}(Y) = \nabla_Y \vec{v} - \langle \vec{v}, Y \rangle x$$

for every vector field Y on \mathbb{S}^{2k+1} , and then

$$\langle d\varphi_t^{\vec{v}}(e_1), e_1 \rangle = \langle e_1 + td\vec{v}(e_1), e_1 \rangle = 1 + t\langle \nabla_{e_1} \vec{v}, e_1 \rangle$$

Analogously, we can conclude that

$$\begin{aligned} \langle d\varphi_t^{\vec{v}}(e_i), e_i \rangle &= 1 + t\langle \nabla_{e_i} \vec{v}, e_i \rangle && \text{for } i \in \{1, \dots, 2k\}, \\ \langle d\varphi_t^{\vec{v}}(e_i), e_j \rangle &= t\langle \nabla_{e_i} \vec{v}, e_j \rangle && \text{for } i, j \in \{1, \dots, 2k\}, i \neq j, \\ \langle d\varphi_t^{\vec{v}}(e_i), \vec{u} \rangle &= 0 && \text{for } i \in \{1, \dots, 2k\}, \\ \langle d\varphi_t^{\vec{v}}(\vec{v}), \vec{u} \rangle &= \sqrt{1 + t^2}. \end{aligned}$$

By employing the notation $h_{ij}(\vec{v}) := \langle \nabla_{e_i} \vec{v}, e_j \rangle$ (where $i, j \in \{1, \dots, 2k\}$), we can express the determinant of the Jacobian matrix of $\varphi_t^{\vec{v}}$ in the form

$$\det(d\varphi_t^{\vec{v}}) = \sqrt{1 + t^2} \left(1 + \sum_{i=1}^{2k} \sigma_i(\vec{v}) t^2 \right),$$

where, by definition, the functions σ_i are the i -symmetric functions of the h_{ij} . For instance, if $k = 1$, we have

$$\begin{aligned} \sigma_1(\vec{v}) &:= h_{11}(\vec{v}) + h_{22}(\vec{v}), \\ \sigma_2(\vec{v}) &:= h_{11}(\vec{v})h_{22}(\vec{v}) - h_{12}(\vec{v})h_{21}(\vec{v}). \end{aligned}$$

3. Proof of the Theorem

The energy of the vector field \vec{v} (on K) is given by

$$\mathfrak{E}(\vec{v}) := \frac{1}{2} \int_K \|d\vec{v}\|^2 = \frac{2k+1}{2} \text{vol } K + \frac{1}{2} \int_K \|\nabla \vec{v}\|^2$$

Using the notation above, we have

$$\mathfrak{E}(\vec{v}) = \frac{2k+1}{2} \text{vol } K + \frac{1}{2} \int_K \left(\sum_{i,j=1}^{2k} (h_{ij}(\vec{v}))^2 + \sum_{i=1}^{2k} \langle \nabla_{\vec{v}} \vec{v}, e_i \rangle^2 \right)$$

and then

$$(1) \quad \mathfrak{E}(\vec{v}) \geq \frac{2k+1}{2} \text{vol } K + \frac{1}{2} \int_K \sum_{i,j=1}^{2k} (h_{ij}(\vec{v}))^2.$$

Now observe that

$$\sum_{i < j} (h_{ii} - h_{jj})^2 = (2k - 1) \sum_i h_{ii}^2 - 2 \sum_{i < j} h_{ii} h_{jj}$$

and

$$\sum_{i < j} (h_{ij} + h_{ji})^2 = \sum_{i \neq j} h_{ij}^2 + 2 \sum_{i < j} h_{ij} h_{ji}.$$

If we sum these last two equations, we get

$$(2k - 1) \sum_i h_{ii}^2 + \sum_{i \neq j} h_{ij}^2 \geq 2\sigma_2$$

and then

$$(2) \quad \sum_i h_{ii}^2 + \frac{1}{2k - 1} \sum_{i \neq j} h_{ij}^2 \geq \frac{2}{2k - 1} \sigma_2.$$

Also, we can write

$$\sum_{i,j=1}^{2k} h_{ij}^2 = \sum_{i \neq j} h_{ij}^2 + \sum_i h_{ii}^2 \geq \sum_i h_{ii}^2 + \frac{1}{2k - 1} \sum_{i \neq j} h_{ij}^2.$$

From this and (2), we obtain

$$\sum_{i,j=1}^{2k} (h_{ij}(\vec{v}))^2 \geq \frac{2}{2k - 1} \sigma_2(\vec{v}).$$

But then, using inequality (1), we find that

$$(3) \quad \mathfrak{E}(\vec{v}) \geq \frac{2k + 1}{2} \text{vol } K + \frac{1}{2k - 1} \int_K \sigma_2(\vec{v}).$$

On the other hand, by the change of variables theorem, we obtain

$$\text{vol } \varphi_t^H(K) = \int_K \sqrt{1 + t^2} \left(1 + \sum_{i=1}^{2k} \sigma_i(H) t^i\right)$$

By a straightforward computation shown in [Chacón 2000] and [Brito et al. 2004], we have $\sigma_i(H) = \eta_i$ for all $i \in \{1, \dots, 2k\}$, where

$$\eta_i = \begin{cases} \binom{k}{i/2} & \text{if } i \text{ is even,} \\ 0 & \text{if } i \text{ is odd.} \end{cases}$$

We know that the vector fields \vec{v} and H are the same on ∂K . Thus, $\varphi_t^{\vec{v}}(K)$ and $\varphi_t^H(K)$ are $(2k + 1)$ -submanifolds of $\mathbb{S}^{2k+1}(\sqrt{1 + t^2})$ with the same boundary. We

claim that $\varphi_t^{\vec{v}}(K) = \varphi_t^H(K)$ for all t sufficiently small. In fact, if p is an interior point of K ,

$$\lim_{t \rightarrow 0} \varphi_t^{\vec{v}}(p) = \lim_{t \rightarrow 0} \varphi_t^H(p) = p$$

and then we have necessarily

$$\varphi_t^{\vec{v}}(K) = \varphi_t^H(K)$$

for all t sufficiently small; equivalently,

$$\int_K \sqrt{1+t^2} \left(1 + \sum_{i=1}^{2k} \sigma_i(\vec{v}) t^i \right) = \int_K \sqrt{1+t^2} \left(1 + \sum_{i=1}^{2k} \eta_i t^i \right)$$

for all $t > 0$ sufficiently small. Consequently, after canceling the factor $\sqrt{1+t^2}$ and rearranging the terms, we obtain

$$\left(\int_K [\sigma_1(\vec{v}) - \eta_1] \right) t + \left(\int_K [\sigma_2(\vec{v}) - \eta_2] \right) t^2 + \dots + \left(\int_K [\sigma_{2k}(\vec{v}) - \eta_{2k}] \right) t^{2k} = 0$$

for all sufficiently small t . By identity of polynomials, we conclude

$$\int_K \sigma_i(\vec{v}) = \int_K \eta_i = \eta_i \text{ vol } K \quad \text{for } i \in \{1, \dots, 2k\}.$$

Using this (for $i = 2$) together with (3), we get

$$\mathcal{E}(\vec{v}) \geq \frac{2k+1}{2} \text{ vol } K + \frac{\eta_2}{2k-1} \text{ vol } K = \left(\frac{2k+1}{2} + \frac{k}{2k-1} \right) \text{ vol } K.$$

We can obtain an analogue of this result for volumes using the following inequality (see [Brito et al. 2004] or [Chacón 2000, page 59]):

$$\text{vol } \vec{v} \geq \int_K \left(1 + \sum_{i=1}^k \frac{\binom{k}{i}}{\binom{2k}{2i}} \sigma_{2i}(\vec{v}) \right).$$

But $\int_K \sigma_{2i} = \int_K \eta_{2i} = \eta_{2i} \text{ vol } K$ for all $i \in \{1, \dots, k\}$. Then, we have

$$\text{vol } \vec{v} \geq \left(1 + \sum_{i=1}^k \frac{\binom{k}{i}^2}{\binom{2k}{2i}} \right) \text{ vol } K \geq \frac{4^k}{\binom{2k}{k}} \text{ vol } K$$

4. Final remarks

- (1) If K is a spherical cap (the closure of a connected open set with round boundary of the three unit sphere), the theorem provides a “boundary version” for

the minimalization theorem of energy and volume functionals on [Brito 2000] and [Gluck and Ziller 1986].

- (2) The “Hopf boundary” hypothesis is essential. In fact, if there is no constraint for the unit vector field \vec{v} on ∂K , it is possible to construct vector fields on “small caps” such that $\|\nabla\vec{v}\|$ is small on K (exponential maps may be used on that construction). A consequence of this is that $\mathcal{E}(\vec{v})$ and $\text{vol } \vec{v}$ are less than volume and energy of Hopf vector fields respectively.

Acknowledgements

We express our gratitude to Prof. Jaime Ripoll for helpful conversation concerning the final draft of our paper.

References

- [Asimov 1978] D. Asimov, “Average Gaussian curvature of leaves of foliations”, *Bull. Amer. Math. Soc.* **84**:1 (1978), 131–133. MR 0464257 (57 #4191)
- [Borrelli and Gil-Medrano 2006] V. Borrelli and O. Gil-Medrano, “A critical radius for unit Hopf vector fields on spheres”, *Math. Ann.* **334**:4 (2006), 731–751. MR 2209254 (2007a:53070)
- [Brito 2000] F. G. B. Brito, “Total bending of flows with mean curvature correction”, *Differential Geom. Appl.* **12**:2 (2000), 157–163. MR 1758847 (2001g:53065)
- [Brito et al. 1981] F. Brito, R. Langevin, and H. Rosenberg, “Intégrales de courbure sur des variétés feuilletées”, *J. Differential Geom.* **16**:1 (1981), 19–50. MR 633622 (83a:57032)
- [Brito et al. 2004] F. B. Brito, P. M. Chacón, and A. M. Naveira, “On the volume of unit vector fields on spaces of constant sectional curvature”, *Comment. Math. Helv.* **79**:2 (2004), 300–316. MR 2059434 (2005f:53042)
- [Chacón 2000] P. M. Chacón, *Sobre a energia e energia corrigida de campos unitários e distribuições. Volume de campos unitários*, PhD thesis, Universidade de São Paulo, 2000.
- [Chacón et al. 2001] P. M. Chacón, A. M. Naveira, and J. M. Weston, “On the energy of distributions, with application to the quaternionic Hopf fibrations”, *Monatsh. Math.* **133**:4 (2001), 281–294. MR 1915876 (2003k:53050)
- [Gluck and Ziller 1986] H. Gluck and W. Ziller, “On the volume of a unit vector field on the three-sphere”, *Comment. Math. Helv.* **61**:2 (1986), 177–192. MR 856085 (87j:53063)
- [Milnor 1978] J. Milnor, “Analytic proofs of the “hairy ball theorem” and the Brouwer fixed-point theorem”, *Amer. Math. Monthly* **85**:7 (1978), 521–524. MR 505523 (80m:55001)

Received April 20, 2011. Revised April 20, 2012.

FABIANO G. B. BRITO
 DEPARTAMENTO DE MATEMÁTICA E ESTATÍSTICA
 UNIVERSIDADE FEDERAL DO ESTADO DO RIO DE JANEIRO
 22290-240 RIO DE JANEIRO RJ
 BRAZIL
 brifabiano@gmail.com

ANDRÉ O. GOMES
INSTITUTO DE MATEMÁTICA E ESTATÍSTICA
UNIVERSIDADE DE SÃO PAULO
05508-090 SÃO PAULO SP
BRAZIL
gomes@ime.usp.br

GIOVANNI S. NUNES
INSTITUTO DE FÍSICA E MATEMÁTICA
UNIVERSIDADE FEDERAL DE PELOTAS
96001-970 PELOTAS RS
BRAZIL
giovanni.nunes@ufpel.edu.br

MAPS ON 3-MANIFOLDS GIVEN BY SURGERY

BOLDIZSÁR KALMÁR AND ANDRÁS I. STIPSICZ

Suppose that the 3-manifold M is given by integral surgery along a link $L \subset S^3$. In the following we construct a stable map from M to the plane, whose singular set is canonically oriented. We obtain upper bounds for the minimal numbers of crossing singularities, nonsimple singularities, and connected components of fibers of stable maps from M to the plane in terms of properties of L .

1. Introduction

It is well-known that a continuous map between smooth manifolds can be approximated by a smooth map and any smooth map on a 3-manifold can be approximated by a generic stable map. This line of argument, however, gives no concrete map on a given 3-manifold M even if it is given by some explicit construction. Recall that by [Lickorish 1962; Wallace 1960] a closed oriented 3-manifold M can be given by integral surgery along some link L in S^3 . In the present work we construct an explicit stable map $F : M \rightarrow \mathbb{R}^2$ based on such a surgery presentation of M .

Results of Gromov [2009; 2010] give lower bounds on the topological complexity of the set of critical values of generic smooth maps and on the complexity of the fibers in terms of the topology of the source and target manifolds. In a slightly different direction, [Costantino and Thurston 2008] gives a lower bound for the number of crossing singularities of stable maps from a 3-manifold to \mathbb{R}^2 in terms of the Gromov norm of the 3-manifold. Recently Baykur [2008; 2009] and Gay and Kirby [2007] studied the topology of 4-manifolds through the singularities of their maps into surfaces.

In the present paper we give upper bounds on the minimal numbers of the crossing and nonsimple singularities and of the connected components of the fibers of stable maps on the 3-manifold M in terms of properties of diagrams of L (e.g., the number of crossings or the number of critical points when projected to \mathbb{R}). As an additional result, these constructions lead to upper bounds on a version of the Thurston–Bennequin number of negative Legendrian knots.

MSC2010: primary 57R45; secondary 57M27.

Keywords: stable map, 3-manifold, surgery, negative knot, Thurston–Bennequin number.

Before stating our main results, we need a little preparation. First of all, a stable map of a 3-manifold into the plane can be easily described by its Stein factorization.

Definition 1.1. Let F be a map of the 3-manifold M into \mathbb{R}^2 . Let us call two points $p_1, p_2 \in M$ equivalent if and only if p_1 and p_2 lie on the same component of an F -fiber. Let W_F denote the quotient space of M with respect to this equivalence relation and $q_F : M \rightarrow W_F$ the quotient map. Then there exists a unique continuous map $\bar{F} : W_F \rightarrow \mathbb{R}^2$ such that $F = \bar{F} \circ q_F$. The space W_F or the factorization of the map F into the composition of q_F and \bar{F} is called the *Stein factorization* of the map F . (Sometimes the map \bar{F} is also called the Stein factorization of F .)

In other words, the Stein factorization W_F is the space of connected components of fibers of F . Its structure is strongly related to the topology of the 3-manifold M . For example, an immediate observation is that the quotient map $q_F : M \rightarrow W_F$ induces an epimorphism between the fundamental groups since every loop in W_F can be lifted to M . If $F : M \rightarrow \mathbb{R}^2$ is a stable map, then its Stein factorization W_F is a 2-dimensional CW complex. The local forms of Stein factorizations of proper stable maps of orientable 3-manifolds into surfaces are described in [Kushner et al. 1984; Levine 1985]; see Figure 1. Indeed, let F be a stable map of the closed orientable 3-manifold M into \mathbb{R}^2 . We say that a singular point $p \in M$ of F is of type (A), ..., (E), respectively, if the Stein factorization \bar{F} at $q_F(p)$ looks locally like (a), ..., (e) of Figure 1, respectively. We will call a point $w \in W_F$ a singular point of type (A), ..., (E), respectively, if $w = q_F(p)$ for a singular point $p \in M$ of type (A), ..., (E), respectively. According to [Kushner et al. 1984; Levine 1985] we give the following characterization of the singularities of F : The singular point p is a *cusplike* point if and only if it is of type (C), the singular point p is a *definite fold* point if and only if it is of type (A) and p is an *indefinite fold* point if and only if it is of type (B), (D) or (E). Singular points of types (D) and (E) are called *nonsimple*, while the others are called *simple*. A double point in \mathbb{R}^2 of two crossing

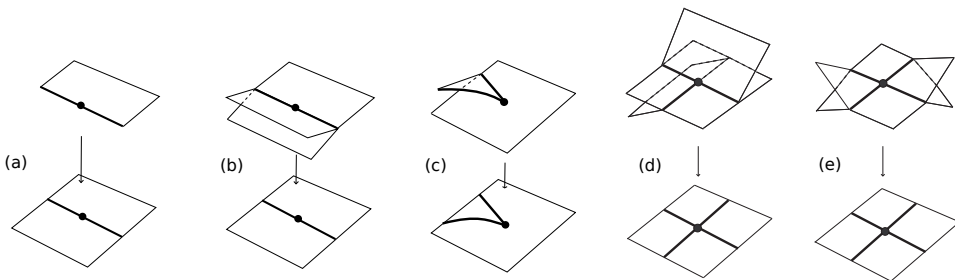


Figure 1. The local forms of Stein factorizations of stable maps from orientable 3-manifolds to surfaces. The map (symbolized by an arrow) maps from the CW complex W_F to \mathbb{R}^2 .

images of singular curves which is not an image of a nonsimple singularity is called a *simple singularity crossing*. A simple singularity crossing or an image in \mathbb{R}^2 of a nonsimple singularity is called a *crossing singularity*. A stable map is called a *fold map* if it has no cusp singularities.

Let $L \subset \mathbb{R}^3 \subset S^3$ be a given link, and let \bar{L} denote a generic projection of it to the plane. Let $n(L)$ and $\text{cr}(\bar{L})$ denote the number of components of L and the number of crossings of \bar{L} , respectively.

Choose a direction in \mathbb{R}^2 , which we represent by a vector $v \in \mathbb{R}^2$. We can assume that v satisfies the condition that the projection of the diagram \bar{L} to $\mathbb{R}v^\perp$ along v yields only non-degenerate critical points. Let $t(\bar{L}) = t_v(\bar{L})$ denote the number of times \bar{L} is tangent to v . Suppose at each v -tangency p the half line emanating from p in the direction of v avoids the crossings of \bar{L} and intersects \bar{L} transversally (at the points different from p). Denote the number of transversal intersections by $\ell(\bar{L}, v, p)$. Let $\ell(\bar{L}, v)$ denote the maximum of the values $\ell(\bar{L}, v, p)$, where p runs over the v -tangencies. With these definitions in place now we can state the main result of the paper.

Theorem 1.2. *Suppose that the 3-manifold M is obtained by integral surgery on the link $L \subset S^3$. Then there is a stable map $F : M \rightarrow \mathbb{R}^2$ such that*

- (1) *the Stein factorization W_F is homotopy equivalent to the bouquet $\bigvee_{i=1}^{n(L)} S^2$,*
- (2) *the number of cusps of F is equal to $t_v(\bar{L})$,*
- (3) *all the nonsimple singularities of F are of type (D), and their number is equal to $\text{cr}(\bar{L}) + \frac{3}{2}t_v(\bar{L}) - n(L)$,*
- (4) *the number of nonsimple singularities which are not connected by any singular arc of type (B) to any cusp is equal to $\text{cr}(\bar{L}) + \frac{1}{2}t_v(\bar{L}) - n(L)$,*
- (5) *the number of simple singularity crossings of F in \mathbb{R}^2 is no more than*

$$8 \text{cr}(\bar{L}) + 6\ell(\bar{L}, v)t_v(\bar{L}) + t_v(\bar{L})^2,$$

- (6) *the number of connected components of the singular set of F is no more than $n(L) + \frac{3}{2}t_v(\bar{L}) + 1$, and*
- (7) *the maximal number of the connected components of any fiber of F is no more than $t_v(\bar{L}) + 3$.*
- (8) *Suppose we got M by cutting out and gluing back the regular neighborhood N_L of L from S^3 . Then the indefinite fold singular set of F contains a link in $S^3 - N_L$, which is isotopic to L in S^3 and whose F -image coincides with \bar{L} .*

Remarks 1.3. (1) Let Y be a closed orientable 3-manifold, f a given smooth map of Y into \mathbb{R}^2 and $L \subset Y$ a link disjoint from the singular set of f . Suppose furthermore that $f|_L$ is an immersion. Let M denote the 3-manifold obtained

by some integral surgery along L . Then the method developed in the proof of Theorem 1.2 provides a stable map of M into \mathbb{R}^2 (relative to f).

(2) In constructing the map F , the proof of Theorem 1.2 provides a sequence of stable maps f_0, f_1, \dots, f_6 of S^3 into \mathbb{R}^2 , where each f_i is obtained from f_{i-1} by some deformation, $i = 1, \dots, 6$. Finally, the map F is obtained from f_6 . Suppose that X is a compact 4-manifold which admits a handle decomposition with only 0- and 2-handles; i.e., X can be given by attaching 4-dimensional 2-handles to D^4 along S^3 . Using our method we can construct a stable map G of X into $\mathbb{R}^2 \times [0, 1]$.

Recall that according to [Burlet and de Rham 1974] a closed orientable 3-manifold M has a stable map into \mathbb{R}^2 without singularities of types (B), (C), (D) and (E) if and only if M is a connected sum of finitely many copies of $S^1 \times S^2$. According to [Saeki 1996] a closed orientable 3-manifold M has a stable map into \mathbb{R}^2 without singular points of types (C), (D) and (E) if and only if M is a graph manifold. By [Levine 1965] a 3-manifold always has a stable map into \mathbb{R}^2 without singular points of type (C). Our arguments imply a constructive proof for

Theorem 1.4. *Every closed orientable 3-manifold has a stable map into \mathbb{R}^2 without singular points of types (C) and (E).*

Remarks 1.5. (1) One cannot expect to eliminate the singular points of types (A), (B) or (D) of stable maps from arbitrary closed orientable 3-manifolds to \mathbb{R}^2 . In this sense our Theorem 1.4 gives the best possible elimination on 3-manifolds.

(2) By taking an embedding $\mathbb{R}^2 \subset S^2$ we get for every closed orientable 3-manifold a stable map into S^2 as well without singular points of types (C) and (E). Then by using the method of [Saeki 2006], for example, for eliminating the singular points of type (A), we get a stable map, which is a direct analogue of the indefinite generic maps appearing in [Baykur 2008; 2009; Gay and Kirby 2007].

The construction also implies certain relations between quantities one can naturally associate to stable maps and to surgery diagrams.

Definition 1.6. Suppose that M is a fixed closed, oriented 3-manifold and that $F : M \rightarrow \mathbb{R}^2$ is a stable map with singular set Σ .

- Let $s(F)$ denote the number of simple singularity crossings of F .
- Let $ns(F)$ denote the number of nonsimple singularities of F .
- Let $d(F)$ denote the number of crossing singularities of F . Clearly $s(F) + ns(F) = d(F)$.
- Let $nsnc(F)$ denote the number of nonsimple singularities of F which are not connected by any singular arc of type (B) to any cusp.
- Let $c(F)$ denote the number of cusps of F . Clearly $nsnc(F) + c(F) \geq ns(F)$.

- Let $cc(F)$ denote the number of connected components of $F(\Sigma)$. Clearly it is no more than the number of connected components of Σ .
- Let $cf(F)$ denote the maximum number of connected components of the fibers of F .

The inequality

$$\text{rank } H_*(M) \leq 2d(F) + c(F) + 2cc(F)$$

has been shown to hold in [Gromov 2009, Section 2.1].¹ In addition, by [Costantino and Thurston 2008, Theorem 3.38] we have $d(F) \geq \|M\|/10$, where $\|M\|$ is the Gromov norm of M ; see also [Gromov 2009, Section 3].

Theorem 1.2 provides several estimates for upper bounds on the topological complexity of smooth maps of a 3-manifold given by surgery. For example, by summing quantities in Definiton 1.6 and their estimates in Theorem 1.2, we immediately obtain

Corollary 1.7. *Suppose that the 3-manifold M is obtained by integral surgery on the link $L \subset S^3$. Let \bar{L} be any diagram of L and v a general position vector in \mathbb{R}^2 . Then*

- $\min d(F) \leq 9cr(\bar{L}) + (6\ell(\bar{L}, v) + \frac{3}{2})t_v(\bar{L}) + t_v(\bar{L})^2 - n(L)$,
- $\min cf(F) \leq t_v(\bar{L}) + 3$,
- $\min\{2d(F) + c(F) + 2cc(F)\} \leq 18cr(\bar{L}) + (12\ell(\bar{L}, v) + 7)t_v(\bar{L}) + 2t_v(\bar{L})^2 + 2$,

where the minima are taken for all the stable maps F of M into \mathbb{R}^2 . Evidently, we can estimate other properties in Definiton 1.6 of stable maps on M as well.

These expressions can be simplified by estimating $\ell(\bar{L}, v)$ as

$$(1-1) \quad \ell(\bar{L}, v) \leq t_v(\bar{L}) - 1;$$

see Lemma 3.7.

The number of tangencies of a projection of a knot in a fixed direction is reminiscent to the number of cusp singularities of a front projection of a Legendrian knot in the standard contact 3-space. Based on this analogy, our previous results imply an estimate on a quantity attached to a Legendrian knot in the following way.

Recall first that the standard contact structure ξ_{st} on \mathbb{R}^3 is the 2-plane field given by the kernel of the 1-form $\alpha = dz + xdy$. A knot \mathcal{L} is *Legendrian* if the tangent vectors of \mathcal{L} are in ξ_{st} . (To indicate the Legendrian structure on the knot, we will denote it by \mathcal{L} and reserve the notation L for smooth knots and links.) If \mathcal{L} is chosen generically within its Legendrian isotopy class, its projection to the (y, z) plane will have no vertical tangencies, and at any crossing the strand with smaller slope will

¹The paper [Motta et al. 1995] is also closely related.

be over the one with higher slope. Consider now a Legendrian knot \mathcal{L} and let $\bar{\mathcal{L}}$ denote such a projection (called a *front projection*) of \mathcal{L} . The *Thurston–Bennequin number* $\text{tb}(\mathcal{L})$ of \mathcal{L} is given by the formula $w(\bar{\mathcal{L}}) - \frac{1}{2}\#\text{cusps}(\bar{\mathcal{L}})$, where $w(\bar{\mathcal{L}})$ stands for the *writhe* (i.e., the signed sum of the double points) of the projection. Although the definition of $\text{tb}(\mathcal{L})$ uses a projection of the Legendrian knot \mathcal{L} , it is not hard to show that the resulting number is an invariant of the Legendrian isotopy class of \mathcal{L} .

If the projection has only negative crossings, we have that $w(\bar{\mathcal{L}}) = -\text{cr}(\bar{\mathcal{L}})$, hence the resulting Thurston–Bennequin number can be identified with $-\text{cr}(\bar{\mathcal{L}}) - \frac{1}{2}t_v(\bar{\mathcal{L}})$ after choosing v appropriately; cf. [Geiges 2008; Ozbagci and Stipsicz 2004]. (In this case the generic projection \bar{L} used in the definitions of $t_v(\bar{L})$ and $\text{cr}(\bar{L})$ is derived from the front projection $\bar{\mathcal{L}}$ by rounding the cusps.)

As it is customary, we define $\text{TB}(L)$ as the maximum of all Thurston–Bennequin numbers of Legendrian knots smoothly isotopic to L . (It is a nontrivial fact, and follows from the tightness of ξ_{st} that this maximum exists.) A modification of this definition for negative knots (i.e., for knots admitting projections with only negative crossings) provides

Definition 1.8. For a negative knot $L \subset \mathbb{R}^3$ let $\text{TB}^-(L)$ denote the value $\max\{\text{tb}(\mathcal{L})\}$ where \mathcal{L} runs over those Legendrian knots smoothly isotopic to L which admit front diagrams with only negative crossings.

It is rather easy to see that if the knot L admits a projection with only negative crossings, then it also has a front projection with the same property. Clearly $\text{TB}^-(L) \leq \text{TB}(L)$.

Theorem 1.9. For a negative knot $L \subset \mathbb{R}^3$ and any 3-manifold M obtained by an integral surgery along L we have

$$\text{TB}^-(L) \leq -\min \frac{\sqrt{s(F)}}{2\sqrt{7}},$$

$$\text{TB}^-(L) \leq -\min \frac{\sqrt{d(F)}}{2\sqrt{7}},$$

$$\text{TB}^-(L) \leq -\min \text{nsc}(F) - 1,$$

where the minima are taken for all the stable maps F of M into \mathbb{R}^2 .

By Theorem 1.9 and [Costantino and Thurston 2008, Theorem 3.38] we obtain:

Corollary 1.10. For a negative knot $L \subset \mathbb{R}^3$ and any 3-manifold M obtained by an integral surgery along L , we have

$$\text{TB}^-(L) \leq -\frac{\sqrt{\|M\|}}{2\sqrt{70}}.$$

2. Preliminaries

In this section, we recall and summarize some technical tools. First, we show that a cusp can be pushed through an indefinite fold arc as in Figure 2.

Lemma 2.1 (moving cusps). *Suppose that in a neighborhood U of a point $p \in M$ the Stein factorization of a map $f : M \rightarrow \mathbb{R}^2$ is given by Figure 2(a). Then f can be deformed in this neighborhood to a map f' so that the Stein factorization of f' is as the diagram of Figure 2(b).*

Proof. Suppose $q \in M$ is the cusp singular point and $\alpha \subset M$ is the indefinite fold arc at hand. Let $x \in \mathbb{R}^2$ be a point on the other side of $f(\alpha)$ in $f(U)$. Connect $f(q)$ and x by an embedded arc β' . Then there is an arc $\beta \subset M$ such that $f(\beta) = \beta'$, β starts at q , and β and α do not intersect. By using the technique of [Levine 1965] we can now deform f in a small tubular neighborhood of β to achieve the claimed map f' . Note that during this move one singular point of type (D) appears. \square

An analogous statement holds if we move a cusp from a 1-sheeted region to a 2-sheeted region.

According to the next result, two cusps can be eliminated as in Figure 3.

Lemma 2.2 (eliminating cusps). *Suppose that in a neighborhood U of a point $p \in M$ the Stein factorization of a map $f : M \rightarrow \mathbb{R}^2$ is given by Figure 3(a). Then f can be deformed in this neighborhood to a map f' so that the Stein factorization of f' is as the diagram of Figure 3(b).*

Proof. This statement is the elimination in [Levine 1965, pages 285–295] for 3-dimensional source manifolds. \square

Recall that if $f : M \rightarrow \mathbb{R}^2$ is a stable map and $S_f \subset M$ denotes its singular set, then $f|_{S_f}$ is a generic immersion with cusps; i.e., if $C_f \subset M$ denotes the set of

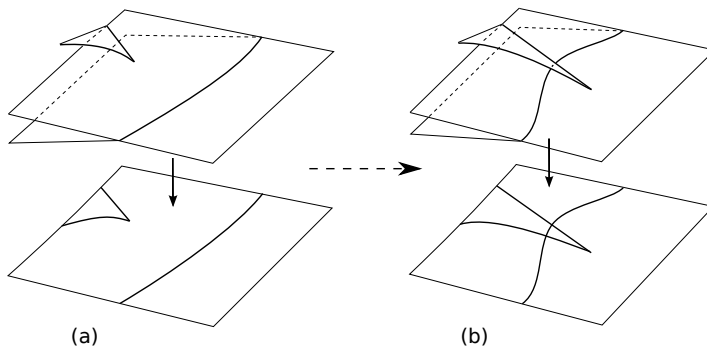


Figure 2. Moving cusps. A map can be deformed so that the image of a cusp point goes to the other side of the image of an indefinite fold arc.

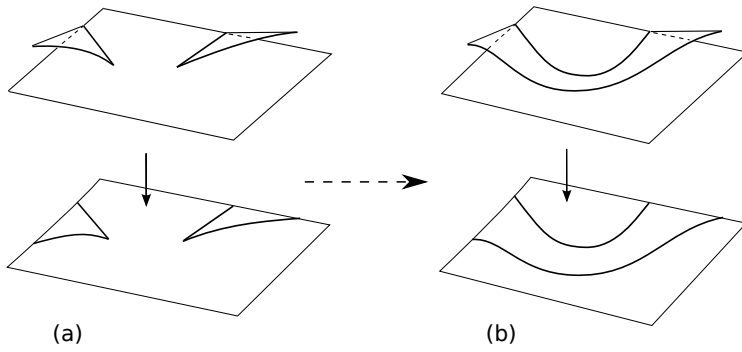


Figure 3. Eliminating cusps.

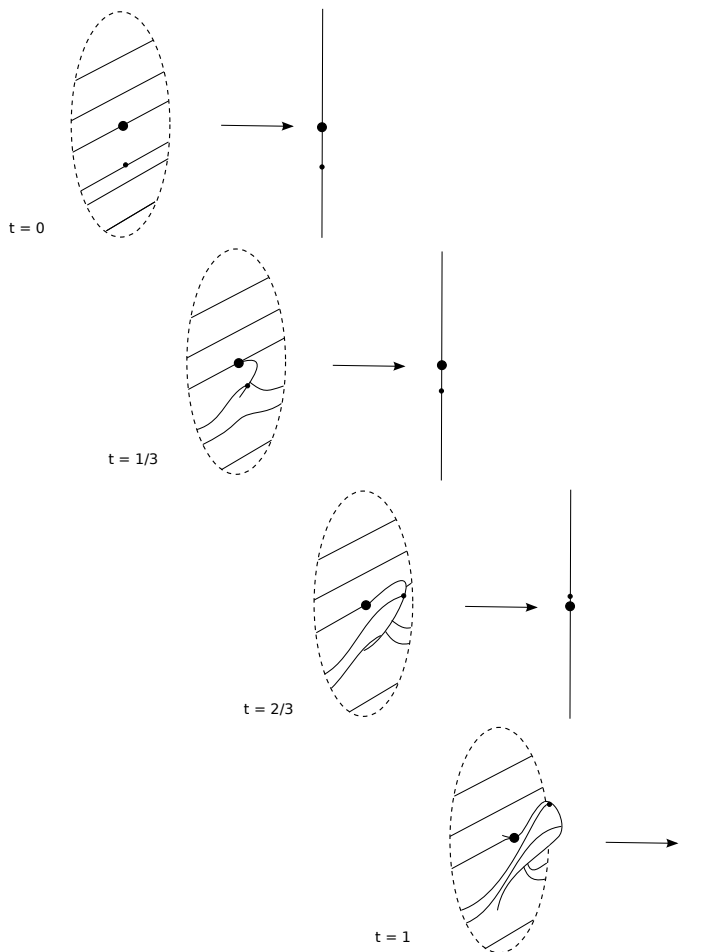


Figure 4. The deformation of f to f' in a fiber of N_L .

cusp points, then $f|_{S_f - C_f}$ is a generic immersion with finitely many double points and $f|_{C_f}$ is disjoint from $f|_{S_f - C_f}$.

The following result will be the key ingredient in our subsequent arguments for proving Theorem 1.2.

Lemma 2.3 (making wrinkles). *Suppose that $f : M \rightarrow \mathbb{R}^2$ is a stable map and let $L \subset M$ denote an embedded closed 1-dimensional manifold such that L is disjoint from the singular set S_f , $f|_L$ is a generic immersion and $f|_{L \cup S_f}$ is a generic immersion with cusps. Let N_L be a small tubular neighborhood of L disjoint from S_f and fix an identification of N_L with the normal bundle of L . Let $s : L \rightarrow N_L$ be a non-zero section such that $f(s(x)) \neq f(x)$ for any $x \in L$. Then f is homotopic to a smooth stable map f' such that*

- (1) $f = f'$ outside N_L ,
- (2) the singular set of f' is $S_f \cup L \cup s(L)$,
- (3) f' has indefinite fold singularities along L ,
- (4) f' has definite fold singularities along $s(L)$,
- (5) $f'|_L = f|_L$,
- (6) $f'|_{s(L)}$ is an immersion parallel to $f|_L$ and
- (7) if for a double point y of $f|_L$ the two points in $f^{-1}(y) \cap L$ lie in the same connected component of the fiber $f^{-1}(y)$, then the double point y of $f'|_L$ correspond to a singularity of type (D).

Proof. We perform the homotopy inside N_L fiberwise as shown by Figure 4 (see previous page). Since N_L is the trivial bundle, the homotopy of the fibers yields a homotopy of the entire N_L . \square

Remark 2.4. If the submanifold L has boundary, we can still get something similar. In this case the section s should be zero at the boundary points of L , and the homotopy yields a stable map f' with cusps at ∂L .

3. Construction of the stable map on M

Proof of Theorem 1.2. We will prove the theorem by presenting an algorithm which produces the map F on M with the desired properties. This algorithm will be given in seven steps; the first six of these steps are concerned with maps on S^3 . Let us start with a fold map $f_0 : S^3 \rightarrow \mathbb{R}^2$ with one unknotted circle $C \subset S^3$ as singular set such that $f_0|_C$ is an embedding and $f_0^{-1}(p)$ is a circle for each regular point $p \in f_0(S^3)$. Then the Stein factorization of f_0 is a disk together with its embedding into \mathbb{R}^2 . By cutting out the interior of a sufficiently small tubular neighborhood N_C of C from S^3 , we get a solid torus $S^3 - \text{int } N_C$ whose boundary is mapped into \mathbb{R}^2 by f_0 as a circle fibration over a circle parallel to $f_0(C)$, and $f_0|_{S^3 - \text{int } N_C}$

is a trivial circle bundle $S^1 \times D^2 \rightarrow D^2$. Suppose the link $L \subset S^3$ is disjoint from $N_C \cup \{1\} \times D^2$. Then by identifying $S^3 - (N_C \cup \{1\} \times D^2)$ with \mathbb{R}^3 and $f_0|_{S^3 - (N_C \cup \{1\} \times D^2)}$ with the projection onto \mathbb{R}^2 , we get a link diagram $\bar{L} = f_0(L)$. Now we start modifying this map f_0 . In Steps 1 through 6 we will deal with maps on S^3 , and the goal will be to obtain a map which is suitable with respect to the fixed surgery link L . In particular, we aim to find a map on S^3 with the property that its restriction to any component of L is an embedding into \mathbb{R}^2 . We suppose that the modifications through Step 1, \dots , Step 6 happen so that all the images of the maps f_1, \dots, f_6 lie completely inside the disk determined by the (unchanged) circle $f_i(C)$, $i = 1, \dots, 6$. This can be reached easily by choosing $f_0(C)$ to bound an area “large enough” in \mathbb{R}^2 and supposing that the diameter of \bar{L} is small.

Step 1. Our first goal is to deform f_0 so that the resulting map f_1 has fold singularities along L . Apply Lemma 2.3 to the map $f_0 : S^3 \rightarrow \mathbb{R}^2$ and the embedded 1-dimensional manifold $L \subset S^3$, and denote the resulting stable map by f_1 . It is a fold map, its indefinite fold singular set is L and its definite fold singular set is $C \cup L'$, where $L' = s(L)$ is isotopic to L ; for an example see Figure 5.

Since L' is isotopic to L , the integral surgery along L giving M can be equally performed along L' . Recall that doing surgery along L' simply means that we cut out a tubular neighborhood of the definite fold curve L' (which is diffeomorphic to $L' \times D^2$), and glue it back by a diffeomorphism of its boundary $L' \times S^1$. If the image $f_1(L')$ was an embedding of circles, then it would be easy to construct the claimed map F on the 3-manifold given by the integral surgery. Since this is not the case in general, we need to further deform the map f_1 .

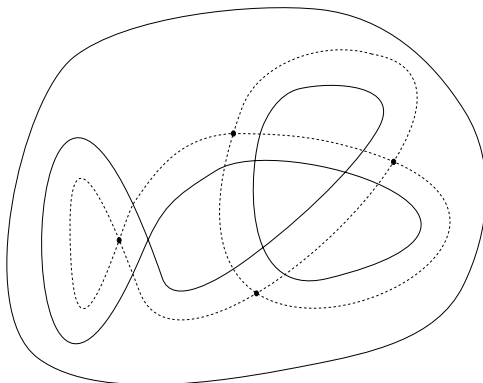


Figure 5. The image of the singular set of the map $f_1 : S^3 \rightarrow \mathbb{R}^2$, where L is the trefoil knot. The outer circle is $f_1(C)$, the inner solid curve is $f_1(L')$ and the dashed curve is $f_1(L)$. The double points of $f_1(L)$ correspond to singularities of type (D).

Let B denote the interior of the bands (one for each component of L) bounded by $q_{f_1}(L)$ and $q_{f_1}(L')$ in the Stein factorization W_{f_1} . Then B is immersed into \mathbb{R}^2 by \bar{f}_1 . The Stein factorizations of the maps f_2, \dots, f_6 in the next steps will be built on B . Let B' denote the surface $W_{f_1} - \text{cl } B$.

Step 2. Now, our goal is to deform f_1 so that the Stein factorization of the resulting map f_2 has small “flappers” near $q_{f_2}(L')$ at the points where $\bar{f}_2(q_{f_2}(L'))$ is tangent to the general position vector v . These “flappers” will help us to move the image of L so that it will become an embedding into \mathbb{R}^2 .

First, we use Lemma 2.3 together with Remark 2.4 as follows. Let T be the set of points in $q_{f_1}(L')$ such that for each $p \in T$ the direction v is tangent to $f_1(L')$ at $\bar{f}_1(p)$. For each $p \in T$ take a small embedded arc α_p in a small neighborhood of p in B such that $\bar{f}_1|_{\alpha_p}$ is an embedding parallel to $f_1(L)$. For each arc α_p there exists an embedded arc $\tilde{\alpha}_p$ in S^3 such that $q_{f_1}|_{\tilde{\alpha}_p}$ is an embedding onto α_p . See, for example, the upper picture of Figure 6, where the small dashed arcs having cusp endpoints represent the arcs $f_1(\tilde{\alpha}_p) = \bar{f}_1(\alpha_p)$ for all $p \in T$.

Apply Lemma 2.3 and Remark 2.4 to the map $f_1 : S^3 \rightarrow \mathbb{R}^2$ and the arcs $\{\tilde{\alpha}_p \subset S^3 : p \in T\}$ to obtain a map f'_1 . The section s in Lemma 2.3 is chosen so that if we project the f'_1 -images of the arising new definite fold curves in \mathbb{R}^2 to $\mathbb{R}v$, then for each curve there is only one critical point, which is a maximum. An example for the resulting map f'_1 can be seen in the upper picture of Figure 6. Note that the deformation yielded small “flappers” in $W_{f'_1}$ attached to B along the arcs $\{\alpha_p : p \in T\}$. Next, for each $p \in T$ take small arcs β_p in $W_{f'_1}$ which intersect generically the previous arcs $\{\alpha_p : p \in T\}$, lie in B and on the “flappers” and are mapped into \mathbb{R}^2 almost parallel to v . See the new small dashed arcs in the lower picture of Figure 6. Once again, there are small arcs $\{\tilde{\beta}_p : p \in T\}$ embedded in S^3 mapped by f'_1 onto $\{\beta_p : p \in T\}$, respectively.

The application of Lemma 2.3 and Remark 2.4 for these arcs provides us a map, which we denote by f_2 . This map will have one additional flapper for every flapper of f'_1 . We choose the section s in Lemma 2.3 so that the f_2 -images of the arising new definite fold curves lie inward² from the arcs $\{\bar{f}'_1(\beta_p) : p \in T\}$, respectively, in the \bar{f}_2 -image of B and the previous flappers. For an enlightening example, see the lower picture of Figure 6. Note that after this step $|T|$ new singular points of type (D) appeared. Also note that for each $p \in T$ we have four cusp singular points in S^3 , three of which are mapped by q_{f_2} into B . We denote the set of these three cusps by C_p . For each $p \in T$ the f_2 -images of two of these three cusps in C_p point to the direction $-v$. We denote the set of these two cusps by D_p . Note that the definite fold curves in the images of the two cusps in D_p are on opposite sides.

²At a point of $\{\bar{f}'_1(p) : p \in T\}$ let us call the direction which is perpendicular to $f'_1(L')$ and points toward the direction where locally $f'_1(L')$ lies “inward”.

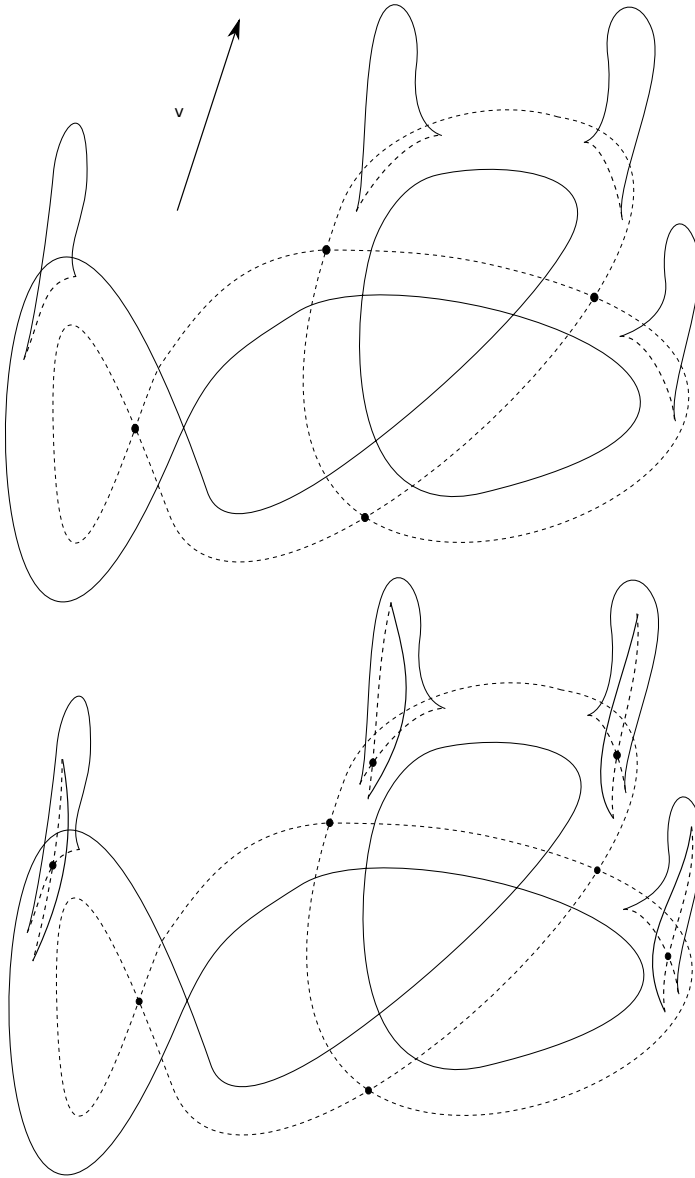


Figure 6. We obtain the upper picture by applying Lemma 2.3 and Remark 2.4 to the small arcs $\{\tilde{\alpha}_p : p \in T\}$ in S^3 which are mapped by f_1 to the dashed arcs near the points of the diagram \bar{L} where it is tangent to v . We obtain the lower picture by applying Lemma 2.3 and Remark 2.4 to the new arcs added to the upper picture. The solid arcs correspond to singularities of type (A) and the black double points of the dashed arcs correspond to singularities of type (D).

Step 3. Now our goal is to obtain definite fold arcs connecting points of S^3 where f_2 had cusps. Moreover these definite fold arcs will be mapped into \mathbb{R}^2 parallel to the diagram \bar{L} . (These curves will be translated in the next step so that later they will lead to an embedding of L into \mathbb{R}^2 .)

In order to reach this goal, we deform the map $f_2 : S^3 \rightarrow \mathbb{R}^2$ further by eliminating half of the cusps as follows. We proceed for each component of L separately and in the same way, thus in the following we can suppose that L is connected. Take a cusp $q_0 \in S^3$ which is in $C_x - D_x$ for an $x \in T$ such that the entire $f_2(L')$ lies to the right hand side of its tangent at $\bar{f}_2(x)$. By going along the band B in W_{f_2} in the direction to which the f_2 -image of this cusp q_0 points, we reach another cusp q_1 in C_p for some $p \in T$ at the next v -tangency of $f_2(L')$. If this cusp does not belong to D_p , then it is possible to apply Lemma 2.2 and eliminate these two cusps, since they are in the position of Figure 3. Then we continue by taking the cusp in D_p whose Stein factorization is folded inward. If the cusp q_1 does belong to D_p , then we choose that cusp from D_p which can be used to eliminate q_0 (it is easy to see that this is exactly the cusp in D_p whose Stein factorization is folded inward), we eliminate them, then we continue by taking the cusp belonging to $C_p - D_p$. This procedure goes all along the band B , meets all $p \in T$ and eliminates half of the cusps. After finishing this process, we obtain a stable map, which we denote by f_3 ; see Figure 7 for an example.

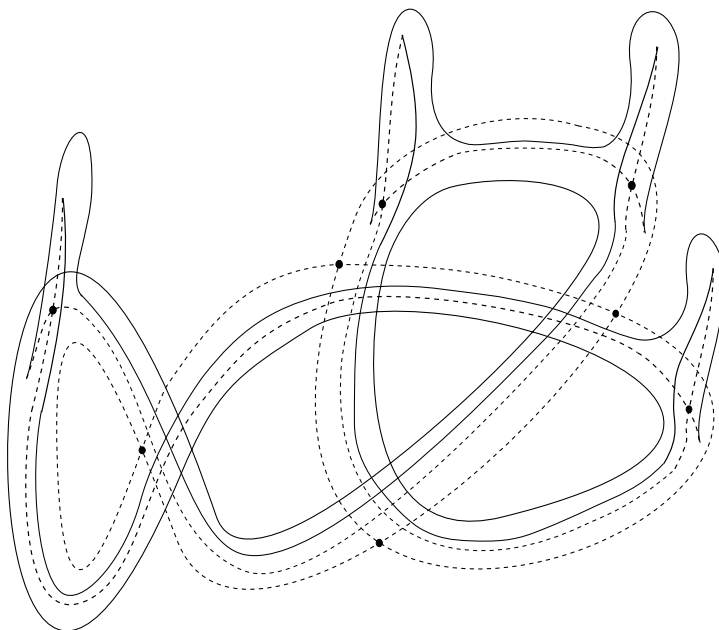


Figure 7. Eliminating half the cusps in the lower part of Figure 6. The black double points correspond to singularities of type (D).

The cusp elimination results new definite fold curves whose f_3 -image is an immersion, and which have double points near the crossings of the diagram \bar{L} . In the next step we will deform f_3 so that the double points of these new curves will be localized near the images of the remaining cusps.

Step 4. Now our goal is to deform f_3 to a map f_4 such that the definite fold arcs obtained in the previous step will be mapped into \mathbb{R}^2 far from the diagram \bar{L} . (Informally, we will “lift” some of the arcs in the direction of v .) Moreover, the immersion of these definite fold arcs into \mathbb{R}^2 will have double points only near some cusps of f_4 . This brings us closer to the original goal to have a map which embeds a link isotopic to L into the plane.

The cusp eliminations above affect only small tubular neighborhoods of curves connecting cusps in S^3 . Denote by $\delta \subset S^3$ the new definite fold arcs which appear in these tubular neighborhoods after the eliminations. Note that by the algorithm above, the arcs δ are mapped into \mathbb{R}^2 so that by an elementary deformation they can be moved “upward” in the direction of v , see Figure 7.

So we further deform $f_3 : S^3 \rightarrow \mathbb{R}^2$ to get a stable map denoted by f_4 as indicated in Figure 8: as it is shown by the picture, the arcs are “lifted”. In fact, we deform \bar{f}_3 : we move the top of the “flappers” corresponding to the α -curves of Step 2 and the \bar{f}_3 -image of the curves $q_{f_3}(\delta)$ in the direction of v and far from $f_3(L)$. We proceed for each component of L separately and in the same way, thus in the following we can suppose that L is connected. First we choose a point $x \in T$ such that the entire $f_3(L')$ lies to the right hand side from its tangent at $\bar{f}_3(x)$. Then, by walking along the band $B \subset W_{f_3}$ starting from x , we deform the flappers and the curves $\bar{f}_3(q_{f_3}(\delta))$ to be mapped into the plane as a “zigzag” far away from the diagram \bar{L} . More precisely, consider the coordinate system in \mathbb{R}^2 with origin x and coordinate axes $\mathbb{R}v^\perp$ and $\mathbb{R}v$, respectively, where v^\perp denotes the vector obtained by rotating v clockwise by 90 degrees. By extending the \bar{f}_3 -image of the flappers in the direction of v deform the \bar{f}_3 -image of the curves $q_{f_3}(\delta)$ so that by going along B between the points $p_i, p_{i+1} \in T$, where $1 \leq i \leq |T| - 1$ and $p_1 = x$, the corresponding component of the curve $f_3(\delta)$ is mapped into a small tubular neighborhood of a line with slope $(-1)^{i+1}$ for $i = 1, \dots, |T| - 1$. Finally, arrange the last component of $f_3(\delta)$ starting with slope -1 and ending at the first (extended) flapper belonging to x , see Figure 8.

As a result the double points of the immersion of the deformed curves $f_4(\delta)$ are in a small neighborhood of the cusps mapped close to the tops of the flappers.

Step 5. In this step, we modify the stable map f_4 so that the cusps of the resulting map f_5 will be easy to eliminate in the next step. Let $l \subset \mathbb{R}^2$ be a line perpendicular to v located near $\bar{f}_4(B)$, separating it from the other parts moved to the direction of v in Step 4, as indicated in Figure 8.

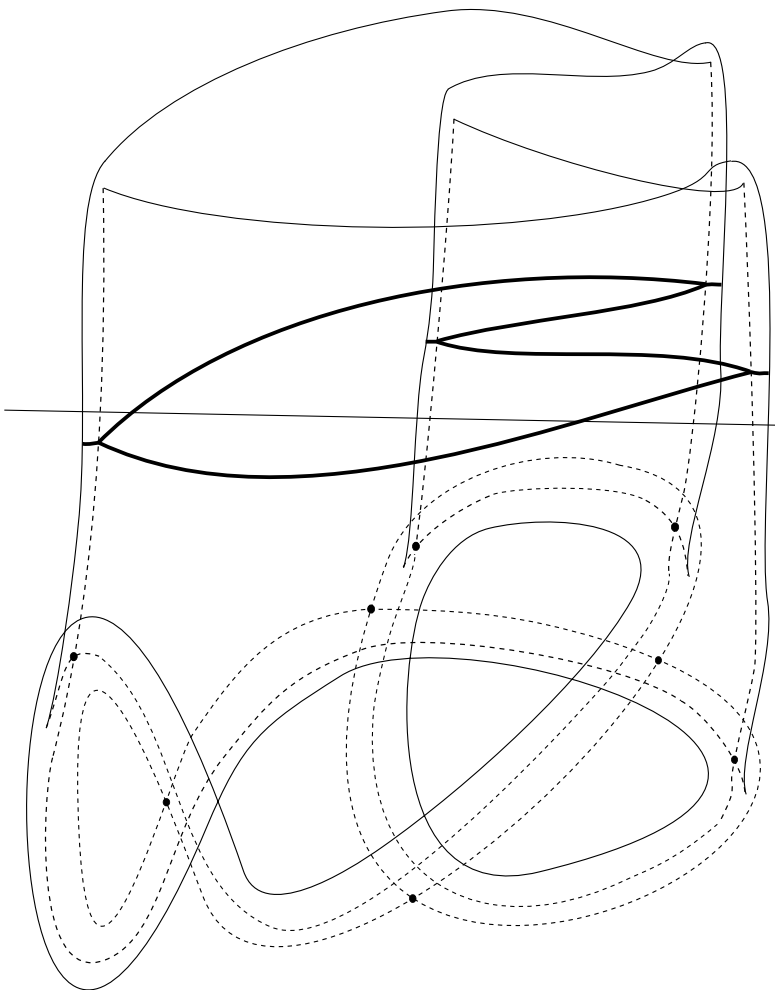


Figure 8. The Stein factorization of f_4 , i.e., the deformation of f_3 of Figure 7. (The straight line represents the line l used to cut W_{f_4} in Step 5.) The upper part of W_{f_4} from the bold 1-complex is denoted by A . (As usual, the circle $f_4(C)$ is omitted.)

Now, we cut the 2-complex $W_{f_4} - B'$ (recall that B' denotes $W_{f_1} - \text{cl } B$; see Step 1) along the \tilde{f}_4 -preimage of the line l , thus we obtain the decomposition

$$W_{f_4} = A \cup_{\tilde{f}_4^{-1}(l) \cap (W_{f_4} - B')} A',$$

where A' denotes the 2-dimensional CW complex containing $q_{f_4}(L)$ and A denotes the closure of $W_{f_4} - A'$. Then $q_{f_4}^{-1}(A)$ is a 3-manifold with boundary. Let us denote the 1-complex $q_{f_4}(\partial q_{f_4}^{-1}(A))$ by ∂A . In order to visualize ∂A in Figure 8, we suppose that the cutting of W_{f_4} along $\tilde{f}_4^{-1}(l) \cap (W_{f_4} - B')$ is a little bit perturbed

and thus the bold 1-complex in Figure 8 represents ∂A . Before proceeding further, we need a better understanding of the q_{f_4} -preimages of the sets appearing in the above decomposition. The preimage $q_{f_4}^{-1}(\partial A)$ is clearly diffeomorphic to $J \times S^1$ for a link $J \subset S^3$. The following statements show much more about $q_{f_4}^{-1}(\partial A)$. It is easy to see that the numbers of components of J and L are equal. However, we have a stronger result:

Lemma 3.1. *A longitudinal curve in $q_{f_4}^{-1}(\partial A)$ is isotopic to L .*

Proof. The 1-complex ∂A decomposes as a union of 1-cells: some of them (which we depict as “small 1-cells” in Figure 8) are attached at one of their endpoints to the union of the other 1-cells, we denote these small cells by σ_i for $i = 1, \dots, |T|$. Others are attached by both of their endpoints. Let σ denote the 1-complex $\partial A - \bigcup_{i=1}^{|T|} \sigma_i$. Then the PL embedding $\sigma \subset W_{f_4}$ is isotopic to the subcomplex ι of W_{f_4} formed by the arcs of type (B) in the open bands B connecting the singular points of type (D) in B . Furthermore, the subcomplex ι is isotopic to $q_{f_4}(L')$. Take a small closed regular neighborhood N of $q_{f_4}(L')$. Then $q_{f_4}^{-1}(N)$ is naturally a D^2 -bundle over L' . The boundary of N in W_{f_4} is a 1-manifold isotopic to $q_{f_4}(L')$, and we will denote it by λ . Clearly $q_{f_4}^{-1}(\lambda)$ is diffeomorphic to $L' \times S^1$. Note that any section of $q_{f_4}^{-1}(\lambda)$ is isotopic to L' .

The isotopy between λ and ι and the isotopy between ι and σ can be chosen easily so that they give a PL embedding $\varepsilon : S^1 \times [0, 1] \rightarrow W_{f_4}$ such that $S^1 \times \{0\}$ and $S^1 \times \{1\}$ correspond to λ and σ , respectively. For $j = 1, \dots, |T|$, let U_j denote small regular neighborhoods of the singular points of type (D) located near the cusp points in B in W_{f_4} , such a U_j and the restriction $\tilde{f}_4|_{U_j}$ can be seen in Figure 1(d). Then the intersection

$$\varepsilon(S^1 \times [0, 1]) \cap \left(\bigcup_{j=1}^{|T|} U_j \right)$$

consists of a union of disks, which will be denoted by

$$\bigcup_{j=1}^{|T|} D_j.$$

First, observe that for each $j = 1, \dots, |T|$ there exists a disk \tilde{D}_j embedded into $q_{f_4}^{-1}(U_j)$ in S^3 whose boundary $\partial \tilde{D}_j$ is mapped by q_{f_4} homeomorphically onto the boundary ∂D_j ; i.e., $\partial \tilde{D}_j$ is a lifting of ∂D_j . To see this, consider the 3-manifold $q_{f_4}^{-1}(U_j)$ for each $j = 1, \dots, |T|$. By [Levine 1985] the manifold $q_{f_4}^{-1}(U_j)$ is diffeomorphic to $R \times [0, 1]$, where R is a disk with three holes and it is mapped by f_4 into \mathbb{R}^2 as we can see in Figure 9(a).

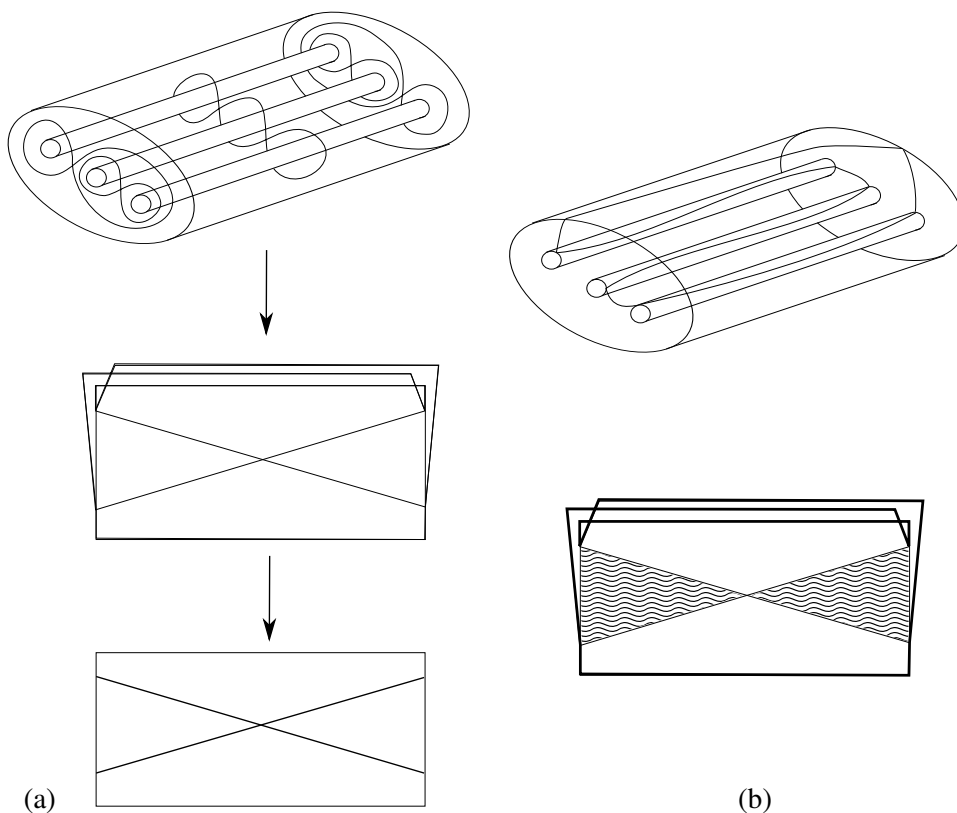


Figure 9. In (a) we can see the manifold $R \times [0, 1]$ and how it is mapped onto the regular neighborhood U_j and into \mathbb{R}^2 ; cf. Figure 1(d). $R \times \{0\}$ is mapped onto the left side of the rectangle $\bar{f}_4(U_j)$ as a proper Morse function with two indefinite critical points. The two “figure eights” in $R \times \{0\}$ are the two singular fibers. $R \times \{1\}$ is mapped similarly onto the right side of $\bar{f}_4(U_j)$. The middle fiber in $R \times [0, 1]$ is mapped to the singular point of type (D). For a detailed analysis see [Levine 1985]. In (b) we can see the boundary $\partial \tilde{D}_j$ in $R \times [0, 1]$ and its image in U_j represented by a bold 1-complex.

Each disk D_j can be located in U_j essentially in four ways, for example the lower picture of Figure 9(b) shows the disk D_j for the leftmost nonsimple singularity crossing of type (D) in Figure 8. We get D_j on the picture by cutting out the two shaded areas from the 2-complex U_j . It is easy to see in the upper picture of Figure 9(b) how to put the disk \tilde{D}_j into $R \times [0, 1]$. The other three possibilities for the location of a disk D_j in U_j and the disk \tilde{D}_j in $q_{f_4}^{-1}(U_j)$ can be described in a similar way.

Now observe that $\varepsilon(S^1 \times [0, 1]) - \bigcup_{j=1}^{|T|} D_j$ can be lifted to S^3 extending $\bigcup_{j=1}^{|T|} \tilde{D}_j$ because of the following. First, the regular neighborhoods of the singular points of type (C) in B (see Figure 1(c)) intersect $\varepsilon(S^1 \times [0, 1])$ in disks which can be lifted to S^3 . Then the intersection of the small regular neighborhoods of the singular curves of type (B) and $\varepsilon(S^1 \times [0, 1])$ can be lifted as well since there is no constraint for the lift at the regular points of f_4 . Finally observe that the rest of $\varepsilon(S^1 \times [0, 1])$ intersects W_{f_4} only in areas of non-singular points which are attached to the boundary of $\varepsilon(S^1 \times [0, 1])$, so the previous lifts extend over the entire $\varepsilon(S^1 \times [0, 1])$.

Hence we obtain an embedding $\tilde{\varepsilon} : S^1 \times [0, 1] \rightarrow S^3$ with $S^1 \times \{0\}$ and $S^1 \times \{1\}$ corresponding to lifts of λ and σ , respectively. Thus we obtain an isotopy between a longitude of $q_{f_4}^{-1}(\partial A)$ and a lift of λ . The fact that any lift of λ is isotopic to L' finishes the proof. \square

Lemma 3.2. *The preimage $q_{f_4}^{-1}(A)$ is isotopic to a regular neighborhood of L .*

Proof. It is enough to show that $q_{f_4}^{-1}(A)$ is diffeomorphic to $L \times D^2$ extending naturally the $L \times S^1$ structure on its boundary since by Lemma 3.1 the union of tori $\partial q_{f_4}^{-1}(A)$ contains a longitude isotopic to L . Moreover it is enough to show that the q_{f_4} -preimage of the part of A homeomorphic to the CW complex in Figure 10 is diffeomorphic to $[0, 1] \times D^2$, where the q_{f_4} -preimage of the two vertical edges on the right-hand side of the 2-complex of Figure 10 corresponds to $\{0, 1\} \times D^2$. Clearly the q_{f_4} -preimage of the two vertical edges on the right-hand side is diffeomorphic to $\{0, 1\} \times D^2$ since $q_{f_4}^{-1}(x)$ is a circle for any x lying in the two vertical edges except if x is one of the two top ends. If x is one of the two top ends, then $q_{f_4}^{-1}(x)$ is one point since it is a definite fold singularity. The q_{f_4} -preimage of the backward sheet in Figure 10 is diffeomorphic to $[0, 1] \times D^2$ minus $I \times D^2$ for an interval I . The q_{f_4} -preimage of the forward sheet is diffeomorphic to $I \times D^2$. \square

Corollary 3.3. *Any longitudinal curve in $q_{f_4}^{-1}(\partial A)$ is isotopic to L .*

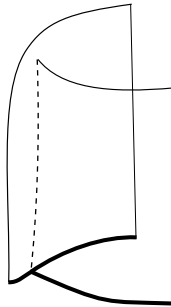


Figure 10

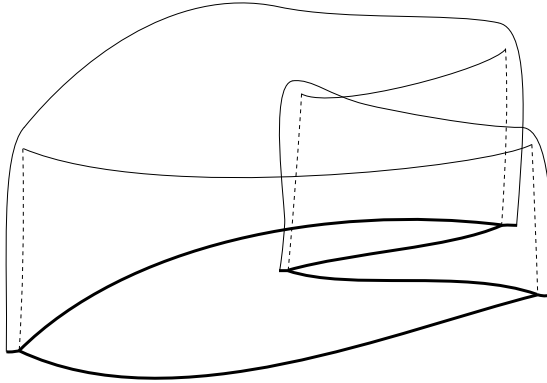


Figure 11. The Stein factorization of

$$f_5|_{q_{f_5}^{-1}(W_{f_5}-A')} : L \times D^2 \rightarrow \mathbb{R}^2.$$

There are two \mathcal{P} -pairs of cusps.

In order to obtain the map f_5 , we modify the map

$$f_4|_{q_{f_4}^{-1}(A)} : L \times D^2 \rightarrow \mathbb{R}^2$$

outside a neighborhood of $q_{f_4}^{-1}(\partial A)$, as shown by Figure 11: our goal is to have the arrangement that if for a cusp singularity $q_1 \in S^3$ the point $q_{f_5}(q_1)$ is connected in $W_{f_5} - A'$ to ∂A by a 1-cell γ mapped into \mathbb{R}^2 parallel to v and γ corresponding to an indefinite fold curve, then a definite fold curve should connect q_1 to another cusp q_2 with the same property for $q_{f_5}(q_2)$. Thus we obtain a map f_5 such that $q_{f_5}^{-1}(W_{f_5} - A')$ is isotopic to a regular neighborhood of L by the same argument as in Lemma 3.2. Also $q_{f_5}^{-1}(W_{f_5} - A')$ coincides with $q_{f_4}^{-1}(A)$ and f_5 coincides with f_4 in a neighborhood of $q_{f_5}^{-1}(A')$.

We arrange the cusps of f_5 in $q_{f_4}^{-1}(A)$ to form pairs as follows. In W_{f_5} sheets are attached to B along arcs of type (B) (possibly containing points of type (C) at some endpoints). Walking along the bands B and restricting ourselves to the intersection of the sheets and $W_{f_5} - A'$, we have that every sheet contains a pair of cusps and every second sheet contains a singular arc of type (A) connecting its pair of cusps; for example, see Figure 11.

A natural pairing is that two cusps form a pair if they are in the same sheet and they are connected by a singular arc of type (A). We refer to this pairing as \mathcal{Q} -pairing. We also define another pairing \mathcal{P} : two cusps form a \mathcal{P} -pair if they are in the same sheet and they are *not* connected by any singular arc of type (A).

Step 6. In this step, we eliminate the cusps of f_5 contained in $q_{f_5}^{-1}(W_{f_5} - A')$. These cusps are mapped by f_5 in the direction of v far from \bar{L} and arranged into

\mathcal{P} -pairs in the previous step. The restriction of the resulting map $f_6 : S^3 \rightarrow \mathbb{R}^2$ to a link isotopic to L will be an embedding. (Hence after this step the construction of the claimed map F on M will be easy.)

We have exactly $|T|/2$ \mathcal{P} -pairs of cusps in $q_{f_5}^{-1}(W_{f_5} - A')$. Observe that for each component of L one \mathcal{P} -pair can be eliminated immediately: for example in Figure 11 the pair on the “highest” sheet is in the sufficient position to eliminate. In the following, we deal with the other \mathcal{P} -pairs.

More concretely, we perform the deformations and the eliminations of the pairs of cusps of f_5 in $q_{f_4}^{-1}(A)$ as shown in Figure 12 as follows.

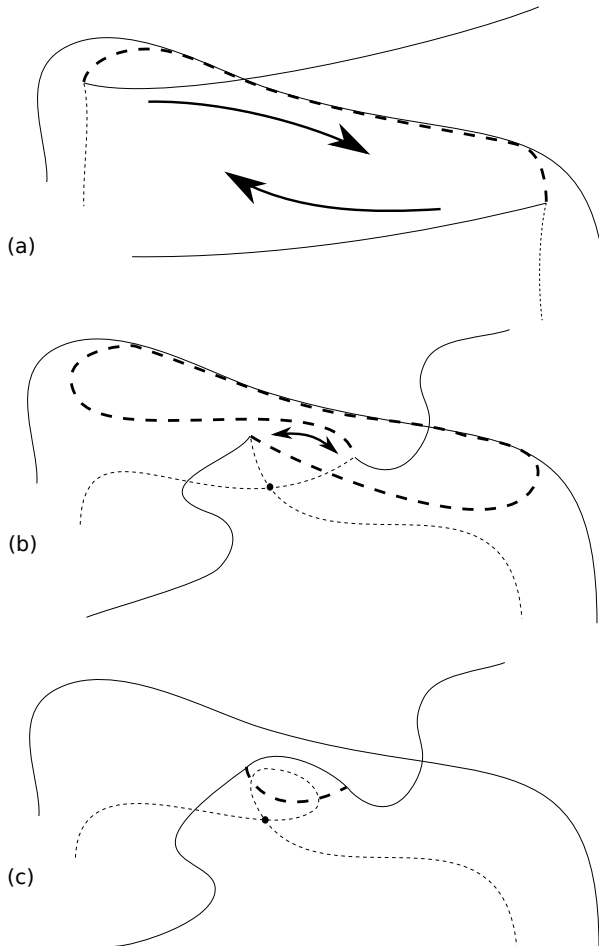


Figure 12. Moving and eliminating the cusps. We move and eliminate the \mathcal{P} -pair of cusps along the arrows. The dashed arcs represent 1-complexes used to deform σ in the proof of Lemma 3.4.

First, by using Lemma 2.1 we move each pair of cusps having the position as in Figure 12(a) to the position as in Figure 12(b) thus creating a singularity of type (D). Then by using Lemma 2.2 we eliminate each pair of cusps, see Figures 12(b) and 12(c).

The resulting map will be denoted by f_6 (see Figure 13). Notice that f_6 and f_5 coincide in a neighborhood of $q_{f_5}^{-1}(A')$. The deformations above yield definite fold curves $K \subset S^3$, whose image under f_6 is an embedding into \mathbb{R}^2 as indicated in Figure 13 by the bold curve.

Lemma 3.4. *The link K is isotopic to L .*

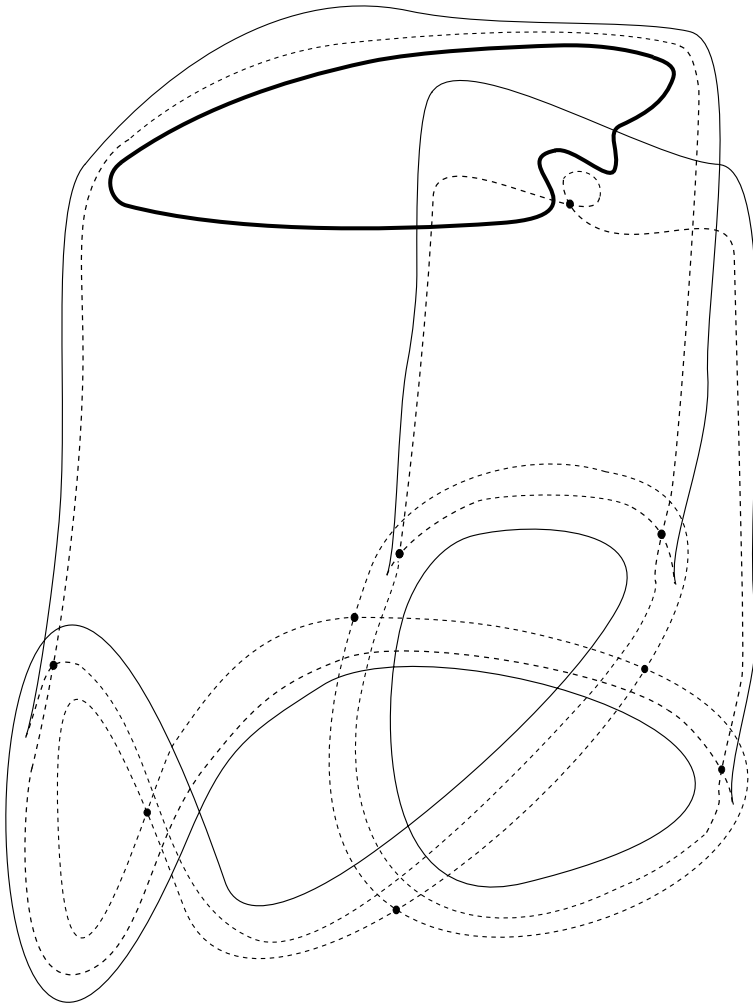


Figure 13. The Stein factorization of the stable map $f_6 : S^3 \rightarrow \mathbb{R}^2$. (The circle $f_6(C)$ is omitted.)

Proof. By Lemma 3.1 the link L is isotopic to a longitude of the union of tori $q_{f_4}^{-1}(\partial A)$. In Step 6 we modify f_5 only inside $q_{f_4}^{-1}(A)$. The subcomplex σ of ∂A used in the proof of Lemma 3.1 is PL-isotopic to a 1-dimensional PL submanifold σ' of $W_{f_5} - A'$ such that σ' goes through the singular curves of type (A) appearing in the 2-pairing at the end of Step 5 and goes through the top of $W_{f_5} - A'$, i.e., the top of the 2-complex in Figure 11. To be more precise, in Figure 12(a) the part of σ' connecting the two cusp endpoints of the singular arcs of type (A) is represented by a bold dashed arc and denoted by σ'' . During the moving of the pair of cusps as depicted by the arrows in Figure 12(a), σ'' is deformed to the curve σ''' represented by a bold dashed arc in Figure 12(b). This deformation gives an isotopy between some liftings to S^3 of σ'' and σ''' . Since a part of σ'' is collinear to a singular arc of type (A) as we can see in Figure 12(a), any lifting to S^3 of σ'' is isotopic to any other lifting. Hence further deforming σ''' to σ'''' represented by the bold dashed curve in Figure 12(c) yields an isotopy between some liftings of σ'' and σ'''' . Finally, changing again the lifting to S^3 of σ'''' if necessary, we eliminate the pair of cusps as indicated in Figure 12(b) and deform σ'''' to be identical to the type (A) singular arc appearing at the elimination in Figure 12(c). All this process gives an isotopy in S^3 between K and a lifting of σ , hence an isotopy between K and L . \square

Step 7. As a final step, we perform the given surgeries along K with the appropriate coefficients. Since $f_6|_K$ is an embedding into \mathbb{R}^2 on each component of K , and K consists of definite fold singular curves such that the local image of a small neighborhood of the definite fold curve is situated “outside” of the image of the definite fold curve, a map of M is particularly easy to construct: a small tubular neighborhood N_K of K , which is diffeomorphic to $K \times D^2$, is glued back to $S^3 - \text{int } N_K$ such that $\{pt.\} \times \partial D^2$ maps to a longitude in $\partial(M - \text{int } N_K)$, hence N_K can be mapped into \mathbb{R}^2 as the projection $\pi : K \times D^2 \rightarrow D^2$. This π extends over $M - \text{int } N_K$ and the resulting map $M \rightarrow \mathbb{R}^2$ is stable. Let us denote it by F .

It is easy to see that F has the claimed properties:

The Stein factorization W_F is homotopy equivalent to the bouquet $\bigvee_{i=1}^{n(L)} S^2$. The Stein factorization W_{f_4} is clearly contractible. The CW-complexes W_{f_5} and W_{f_6} are still contractible since the corresponding steps do not change the homotopy type. At the final surgery we attach a 2-disk to W_{f_6} for each component of L .

The number of cusps of F is equal to $t_v(\bar{L})$. Each point in $f_1(L')$ at which $f_1(L')$ is tangent to the chosen general position vector v (these are exactly the points of the set $\tilde{f}_1(T)$) corresponds to a cusp of F by the construction and there are no other cusps. $|T| = t_v(\bar{L})$ hence we get the statement.

All the nonsimple singularities of F are of type (D). This follows from the fact that singularities of type (E) never appear during the construction.

The number of the nonsimple singularities of F is equal to $\text{cr}(\bar{L}) + \frac{3}{2}t_v(\bar{L}) - n(L)$. Each crossing of the diagram \bar{L} gives a singularity of type (D). Also each point in T gives a singularity of type (D) by the construction. Finally, the movement illustrated in Figure 12(b) gives one singular point of type (D) for each pair of points in T except one pair for each component of L .

The number of nonsimple singularities which are not connected by any singular arc of type (B) to any cusp is equal to $\text{cr}(\bar{L}) + \frac{1}{2}t_v(\bar{L}) - n(L)$.

In the previous argument, if we do not count the singularities of type (D) corresponding to the v -tangencies of $f_1(L')$, then we get the statement.

The number of simple singularity crossings of F in \mathbb{R}^2 is no more than

$$8\text{cr}(\bar{L}) + 6\ell(\bar{L}, v)t_v(\bar{L}) + t_v(\bar{L})^2.$$

We can suppose that the number of simple singularity crossings of $f_4|_{q_{f_4}^{-1}(A')}$ is at most $8\text{cr}(\bar{L}) + 2t_v(\bar{L}) + 6\ell(\bar{L}, v)t_v(\bar{L})$. The maps f_4, f_5, f_6 and F coincide in a neighborhood of $q_{f_4}^{-1}(A')$ and also their images coincide in the half plane bounded by the line l and lying in the direction $-v$ (for the notations, see Step 5). The simple singularity crossings of F in $F(q_{f_4}^{-1}(A))$ come from the intersections of the \bar{F} -images of the “sheets” attached to the bands $B \subset W_F$ (for the notation, see Step 2). For example, in Figure 13, two such sheets intersect on the right-hand side in four simple singularity crossings. Hence we obtain an upper bound for the number of simple singularity crossings of F in $F(q_{f_4}^{-1}(A))$ if we suppose that all the sheets intersect each other in eight crossings. This gives the upper bound

$$8\left(\frac{t_v(\bar{L})}{2} - 1 + \frac{t_v(\bar{L})}{2} - 2 + \dots + 1\right) = 4\frac{t_v(\bar{L})}{2}\left(\frac{t_v(\bar{L})}{2} - 1\right) = t_v(\bar{L})^2 - 2t_v(\bar{L}).$$

Thus we obtain the upper bound

$$8\text{cr}(\bar{L}) + 2t_v(\bar{L}) + 6\ell(\bar{L}, v)t_v(\bar{L}) + t_v(\bar{L})^2 - 2t_v(\bar{L}) = 8\text{cr}(\bar{L}) + 6\ell(\bar{L}, v)t_v(\bar{L}) + t_v(\bar{L})^2$$

for all the simple singularity crossings of F .

The number of connected components of the singular set of F is no more than $n(L) + \frac{3}{2}t_v(\bar{L}) + 1$. The curve C is a component and the links L and L' give singular set components as well. Also the cusp elimination in Step 3 gives additional $t_v(\bar{L})$ components. Steps 4 and 5 clearly do not increase more the number of singular set components. In Step 6 the changings showed in Figure 12 increase the number of components by at most $\frac{1}{2}t_v(\bar{L})$. Finally Step 7 decreases it by $n(L)$.

The maximal number of the connected components of any fiber of F is no more than $t_v(\bar{L}) + 3$. The maximal number of the connected components of any fiber of f_1 is 3. This value is no more than $3 + t_v(\bar{L})$ for f_2, \dots, f_5 and also for f_6 . When

we perform the surgery in Step 7, $3 + t_v(\bar{L})$ is still an upper bound hence we get the statement.

The indefinite fold singular set of F . Finally the statement of (8) about the indefinite fold singular set of F is obvious from the construction. This finishes the proof of Theorem 1.2. \square

Remark 3.5. Suppose we have two links in S^3 . If the projections of the two links coincide, then the resulting stable maps on the two 3-manifolds in the construction described above will have the same Stein factorizations. Therefore only the Stein factorization itself is a very weak invariant of the 3-manifold.³

Proof of Theorem 1.4. Let M be a closed orientable 3-manifold obtained by an integral surgery along a link in S^3 . Theorem 1.2 gives a stable map F of M into \mathbb{R}^2 without singularities of type (E). We can eliminate the cusps of F without introducing any singularities of type (E). Indeed, the map constructed by Theorem 1.2 has an even number of cusps, whose q_F -image is situated in $B \subset W_F$. Moreover since the locations of the F -images of the cusps are at the v -tangencies of \bar{L} , each cusp c has a pair c' which can be moved close to c (thus possibly creating new singular points of type (D)) and can be used to eliminate these pairs in the sense of Lemmas 2.1 and 2.2. \square

Remark 3.6. By results from [Eliashberg and Mishachev 1997], every closed orientable 3-manifold has a wrinkled map into \mathbb{R}^2 since any orientable 3-manifold is parallelizable. This argument leads to another proof of Theorem 1.4. However, the h -principle used in the proof of the results cited does not provide any construction for the wrinkled map.

Next we give the proof of the estimate given in (1-1) in Section 1.

Lemma 3.7. $\ell(\bar{L}, v) \leq t_v(\bar{L}) - 1$.

Proof. For any v -tangency p we have $\ell(\bar{L}, v, p) \leq t_v(\bar{L}) - 1$ since by going along the components of L in the diagram \bar{L} , in order to pass through the intersections of the half line emanating from p in the direction of v , for each intersection one needs to pass through a v -tangency as well. \square

4. Estimates for TB^-

Recall that the Thurston–Bennequin number $\text{tb}(\mathcal{L})$ of a Legendrian knot \mathcal{L} can be computed through the simple formula

$$\text{tb}(\mathcal{L}) = w(\bar{\mathcal{L}}) - \frac{1}{2}\#\text{cusps}(\bar{\mathcal{L}}).$$

³The paper [Motta et al. 1995] is closely related to this remark.

Proof of Theorem 1.9. By Theorem 1.2(5) and Lemma 3.7 we have

$$s(F) \leq 8\text{cr}(\bar{L}) + 7t_v(\bar{L})^2 - 6t_v(\bar{L})$$

for the constructed stable map F . (Here, again, \bar{L} denotes the generic projection of the knot L we get from the front projection of the Legendrianization \mathcal{L} of L by rounding the cusps.) Since $d(F) = s(F) + ns(F)$, by Theorem 1.2 (3), (5) and Lemma 3.7 we have

$$d(F) \leq 9\text{cr}(\bar{L}) + 7t_v(\bar{L})^2 - \frac{9}{2}t_v(\bar{L}) - n(L).$$

If $\bar{\mathcal{L}}$ has only negative crossings, then the Thurston–Bennequin number $\text{tb}(\mathcal{L})$ is equal to $-\text{cr}(\bar{L}) - \frac{1}{2}t_v(\bar{L})$, where v is the vector in which the front projection has no tangency.

Hence

$$28\text{tb}(\mathcal{L})^2 = 28\text{cr}(\bar{L})^2 + 28\text{cr}(\bar{L})t_v(\bar{L}) + 7t_v(\bar{L})^2$$

and

$$28\text{cr}(\bar{L})^2 + 28\text{cr}(\bar{L})t_v(\bar{L}) + 7t_v(\bar{L})^2 \geq 9\text{cr}(\bar{L}) + 7t_v(\bar{L})^2 - \frac{9}{2}t_v(\bar{L}) - n(L).$$

Thus $|\text{tb}(\mathcal{L})| \geq \sqrt{d(F)}/\sqrt{28}$, implying (by the fact that $\text{tb}(\mathcal{L})$ is negative for a knot admitting a projection with only negative crossings)

$$(4-1) \quad \text{tb}(\mathcal{L}) \leq -\frac{\sqrt{d(F)}}{\sqrt{28}}.$$

Also by Theorem 1.2 (4), we have

$$|\text{tb}(\mathcal{L})| = \text{cr}(\bar{L}) + \frac{1}{2}t_v(\bar{L}) \geq \text{nsnc}(F) + 1,$$

which gives

$$(4-2) \quad \text{tb}(\mathcal{L}) \leq -\text{nsnc}(F) - 1.$$

Finally note that $d(F) \geq s(F)$ for any stable map F , and by taking the minimum for all the stable maps in (4-1) and (4-2), we get the statement. \square

Acknowledgements

The authors were supported by OTKA NK81203 and by the *Lendület program* of the Hungarian Academy of Sciences. The first author was partially supported by Magyary Zoltán Postdoctoral Fellowship. The authors thank the anonymous referee for the comments which improved the paper.

References

- [Baykur 2008] R. İ. Baykur, “Existence of broken Lefschetz fibrations”, *Int. Math. Res. Not.* **2008** (2008), Art. ID rnn 101. MR 2010b:57026
- [Baykur 2009] R. İ. Baykur, “Topology of broken Lefschetz fibrations and near-symplectic four-manifolds”, *Pacific J. Math.* **240**:2 (2009), 201–230. MR 2010c:57035 Zbl 1162.57011
- [Burllet and de Rham 1974] O. Burllet and G. de Rham, “Sur certaines applications génériques d’une variété close à 3 dimensions dans le plan”, *Enseignement Math. (2)* **20** (1974), 275–292. MR 51 #1846 Zbl 0299.58005
- [Costantino and Thurston 2008] F. Costantino and D. Thurston, “3-manifolds efficiently bound 4-manifolds”, *J. Topol.* **1**:3 (2008), 703–745. MR 2009g:57034 Zbl 1166.57016
- [Eliashberg and Mishachev 1997] Y. Eliashberg and N. M. Mishachev, “Wrinkling of smooth mappings and its applications. I”, *Invent. Math.* **130**:2 (1997), 345–369. MR 99d:57021 Zbl 0896.58010
- [Gay and Kirby 2007] D. T. Gay and R. Kirby, “Constructing Lefschetz-type fibrations on four-manifolds”, *Geom. Topol.* **11** (2007), 2075–2115. MR 2009b:57048 Zbl 1135.57009
- [Geiges 2008] H. Geiges, *An introduction to contact topology*, Cambridge Studies in Adv. Math. **109**, Cambridge University Press, Cambridge, 2008. MR 2008m:57064 Zbl 1153.53002
- [Gromov 2009] M. Gromov, “Singularities, expanders and topology of maps, I: Homology versus volume in the spaces of cycles”, *Geom. Funct. Anal.* **19**:3 (2009), 743–841. MR 2012a:58062
- [Gromov 2010] M. Gromov, “Singularities, expanders and topology of maps, 2: From combinatorics to topology via algebraic isoperimetry”, *Geom. Funct. Anal.* **20**:2 (2010), 416–526. MR 2012a:58063 Zbl 05800304
- [Kushner et al. 1984] L. Kushner, H. Levine, and P. Porto, “Mapping three-manifolds into the plane, I”, *Bol. Soc. Mat. Mexicana (2)* **29**:1 (1984), 11–33. MR 86j:58011 Zbl 0586.57018
- [Levine 1965] H. I. Levine, “Elimination of cusps”, *Topology* **3**:suppl. 2 (1965), 263–296. MR 31 #756 Zbl 0146.20001
- [Levine 1985] H. Levine, *Classifying immersions into \mathbf{R}^4 over stable maps of 3-manifolds into \mathbf{R}^2* , Lecture Notes in Mathematics **1157**, Springer, Berlin, 1985. MR 88f:57056 Zbl 0567.57001
- [Lickorish 1962] W. B. R. Lickorish, “A representation of orientable combinatorial 3-manifolds”, *Ann. of Math. (2)* **76** (1962), 531–540. MR 27 #1929 Zbl 0106.37102
- [Motta et al. 1995] W. Motta, P. Porto, Jr., and O. Saeki, “Stable maps of 3-manifolds into the plane and their quotient spaces”, *Proc. London Math. Soc. (3)* **71**:1 (1995), 158–174. MR 96a:57067 Zbl 0845.57025
- [Ozbagci and Stipsicz 2004] B. Ozbagci and A. I. Stipsicz, *Surgery on contact 3-manifolds and Stein surfaces*, Bolyai Soc. Math. Studies **13**, Springer, Berlin, 2004. MR 2005k:53171 Zbl 1067.57024
- [Saeki 1996] O. Saeki, “Simple stable maps of 3-manifolds into surfaces”, *Topology* **35**:3 (1996), 671–698. MR 97m:57047 Zbl 0864.57028
- [Saeki 2006] O. Saeki, “Elimination of definite fold”, *Kyushu J. Math.* **60**:2 (2006), 363–382. MR 2007g:57050 Zbl 1113.57016
- [Wallace 1960] A. H. Wallace, “Modifications and cobounding manifolds”, *Canad. J. Math.* **12** (1960), 503–528. MR 23 #A2887 Zbl 0108.36101

Received April 30, 2011. Revised May 7, 2012.

BOLDIZSÁR KALMÁR
ALFRÉD RÉNYI INSTITUTE OF MATHEMATICS
HUNGARIAN ACADEMY OF SCIENCES
REÁLTANODA UTCA 13-15
1053 BUDAPEST
HUNGARY
bkalmar@renyi.hu

ANDRÁS I. STIPSICZ
ALFRÉD RÉNYI INSTITUTE OF MATHEMATICS
HUNGARIAN ACADEMY OF SCIENCES
REÁLTANODA UTCA 13-15
1053 BUDAPEST
HUNGARY

and

INSTITUTE FOR ADVANCED STUDY
PRINCETON, NJ 08540
UNITED STATES
stipsicz@renyi.hu

STRONG SOLUTIONS TO THE COMPRESSIBLE LIQUID CRYSTAL SYSTEM

YU-MING CHU, XIAN-GAO LIU AND XIAO LIU

We prove the existence of local strong solutions of the compressible liquid crystal system.

1. Introduction

We consider the following simplified system of Ericksen–Leslie equations:

$$(1.1) \quad \rho_t + \operatorname{div}(\rho u) = 0,$$

$$(1.2) \quad \rho u_t + \rho u \cdot \nabla u + \nabla p - \mu \Delta u + \lambda \left(\operatorname{div}(\nabla n \otimes \nabla n) - \nabla \frac{|\nabla n|^2}{2} \right) = 0,$$

$$(1.3) \quad \frac{\partial n}{\partial t} + u \cdot \nabla n - \nu(\Delta n + |\nabla n|^2 n) = 0,$$

with the following initial and boundary conditions:

$$(1.4) \quad (\rho, u, n)|_{t=0} = (\rho_0, u_0, n_0), \quad x \in \Omega,$$

$$(1.5) \quad u(x, t) = u_0(x) = 0, \quad n(x, t) = n_0(x), \quad x \in \partial\Omega,$$

where u is the velocity field, n the macroscopic average of the nematic liquid crystal orientation field, $\rho_0 \geq 0$, $|n_0| = 1$, and pressure $p = a\rho^\gamma$ with $\gamma > 1$, where γ is the adiabatic constant (in the physically relevant case of a monoatomic gas, $\gamma = \frac{5}{3}$). This system is modeled after the theory of Oseen [1933] and Frank [1958]; see the articles [Ericksen 1962; Forster et al. 1971; Leslie 1966; 1968] or the books [Ericksen and Kinderlehrer 1987; Gennes and Prost 1993; Pasechnik et al. 2009; Stephen 1970; Xie 1988].

The system (1.1)–(1.3) is much more complicated than the compressible Navier–Stokes equations, because equation (1.3), like the situation with heat flow into a sphere, makes the strongly coupling term $\operatorname{div}(\nabla n \otimes \nabla n) - \nabla \frac{|\nabla n|^2}{2}$ have a weak convergence. So far, the existence of weak solutions to the system remains open, though there are celebrated contributions by Lions [1998]; see also [Feireisl 2004;

This work was supported partly by NSFC grant 11071043, 11131005, and 11071069.

MSC2010: 76N10, 35Q35, 35Q30.

Keywords: strong solutions, compressible liquid crystals, local existence.

Feireisl et al. 2001]. Liu and Qing [2011] proved the global existence of finite energy weak solutions to the case where the free energy is replaced by the Ginzburg–Landau approximation energy,

$$\min_{n \in H^1(\Omega; \mathbb{R}^3)} \int_{\Omega} \frac{1}{2} |\nabla n|^2 + \frac{1}{4\sigma^2} (|n|^2 - 1)^2 dx.$$

In the incompressible case, F. H. Lin and C. Liu, among others [Lin 1989; Lin and Liu 1995; Lin and Liu 2001; Lin and Liu 2000; Lin and Liu 1996; Calderer and Liu 2000], systematically studied the incompressible liquid crystal dynamics system based on the Ericksen–Leslie model (that is, the Ginzburg–Landau approximation case with ρ being a constant in system (1.1) makes the velocity field divergence free) and proved the global existence of weak solutions, classical solutions, and partial regularity. Liu and Zhang [2009] also studied the existence of weak solutions to the incompressible liquid crystal system with the Ginzburg–Landau approximation and ρ nonconstant.

It is well known that there exist no global solutions to the system (1.1)–(1.3) even in the incompressible case. Surprisingly, we can prove the local existence of a strong solution to the compressible liquid crystal system with initial density $\rho_0 \geq 0$. We gained enlightenment from the corresponding results of the compressible Navier–Stokes equations. There is a huge literature on the compressible Navier–Stokes equations, under the crucial assumption that the initial density ρ_0 is bounded below away from zero. The existence results were obtained by Nash, Itaya, Tani, Matsumura, and Nishida, among others. For general nonnegative initial density, Cho, Kim, and Choe [Choe and Kim 2003; Cho et al. 2004; Cho and Kim 2006] obtained the existence of a local strong solution to a compressible Navier–Stokes equation.

We first have the energy law

$$\frac{dE}{dt} + \int_{\Omega} \mu |\nabla u|^2 + \lambda \nu |\Delta n + |\nabla n|^2 n|^2 = 0$$

with

$$E(t) = \int_{\Omega} \left(\frac{1}{2} \rho u^2 + \frac{\lambda}{2} |\nabla n|^2 + \frac{a}{\gamma - 1} \rho^\gamma \right).$$

From the definition of velocity,

$$(1.6) \quad \frac{dx(X, t)}{dt} = u(x(X, t), t),$$

$$(1.7) \quad x(X, 0) = X.$$

The continuity equation can be rewritten as

$$\frac{d\rho(x(X, t), t)}{dt} + \rho \operatorname{div} u = 0,$$

that is,

$$(1.8) \quad \rho(x, t) = \rho_0 \exp\left(-\int_0^t \operatorname{div} u\right).$$

We need the following regularity for ρ_0 , n_0 , and u_0 :

$$(1.9) \quad \rho_0 \in W^{1,6}(\Omega), \quad u_0 \in H_0^1(\Omega) \cap H^2(\Omega), \quad n_0 \in H^3(\Omega).$$

We also need some compatibility condition on the initial data: for some $g \in L^2$,

$$(1.10) \quad \mu \Delta u_0 - \lambda \operatorname{div}(\nabla n_0 \otimes \nabla n_0 - \frac{1}{2} |\nabla n_0|^2 I) - a \nabla \rho_0^\gamma = \rho_0^{\frac{1}{2}} g.$$

The following is our main result.

Theorem 1.1. *Assume Ω is a smooth bounded domain in \mathbb{R}^3 and (ρ_0, n_0, u_0) satisfies regularity condition (1.9) and compatibility condition (1.10). Then there exist a small time $T^* > 0$ and a unique strong solution (ρ, n, u) of the compressible liquid crystal system (1.1)–(1.3) in $(0, T^*) \times \Omega$, satisfying initial and boundary conditions (1.4) and (1.5), such that*

$$\begin{aligned} \rho &\in C([0, T^*]; W^{1,6}), & \rho_t &\in C([0, T^*]; L^6), \\ u &\in C([0, T^*]; H_0^1 \cap H^2) \cap L^2(0, T^*; W^{2,6}), & u_t &\in L^2(0, T^*; H_0^1), \\ n &\in C([0, T^*]; H^2) \cap L^2(0, T^*; W^{2,6}), & n_t &\in C([0, T^*]; H_0^1), \\ \sqrt{\rho} u_t &\in C([0, T^*]; L^2). \end{aligned}$$

2. Approximation solutions

We now consider the linearized equations as follows: for fixed smooth functions $v, d : \Omega \times [0, T] \rightarrow \mathbb{R}^3$ with

$$\frac{dx(X, t)}{dt} = v(x(X, t), t)$$

and $x(X, 0) = X$, and $v(x, 0) = u_0(x)$, $d(x, 0) = n_0(x)$,

$$(2.1) \quad \rho_t + \operatorname{div}(\rho v) = 0,$$

$$(2.2) \quad (\rho u)_t + \operatorname{div}(\rho v \otimes v) + a \nabla \rho^\gamma = \mu \Delta u - \lambda \operatorname{div}(\nabla n \otimes \nabla n - \frac{1}{2} |\nabla n|^2 I),$$

$$(2.3) \quad n_t - \gamma \Delta n = \lambda |\nabla d|^2 d - v \cdot \nabla d,$$

with initial and boundary conditions

$$(2.4) \quad (\rho, u, n)|_{t=0} = (\rho_0 + \delta, u_0, n_0), \quad x \in \Omega,$$

$$(2.5) \quad u(x, t) = u_0(x) = 0, \quad n(x, t) = n_0(x), \quad x \in \partial\Omega.$$

Here $\delta > 0$ is a constant, and $\rho_0 \geq 0$, $|n_0| = 1$.

We use the following notations: Suppose Banach spaces

$$\mathcal{A} = L^\infty(0, T; H^2(\Omega)) \cap L^2(0, T; W^{2,6}(\Omega)) \cap W_2^{1,1}((0, T) \times \Omega),$$

$$\mathcal{B} = L^\infty(0, T; W^{2,6}(\Omega)) \cap W_\infty^{1,1}((0, T) \times \Omega) \cap W_2^{2,1}((0, T) \times \Omega)$$

with norm respectively

$$\|v\|_{\mathcal{A}} = \|v\|_{L^\infty(0,T;H^2(\Omega))} + \|v\|_{L^2(0,T;W^{2,6}(\Omega))} + \|v_t\|_{L^2(0,T;H^1(\Omega))},$$

$$\|d\|_{\mathcal{B}} = \|d_t\|_{L^2(0,T;H^2(\Omega))} + \|d_t\|_{L^\infty(0,T;H^1(\Omega))} + \|d\|_{L^\infty(0,T;W^{2,6}(\Omega))}.$$

Lemma 2.1. *For given v with $\|v\|_{\mathcal{A}} \leq A$, the unique solution ρ of (2.1) satisfies*

$$(2.6) \quad \|\rho\|_{L^\infty(0,T;W^{1,6}(\Omega))} \leq cc_0(1 + T^{\frac{1}{2}}A) \exp(cT^{\frac{1}{2}}A),$$

$$(2.7) \quad \|\rho_t\|_{L^\infty(0,T;L^6(\Omega))} \leq cc_0A \exp(cT^{\frac{1}{2}}A).$$

In particular,

$$(2.8) \quad \|p\|_{L^\infty(0,T;W^{1,6}(\Omega))} \leq cc_0(1 + T^{\frac{1}{2}}A) \exp(cT^{\frac{1}{2}}A),$$

$$(2.9) \quad \|p_t\|_{L^\infty(0,T;L^6(\Omega))} \leq cc_0A \exp(cT^{\frac{1}{2}}A),$$

where c is an absolute constant, perhaps dependent on Ω , λ , μ , γ , etc., and c_0 is a constant dependent on initial and boundary data.

Proof. Since

$$\begin{aligned} \nabla \rho &= \nabla \rho_0 \exp\left(-\int_0^t \operatorname{div} v\right) - \rho_0 \int_0^t \nabla \operatorname{div} v \exp\left(-\int_0^t \operatorname{div} v\right), \\ \rho_t &= -\rho_0 \operatorname{div} v \exp\left(-\int_0^t \operatorname{div} v\right), \end{aligned}$$

we have, from the Minkowski inequality,

$$\begin{aligned} \|\nabla \rho\|_{L^6(\Omega)} &\leq c\|\rho_0\|_{W^{1,6}(\Omega)} \left(1 + \left\|\int_0^t \nabla^2 v\right\|_{L^6(\Omega)}\right) \exp\left(\int_0^T \|\operatorname{div} v\|_{L^\infty(\Omega)}\right) \\ &\leq c\|\rho_0\|_{W^{1,6}(\Omega)} \left(1 + \int_0^T \|\nabla^2 v\|_{L^6(\Omega)}\right) \exp\left(\int_0^T \|\operatorname{div} v\|_{L^\infty(\Omega)}\right) \\ &\leq c\|\rho_0\|_{W^{1,6}(\Omega)} (1 + T^{\frac{1}{2}}\|v\|_X) \exp(cT^{\frac{1}{2}}\|v\|_X) \\ &\leq cc_0(1 + T^{\frac{1}{2}}A) \exp(cT^{\frac{1}{2}}A), \\ \|\rho_t\|_{L^6(\Omega)} &\leq c\|\rho_0\|_{L^\infty(\Omega)} \|\nabla v\|_{L^6(\Omega)} \exp\left(\int_0^T \|\operatorname{div} v\|_{L^\infty(\Omega)}\right) \\ &\leq cc_0 \exp(cT^{\frac{1}{2}}A) \|v\|_{H^2(\Omega)} \leq cc_0A \exp(cT^{\frac{1}{2}}A), \end{aligned}$$

where $X = L^2(0, T; W^{2,6}(\Omega))$. □

Lemma 2.2. *Suppose $\|v\|_{\mathcal{A}} \leq A$, $\|d\|_{\mathcal{B}} \leq B$. Then (2.3) with initial condition $n(x, 0) = n_0(x)$ has a unique solution n and a constant K_1 , depending only on n_0 and u_0 , such that, for $T = T(A, B)$ small enough,*

$$(2.10) \quad \|n\|_{\mathcal{B}} = \|n_t\|_{L^2(0,T;H^2(\Omega))} + \|n_t\|_{L^\infty(0,T;H^1(\Omega))} + \|n\|_{L^\infty(0,T;W^{2,6}(\Omega))} \leq K_1.$$

Proof. The existence of a solution to (2.3) is standard. We just give the estimates as follows. Differentiating (2.3) with respect to time t ,

$$n_{tt} - v\Delta n_t = v(|\nabla d|_t^2 d + |\nabla d|^2 d_t) + (v_t \cdot \nabla)d - (v \cdot \nabla)d_t.$$

Multiplying by Δn_t , integrating over Ω , and using the Cauchy inequality, we get

$$(2.11) \quad \begin{aligned} & \frac{1}{2} \frac{d}{dt} \int_{\Omega} |\nabla n_t|^2 + \nu \int_{\Omega} |\Delta n_t|^2 \\ &= - \int_{\Omega} v(|\nabla d|_t^2 d + |\nabla d|^2 d_t) \cdot \Delta n_t + (v_t \cdot \nabla)d \cdot \Delta n_t - (v \cdot \nabla)d_t \cdot \Delta n_t \\ &\leq \int_{\Omega} 2\nu |\nabla d| |\nabla d_t| |d| |\Delta n_t| + \nu |\nabla d|^2 |d_t| |\Delta n_t| \\ &\quad + \int_{\Omega} |\nabla v_t| |\nabla d| |\nabla n_t| + |v_t| |\nabla^2 d| |\nabla n_t| + |v| |\nabla d_t| |\Delta n_t| \\ &= \sum_{i=1}^5 I_i. \end{aligned}$$

We have the following estimates for I_i :

$$\begin{aligned} I_1 &= \int_{\Omega} 2\nu |\nabla d| |\nabla d_t| |d| |\Delta n_t| \leq c \int_{\Omega} |\nabla d|^2 |\nabla d_t|^2 |d|^2 + \frac{\nu}{6} \|\Delta n_t\|_{L^2(\Omega)}^2, \\ I_2 &= \int_{\Omega} \nu |\nabla d|^2 |d_t| |\Delta n_t| \leq c \int_{\Omega} |\nabla d|^4 |d_t|^2 + \frac{\nu}{6} \|\Delta n_t\|_{L^2}^2, \\ I_3 &= \int_{\Omega} |\nabla v_t| |\nabla d| |\nabla n_t| \leq A^{-2} B^{-2} \int_{\Omega} |\nabla v_t|^2 |\nabla d|^2 + A^2 B^2 \int_{\Omega} |\nabla n_t|^2, \\ I_4 &= \int_{\Omega} |v_t| |\nabla^2 d| |\nabla n_t| \leq A^{-2} B^{-2} \int_{\Omega} |v_t|^2 |\nabla^2 d|^2 + A^2 B^2 \int_{\Omega} |\nabla n_t|^2 \\ &\leq c A^{-2} B^{-2} \|\nabla v_t\|_{L^2}^2 \|\nabla^2 d\|_{L^2} \|\nabla^2 d\|_{L^6} + A^2 B^2 \int_{\Omega} |\nabla n_t|^2, \\ I_5 &= \int_{\Omega} |v| |\nabla d_t| |\Delta n_t| \leq \frac{3}{\nu} \int_{\Omega} |v|^2 |\nabla d_t|^2 + \frac{\nu}{6} \|\Delta n_t\|_{L^2}^2. \end{aligned}$$

Substituting all the estimates into (2.11), we get

$$\begin{aligned} \frac{d}{dt} \int_{\Omega} |\nabla n_t|^2 + \nu \int_{\Omega} |\Delta n_t|^2 &\leq c \int_{\Omega} |\nabla d|^2 |\nabla d_t|^2 |d|^2 + c \int_{\Omega} |\nabla d|^4 |d_t|^2 \\ &\quad + c A^{-2} B^{-2} \int_{\Omega} |\nabla v_t|^2 |\nabla d|^2 + c A^2 B^2 \int_{\Omega} |\nabla n_t|^2 \\ &\quad + c \int_{\Omega} |v|^2 |\nabla d_t|^2 + c A^{-2} B^{-2} \|\nabla v_t\|_{L^2}^2 \|\nabla^2 d\|_{L^2} \|\nabla^2 d\|_{L^6}, \end{aligned}$$

that is,

$$\begin{aligned} \int_{\Omega} |\nabla n_t|^2 + \nu \int_0^T \int_{\Omega} |\Delta n_t|^2 \\ \leq cB^6T + cA^2B^2T + c + cA^2B^2 \int_0^T \int_{\Omega} |\nabla n_t|^2 + c(n_0, u_0), \end{aligned}$$

where

$$\begin{aligned} c(n_0, u_0) \\ = c \int_{\Omega} |\Delta \nabla n_0|^2 + |\nabla n_0|^2 |\nabla^2 n_0|^2 + |\nabla n_0|^6 + c \int_{\Omega} |\nabla u_0|^2 |\nabla n_0|^2 + |u_0|^2 |\nabla^2 n_0|^2. \end{aligned}$$

Using Gronwall's inequality, we obtain

$$\int_{\Omega} |\nabla n_t|^2 \leq (cB^6T + cA^2B^2T + c_0) \exp(cA^2B^2T)$$

and

$$\int_{\Omega} |\nabla n_t|^2 + \nu \int_0^T \int_{\Omega} |\Delta n_t|^2 \leq c(B^6T + A^2B^2T + c_0)(1 + \exp(cA^2B^2T)).$$

Taking $T = T(A, B)$ small, we get

$$\int_{\Omega} |\nabla n_t|^2 + \nu \int_0^T \int_{\Omega} |\Delta n_t|^2 \leq c.$$

The elliptic estimates can be deduced from (2.3):

$$\begin{aligned} \|n\|_{W^{2,6}(\Omega)} &\leq \|n_t\|_{L^6} + \|v \cdot \nabla d\|_{L^6} + \| |\nabla d|^2 d \|_{L^6} + \|n_0\|_{W^{2,6}} \\ &\leq \|v \cdot \nabla d\|_{L^6} + \| |\nabla d|^2 d \|_{L^6} + c_0. \end{aligned}$$

We estimate each item:

$$\begin{aligned} \|v \cdot \nabla d\|_{L^6} \\ = \left(\int_{\Omega} |v|^6 |\nabla d|^6 \right)^{\frac{1}{6}} &\leq \left(\int_{\Omega} |v - u_0|^6 |\nabla d|^6 \right)^{\frac{1}{6}} + \|u_0\|_{L^\infty} \left(\int_{\Omega} |\nabla d|^6 \right)^{\frac{1}{6}} \\ &\leq cB \left(\int_{\Omega} |\nabla v - \nabla u_0|^2 \right)^{\frac{1}{2}} + c\|u_0\|_{L^\infty} \left(\int_{\Omega} |\nabla d - \nabla n_0|^6 \right)^{\frac{1}{6}} + c\|u_0\|_{L^\infty} \|\nabla n_0\|_{L^\infty} \\ &\leq cB \left(\int_{\Omega} \left| \int_0^t \nabla v_t \right|^2 \right)^{\frac{1}{2}} + c_0 B^{\frac{2}{3}} \left(\int_{\Omega} \left| \int_0^t \nabla d_t \right|^2 \right)^{\frac{1}{6}} + c_0 \\ &\leq cBT^{\frac{1}{2}} \|\nabla v_t\|_{L^2(Q_T)} + c_0 T^{\frac{1}{3}} B + c_0 \leq cABT^{\frac{1}{2}} + c_0 BT^{\frac{1}{3}} + c_0 \end{aligned}$$

and

$$\begin{aligned}
 \|\nabla d\|^2 d\|_{L^6} &= \left(\int_{\Omega} |\nabla d|^2 d|^6 \right)^{\frac{1}{6}} \leq \left(\int_{\Omega} |\nabla d|^{12} |d - n_0|^6 \right)^{\frac{1}{6}} + c_0 \left(\int_{\Omega} |\nabla d|^{12} \right)^{\frac{1}{6}} \\
 &\leq cB^2 \left(\int_{\Omega} |d - n_0|^6 \right)^{\frac{1}{6}} + c_0 \left(\int_{\Omega} |\nabla d - \nabla n_0|^{12} \right)^{\frac{1}{6}} + c_0 \\
 &\leq cB^2 \left(\int_{\Omega} |\nabla d - \nabla n_0|^2 \right)^{\frac{1}{2}} + c_0 B \left(\int_{\Omega} |\nabla d - \nabla n_0|^6 \right)^{\frac{1}{6}} + c_0 \\
 &\leq cAB^2 T^{\frac{1}{2}} + c_0 B^2 T^{\frac{1}{3}} + c_0.
 \end{aligned}$$

Taking $T = T(A, B)$ small enough, we obtain the desired $\|n\|_{W^{2,6}} \leq c_0$. □

For (2.2) we have following Lemma.

Lemma 2.3. *Under the conditions of Lemma 2.2, suppose n satisfies (2.3) and ρ (2.1). Then there exists a unique solution u satisfying (2.2), and there is a constant K_2 , depending only on n_0 and u_0 , such that, for $T = T(A, B)$ small enough,*

$$(2.12) \quad \|u\|_{\mathcal{A}} \equiv \|u\|_{L^\infty(0,T;H^2(\Omega))} + \|u\|_{L^2(0,T;W^{2,6}(\Omega))} + \|u_t\|_{L^2(0,T;H^1(\Omega))} \leq K_2.$$

Proof. Since

$$\rho \geq \delta \exp\left(-\int_0^T |\nabla v|_{L^\infty((0,T)\times\Omega)}\right) > 0,$$

the standard theory of parabolic equations implies the existence of the solution to (2.2). Differentiating (2.2) with respect to time t , we get

$$(2.13) \quad \rho u_{tt} - \mu \Delta u_t = -\lambda \operatorname{div}((\nabla d \otimes \nabla d)_t - \frac{1}{2} |\nabla d|_t^2 I) - \nabla p_t - (\rho v \cdot \nabla) v_t - (\rho_t v \cdot \nabla) v - (\rho v_t \cdot \nabla) v - \rho_t u_t.$$

Multiplying by u_t , integrating by parts, and using the continuity of (2.1), we get

$$\begin{aligned}
 &\frac{1}{2} \frac{d}{dt} \int_{\Omega} \rho |u_t|^2 + \mu \int_{\Omega} |\nabla u_t|^2 \\
 &= \lambda \int_{\Omega} ((\nabla d \otimes \nabla d)_t - \frac{1}{2} |\nabla d|_t^2 I) \cdot \nabla u_t \\
 &\quad - \int_{\Omega} \nabla p_t \cdot u_t - (\rho v \cdot \nabla) v_t \cdot u_t - (\rho_t v \cdot \nabla) v \cdot u_t - \int_{\Omega} (\rho v_t \cdot \nabla) v \cdot u_t + \rho_t |u_t|^2 \\
 &\leq 3\lambda \int_{\Omega} |\nabla d| |\nabla d_t| |\nabla u_t| + \int_{\Omega} p_t \operatorname{div}(u_t) + \rho |v| |\nabla v_t| |u_t| \\
 &\quad + \int_{\Omega} \rho |v| |\nabla v|^2 |u_t| + \rho |v|^2 |\nabla^2 v| |u_t| + \rho |v| |\nabla v| |\nabla u_t| \\
 &\quad + \int_{\Omega} \rho |v_t| |\nabla v| |u_t| + 2\rho |v| |\nabla u_t| |u_t| \\
 &= \sum_{i=1}^8 I_i.
 \end{aligned}$$

For each I_i we have

$$I_1 = 3\lambda \int_{\Omega} |\nabla d| |\nabla d_t| |\nabla u_t| \leq c \int_{\Omega} |\nabla d|^2 |\nabla d_t|^2 + \frac{\mu}{12} \int_{\Omega} |\nabla u_t|^2 \leq cB^4 + \frac{\mu}{12} \int_{\Omega} |\nabla u_t|^2,$$

$$\begin{aligned} I_2 &= \int_{\Omega} p_t \operatorname{div}(u_t) \leq c \int_{\Omega} |p_t|^2 + \frac{\mu}{12} \int_{\Omega} |\nabla u_t|^2 \\ &\leq c_0 \exp\left(\int_0^T 2\|\nabla v\|_{L^\infty(\Omega)}\right) \int_{\Omega} |\nabla v|^2 + \frac{\mu}{12} \int_{\Omega} |\nabla u_t|^2 \\ &\leq c_0 A^2 \exp(cAT^{\frac{1}{2}}) + \frac{\mu}{12} \int_{\Omega} |\nabla u_t|^2, \end{aligned}$$

$$I_3 = \int_{\Omega} |\rho| |v| |\nabla v_t| |u_t| \leq A^4 \int_{\Omega} \rho |u_t|^2 + c_0 A^{-2} \exp(cAT^{\frac{1}{2}}) \int_{\Omega} |\nabla v_t|^2,$$

$$I_4 = \int_{\Omega} |\rho| |v| |\nabla v|^2 |u_t| \leq A^6 \int_{\Omega} \rho |u_t|^2 + c_0 \exp(cAT^{\frac{1}{2}}),$$

$$I_5 = \int_{\Omega} |\rho| |v|^2 |\nabla^2 v| |u_t| \leq A^6 \int_{\Omega} \rho |u_t|^2 + c_0 \exp(cAT^{\frac{1}{2}}),$$

$$\begin{aligned} I_6 &= \int_{\Omega} \rho |v| |\nabla v| |\nabla u_t| \leq c \int_{\Omega} \rho^2 |v|^2 |\nabla v|^2 + \frac{\mu}{12} \int_{\Omega} |\nabla u_t|^2 \\ &\leq c_0 A^4 \exp(cAT^{\frac{1}{2}}) + \frac{\mu}{12} \int_{\Omega} |\nabla u_t|^2, \end{aligned}$$

$$\begin{aligned} I_7 &= \int_{\Omega} \rho |v_t| |\nabla v| |u_t| \leq A^4 \int_{\Omega} \rho |u_t|^2 + A^{-4} \int_{\Omega} \rho |v_t|^2 |\nabla v|^2 \\ &\leq A^2 \int_{\Omega} \rho |u_t|^2 + c_0 A^{-2} \exp(cAT^{\frac{1}{2}}) \int_{\Omega} |v_t|^2, \end{aligned}$$

$$\begin{aligned} I_8 &= 2 \int_{\Omega} \rho |v| |\nabla u_t| |u_t| \leq c \int_{\Omega} \rho |u_t|^2 (\rho |v|^2) + \frac{\mu}{12} \int_{\Omega} |\nabla u_t|^2 \\ &\leq c_0 A^2 \exp(cAT^{\frac{1}{2}}) \int_{\Omega} \rho |u_t|^2 + \frac{\mu}{12} \int_{\Omega} |\nabla u_t|^2. \end{aligned}$$

From the above estimates, we get

$$\begin{aligned} &\int_{\Omega} \rho |u_t|^2 + \int_0^T \int_{\Omega} |\nabla u_t|^2 \\ &\leq cB^4 T + c_0 A^4 T \exp(cAT^{\frac{1}{2}}) + c_0 + c_0 A^4 \exp(cAT^{\frac{1}{2}}) \int_0^T \int_{\Omega} \rho |u_t|^2, \end{aligned}$$

which implies that

$$\int_{\Omega} \rho |u_t|^2 + \int_0^T \int_{\Omega} |\nabla u_t|^2 \leq (cB^4 T + c_0 A^4 T \exp(cAT^{\frac{1}{2}})) c_0 A^4 T \exp(cAT^{\frac{1}{2}}).$$

Taking $T = T(A, B)$ small enough, we deduce

$$(2.14) \quad \int_{\Omega} \rho |u_t|^2 + \int_0^T \int_{\Omega} |\nabla u_t|^2 \leq C(c_0).$$

Finally, we estimate

$$\|u\|_{L^\infty(0,T;H^2(\Omega))} \quad \text{and} \quad \|u\|_{L^2(0,T;W^{2,6}(\Omega))}.$$

From (2.2), we get

$$\begin{aligned} & \|u\|_{H^2(\Omega)} \\ & \leq c(\|\nabla p\|_{L^2(\Omega)} + \|\rho u_t\|_{L^2(\Omega)} + \|\nabla^2 n \nabla n\|_{L^2(\Omega)}) + c(\|(\rho v \cdot \nabla)v\|_{L^2(\Omega)} + c_0). \end{aligned}$$

Now we have

$$\begin{aligned} \|\nabla p\|_{L^2(\Omega)} & \leq c_0 \exp(cAT^{\frac{1}{2}}) + c_0 AT^{\frac{1}{2}} \exp(cAT^{\frac{1}{2}}), \\ \|\rho u_t\|_{L^2(\Omega)} & \leq c_0 \exp(cAT^{\frac{1}{2}}) \|\sqrt{\rho} u_t\|_{L^2(\Omega)}, \\ \|\nabla^2 n \nabla n\|_{L^2(\Omega)} & \leq \|\nabla^2 n\|_{L^6(\Omega)} \|\nabla n\|_{L^2(\Omega)}^{\frac{1}{2}} \|\nabla n\|_{L^6(\Omega)}^{\frac{1}{2}} \leq K_1^2, \end{aligned}$$

and

$$\begin{aligned} & \|\rho v \cdot \nabla v\|_{L^2(\Omega)}^2 \\ & \leq \|\rho\|_{L^\infty(\Omega)}^2 \int_{\Omega} |v|^2 |\nabla v|^2 \\ & \leq c_0 \exp(cAT^{\frac{1}{2}}) \left(\int_{\Omega} |v - u_0|^2 |\nabla v|^2 + \|u_0\|_{L^\infty}^2 \int_{\Omega} |\nabla v - \nabla u_0|^2 + c_0 \right) \\ & \leq c_0 \exp(cAT^{\frac{1}{2}}) \left(\int_{\Omega} \left| \int_0^t v_t \right|^2 |\nabla v|^2 + c_0 \int_{\Omega} \left| \int_0^t \nabla v_t \right|^2 + c_0 \right) \\ & \leq c_0 \exp(cAT^{\frac{1}{2}}) (A^4 T + c_0 A^2 T + c_0). \end{aligned}$$

Similarly, we have

$$\begin{aligned} \|\nabla p\|_{L^6(\Omega)} & \leq c_0 \exp(cAT^{\frac{1}{2}}) + c_0 AT^{\frac{1}{2}} \exp(cAT^{\frac{1}{2}}), \\ \|\rho u_t\|_{L^2(0,T;L^6(\Omega))} & \leq c_0 \exp(cAT^{\frac{1}{2}}) \|\nabla u_t\|_{L^2(0,T;L^2(\Omega))} \\ & \leq c_0 \exp(cAT^{\frac{1}{2}}) C(c_0), \\ \|\nabla^2 n \nabla n\|_{L^2(0,T;L^6(\Omega))} & \leq \|\nabla^2 n\|_{L^2(0,T;L^6(\Omega))} \|\nabla n\|_{L^\infty(\Omega)} \leq K_1^2, \end{aligned}$$

and

$$\begin{aligned}
& \|\rho v \cdot \nabla v\|_{L^2(0,T;L^6(\Omega))}^2 \\
& \leq \|\rho\|_{L^\infty(\Omega)}^2 \int_0^T \left(\int_\Omega |v|^6 |\nabla v|^6 \right)^{\frac{1}{3}} \\
& \leq c_0 \exp(cAT^{\frac{1}{2}}) \int_0^T \|v\|_{L^\infty(\Omega)}^2 \|\nabla v\|_{L^\infty(\Omega)}^{\frac{4}{3}} \times \left(\int_\Omega |\nabla v - \nabla u_0|^2 + 1 \right)^{\frac{1}{3}} \\
& \leq c_0 \exp(cAT^{\frac{1}{2}}) A^2 \int_0^T \|\nabla v\|_{L^\infty(\Omega)}^{\frac{4}{3}} \times \left(\int_\Omega \left| \int_0^t \nabla v_t \right|^2 + 1 \right)^{\frac{1}{3}} \\
& \leq c_0 \exp(cAT^{\frac{1}{2}}) \left(T \int_0^T \int_\Omega |\nabla v_t|^2 + 1 \right)^{\frac{1}{3}} \times \left(\int_0^T \|v\|_{W^{2,6}(\Omega)}^2 \right)^{\frac{2}{3}} T^{\frac{1}{3}} \\
& \leq c_0 \exp(cAT^{\frac{1}{2}}) (TA^2 + 1)^{\frac{1}{3}} A^{\frac{4}{3}} T^{\frac{1}{3}}.
\end{aligned}$$

Thus

$$\int_\Omega \rho |u_t|^2 dx + \mu \int_0^T \int_\Omega |\nabla u_t|^2 dx dt + \|u\|_{L^\infty(0,T;H^2(\Omega))} + \|u\|_{L^2(0,T;W^{2,6}(\Omega))} \leq C(c_0).$$

This concludes the proof. \square

If (n^δ, u^δ) denotes a unique solution of (2.2) and (2.3) with

$$\rho(x, 0) = \rho_0 + \delta$$

and initial and boundary conditions, then taking $\delta \rightarrow 0$, we obtain a unique solution (n, u) of the linearized system (2.1)–(2.3) with $\rho(x, 0) = \rho_0$ and initial and boundary conditions such that $\|n\|_{\mathfrak{B}} \leq K_1$, $\|u\|_{\mathfrak{A}} \leq K_2$. So we can define a map

$$\mathcal{T} : \mathfrak{W} \rightarrow \mathfrak{W}, \quad (d, v) \mapsto (n, u),$$

where Banach space

$$\mathfrak{W} = (\mathfrak{A} \otimes \mathfrak{B}) \cap \mathfrak{C} = \mathfrak{A} \otimes \mathfrak{B}$$

with

$$\mathfrak{C} = \{(n, u) : \|(n, u)\|_{\mathfrak{C}} = \|n\|_{L^2(0,T;H^2(\Omega))} + \|u\|_{L^2(0,T;H^1(\Omega))} < \infty\}.$$

The following lemma tells us that the map \mathcal{T} is contracted in the sense of weaker norm for $(d, v) \in \mathfrak{W}$.

Lemma 2.4. *There is a constant $0 < \theta < 1$ such that for any $(d^i, v^i) \in \mathfrak{W}$, $i = 1, 2$,*

$$\|\mathcal{T}(d^1, v^1) - \mathcal{T}(d^2, v^2)\|_{\mathfrak{C}} \leq \theta \|(d^1 - d^2, v^1 - v^2)\|_{\mathfrak{C}}$$

for some small $T > 0$.

Proof. Suppose ρ_i , n^i , and u^i are the solutions to (2.1)–(2.3) corresponding to given $(d^i, v^i) \in \mathfrak{W}$. Define $\rho = \rho_2 - \rho_1$, $d = d^2 - d^1$, $v = v^2 - v^1$, $n = n^2 - n^1$, $u = u^2 - u^1$, and

$$\rho_i = \rho_0 \exp\left(-\int_0^t \operatorname{div} v^i\right),$$

$i = 1, 2$. Then

$$(2.15) \quad \rho_t + \operatorname{div}(\rho v^2) = -\operatorname{div}(\rho_1 v),$$

$$(2.16) \quad n_t - v \Delta n = v |\nabla d^2|^2 d^2 - v |\nabla d^1|^2 d^1 - v^2 \nabla d^2 + v^1 \nabla d^1,$$

$$(2.17) \quad \begin{aligned} \rho_2 u_t - \mu \Delta u &= (\rho_1 - \rho_2) u_t^1 + \rho_1 v^1 \nabla v^1 - \rho_2 v^2 \nabla v^2 + \nabla p_1 \\ &\quad - \nabla p_2 - \lambda \nabla \cdot (\nabla n^2 \otimes \nabla n^2 - \frac{1}{2} |\nabla n^2|^2 I) \\ &\quad + \lambda \nabla \cdot (\nabla n^1 \otimes \nabla n^1 - \frac{1}{2} |\nabla n^1|^2 I). \end{aligned}$$

Multiplying (2.16) by n and integrating over Ω , we get

$$(2.18) \quad \begin{aligned} \frac{1}{2} \frac{d}{dt} \int_{\Omega} |n|^2 dx + v \int_{\Omega} |\nabla n|^2 dx \\ \leq \int_{\Omega} |\nabla d^2|^2 d^2 \cdot n - |\nabla d^1|^2 d^1 \cdot n - v \nabla d^2 \cdot n - v^1 \nabla d \cdot n \\ \leq \eta \int_{\Omega} (|\nabla d|^2 + |\nabla v|^2) + c(\eta, A, B) \int_{\Omega} |n|^2, \end{aligned}$$

where $c(\eta, A, B)(s)$ satisfies

$$(2.19) \quad \int_0^T c(\eta, A, B)(s) ds \leq K_3$$

for small $T = T(A, B, \eta)$, where K_3 is a constant dependent on initial and boundary data c_0 .

Differentiating (2.16) with respect to x_i , multiplying by ∇n , and integrating over Ω , we deduce

$$(2.20) \quad \begin{aligned} \frac{1}{2} \frac{d}{dt} \int_{\Omega} |\nabla n|^2 dx + \frac{v}{2} \int_{\Omega} |\nabla^2 n|^2 dx \\ \leq \eta \int_{\Omega} (|\nabla v|^2 + |\nabla d|^2 + |\nabla^2 d|^2) + c(\eta, A, B) \int_{\Omega} |\nabla n|^2, \end{aligned}$$

where $c(\eta, A, B)$ satisfies (2.19), and we have used the following identities and estimates:

$$\begin{aligned} \nabla d^2 \nabla^2 d^2 d^2 - \nabla d^1 \nabla^2 d^1 d^1 &= \nabla d \nabla^2 d^2 d^1 + \nabla d^1 \nabla^2 d d^1 + \nabla d^1 \nabla^2 d^1 d, \\ |\nabla d^2|^2 \nabla d^2 - |\nabla d^1|^2 \nabla d^1 &= |\nabla d^2|^2 \nabla d + (|\nabla d^2|^2 - |\nabla d^1|^2) \nabla d^1, \end{aligned}$$

and

$$\begin{aligned} \int_{\Omega} |\nabla n|^2 |\nabla^2 d^2|^2 &\leq \left(\int_{\Omega} |\nabla^2 d^2|^6 \right)^{\frac{1}{3}} \left(\int_{\Omega} |\nabla n|^3 \right)^{\frac{2}{3}} \\ &\leq cB^2 \|\nabla n\|_{L^2(\Omega)} \|\nabla^2 n\|_{L^2(\Omega)} \leq \frac{\nu}{2} \int_{\Omega} |\nabla^2 n|^2 + cB^4 \int_{\Omega} |\nabla n|^2. \end{aligned}$$

Multiplying (2.15) by ρ and using the Minkowski inequality, we have

$$\begin{aligned} \frac{d}{dt} \int_{\Omega} \frac{1}{2} |\rho|^2 &= \int_{\Omega} -\frac{1}{2} |\rho|^2 \operatorname{div} v^2 - \int_{\Omega} \rho (\nabla \rho_1 v + \rho_1 \operatorname{div} v) \\ &\leq c \int_{\Omega} |\rho|^2 |\nabla v^2| + c \|\rho\|_{L^2(\Omega)} \|\nabla \rho_1\|_{L^3(\Omega)} \|v\|_{L^6(\Omega)} \\ &\quad + c \|\rho\|_{L^2(\Omega)} \|\rho_1\|_{L^\infty(\Omega)} \|\nabla v\|_{L^2(\Omega)} \\ &\leq c \|v^2\|_{W^{2,6}(\Omega)} \|\rho\|_{L^2(\Omega)}^2 + \eta \|\nabla v\|_{L^2(\Omega)}^2 \\ &\quad + c_0 \eta^{-1} \exp(cAT^{\frac{1}{2}}) \left(1 + \left\| \int_0^t \nabla^2 v^1 \right\|_{L^3(\Omega)}^2 \right) \|\rho\|_{L^2(\Omega)}^2 \\ &\leq \eta \|\nabla v\|_{L^2(\Omega)}^2 + c \|v^2\|_{W^{2,6}(\Omega)} \|\rho\|_{L^2(\Omega)}^2 \\ &\quad + c_0 \eta^{-1} \exp(cAT^{\frac{1}{2}}) (1 + T \|\nabla^2 v^1\|_{L^2(0,T;L^6(\Omega))}^2) \|\rho\|_{L^2(\Omega)}^2 \\ &\leq c_0 \eta^{-1} \exp(cAT^{\frac{1}{2}}) (1 + TA^2 + \|v^2\|_{W^{2,6}(\Omega)}) \|\rho\|_{L^2(\Omega)}^2 + \eta \|\nabla v\|_{L^2(\Omega)}^2, \end{aligned}$$

that is,

$$(2.21) \quad \frac{d}{dt} \int_{\Omega} \frac{1}{2} |\rho|^2 \leq \eta \|\nabla v\|_{L^2(\Omega)}^2 + c(\eta, A, T) \|\rho\|_{L^2(\Omega)}^2,$$

where $c(\eta, A, T)$ satisfies (2.19).

Multiplying (2.17) by u and integrating over Ω , we deduce

$$\begin{aligned} (2.22) \quad &\frac{1}{2} \frac{d}{dt} \int_{\Omega} \rho_2 |u|^2 dx + \mu \int_{\Omega} |\nabla u|^2 dx \\ &= \int_{\Omega} -\rho_2 v^2 u \nabla u + (\rho_1 - \rho_2) u_t^1 \cdot u + \rho_1 v^1 \nabla v^1 \cdot u - \rho_2 v^2 \nabla v^2 \cdot u + (p_2 - p_1) \operatorname{div} u \\ &\quad + \lambda (\nabla n^2 \otimes \nabla n^2 - \frac{1}{2} |\nabla n^2|^2 I) \nabla u - \lambda (\nabla n^1 \otimes \nabla n^1 - \frac{1}{2} |\nabla n^1|^2 I) \nabla u \\ &= \int_{\Omega} -\rho_2 v^2 u \nabla u + (\rho_1 - \rho_2) (u_t^1 + v^1 \nabla v^1) \cdot u \\ &\quad - \rho_2 (v \nabla v^2 + v^1 \nabla v) \cdot u + (p_1 - p_2) \operatorname{div} u \\ &\quad + \lambda (\nabla n^2 \otimes \nabla n^2 - \frac{1}{2} |\nabla n^2|^2 I) \nabla u - \lambda (\nabla n^1 \otimes \nabla n^1 - \frac{1}{2} |\nabla n^1|^2 I) \nabla u \\ &\leq \eta \int_{\Omega} |\nabla v|^2 + \frac{2\mu}{3} \int_{\Omega} |\nabla u|^2 + c(\eta, A, B) \int_{\Omega} \rho_2 |u|^2 + |\rho|^2 + |\nabla n|^2, \end{aligned}$$

where $c(\eta, A, B)$ satisfying (2.19). Here we have used the key estimates

$$\begin{aligned}
 \int_{\Omega} \rho_2 |v \nabla v^2 + v^1 \nabla v| |u| &\leq \|\nabla v^2\|_{L^6(\Omega)} \|\rho_2 u\|_{L^2(\Omega)} \|v\|_{L^6(\Omega)} \\
 &\quad + \|\nabla v\|_{L^2(\Omega)} \|\rho_2 u\|_{L^2(\Omega)} \|v^1\|_{L^\infty(\Omega)} \\
 &\leq c_0 \exp(cAT^{\frac{1}{2}}) \|\sqrt{\rho_2} u\|_{L^2(\Omega)} \|\nabla v^2\|_{H^2(\Omega)} \|\nabla v\|_{L^2(\Omega)} \\
 &\quad + c_0 \exp(cAT^{\frac{1}{2}}) \|\sqrt{\rho_2} u\|_{L^2(\Omega)} \|v^1\|_{H^2(\Omega)} \|\nabla v\|_{L^2(\Omega)} \\
 &\leq \eta \|\nabla v\|_{L^2(\Omega)}^2 + c\eta^{-1} A^2 \exp(cAT^{\frac{1}{2}}) \|\sqrt{\rho_2} u\|_{L^2(\Omega)}^2, \\
 \int_{\Omega} |\nabla n| |\nabla u| |\nabla n^2| &\leq \eta \int_{\Omega} |\nabla u|^2 + c\eta^{-1} \|\nabla n^2\|_{L^\infty(\Omega)}^2 \int_{\Omega} |\nabla n|^2 \\
 &\leq \frac{\mu}{3} \int_{\Omega} |\nabla u|^2 + cB^2 \int_{\Omega} |\nabla n|^2,
 \end{aligned}$$

and

$$\begin{aligned}
 \int_{\Omega} (\rho_1 - \rho_2)(u_t^1 + v^1 \nabla v^1) \cdot u &\leq \|\rho\|_{L^{\frac{3}{2}}(\Omega)} \|u_t^1 + v^1 \nabla v^1\|_{L^6(\Omega)} \|u\|_{L^6(\Omega)} \\
 &\leq c \|\rho\|_{L^2(\Omega)} \|u_t^1 + v^1 \nabla v^1\|_{H^1(\Omega)} \|\nabla u\|_{L^2(\Omega)} \\
 &\leq \frac{\mu}{3} \|\nabla u\|_{L^2(\Omega)}^2 + c(A, T)(t) \|\rho\|_{L^2(\Omega)}^2,
 \end{aligned}$$

where $c(\eta, A, T)(t)$ satisfies (2.19).

Summing inequalities (2.18) and (2.20)–(2.22), we obtain

$$\begin{aligned}
 \frac{d}{dt} \int_{\Omega} |\rho|^2 + |n|^2 + |\nabla n|^2 + \rho_2 |u|^2 + \int_{\Omega} |\nabla n|^2 + |\nabla^2 n|^2 + |\nabla u|^2 \\
 \leq c\eta \int_{\Omega} |\nabla v|^2 + |\nabla d|^2 + |\nabla^2 d|^2 + c(\eta, A, B, T) \int_{\Omega} |\rho|^2 + |n|^2 + |\nabla n|^2 + \rho_2 |u|^2,
 \end{aligned}$$

which implies, by (2.19) and taking $T = T(\eta, A, B)$ small enough,

$$\begin{aligned}
 \int_{\Omega} |\rho|^2 + |n|^2 + |\nabla n|^2 + \rho_2 |u|^2 \\
 \leq \eta \exp\left(\int_0^T c(\eta, A, B)(s) ds\right) \int_0^T \int_{\Omega} |\rho|^2 + |n|^2 + |\nabla n|^2 + \rho_2 |u|^2 \\
 \leq c\eta \int_0^T \int_{\Omega} |\rho|^2 + |n|^2 + |\nabla n|^2 + \rho_2 |u|^2.
 \end{aligned}$$

Thus, taking η small, we obtain

$$(2.23) \quad \|\rho\|_{L^\infty(0, T; L^2(\Omega))} + \|n\|_{L^\infty(0, T; H^1(\Omega))} + \|\sqrt{\rho_2} u\|_{L^\infty(0, T; L^2(\Omega))} \leq c$$

and

$$\int_0^T \int_{\Omega} |\nabla n|^2 + |\nabla^2 n|^2 + |\nabla u|^2 \leq \theta \int_0^T \int_{\Omega} |\nabla d|^2 + |\nabla^2 d|^2 + |\nabla v|^2$$

with $0 < \theta < 1$. Since n and u are zero on boundary, we finish the proof. \square

3. Proof of Theorem 1.1

Proof. By the contractibility of \mathcal{T} , we can easily obtain a unique solution (n, u) of (1.3) and (1.2), and ρ is from u by formula (1.8), that is, ρ is a unique solution of (1.1). Lemmas 2.1–2.3 and the lower semicontinuity of norms imply that the solutions (ρ, n, u) satisfy the same estimates. Multiplying (1.3) by n , we get

$$|n|_t^2 + (u \cdot \nabla)|n|^2 = v \Delta |n|^2 + (|n|^2 - 1)|\nabla n|^2,$$

that is,

$$(|n|^2 - 1)_t + (u \cdot \nabla)(|n|^2 - 1) = v \Delta (|n|^2 - 1) + (|n|^2 - 1)|\nabla n|^2.$$

Define $D = (|n|^2 - 1) \exp(\|\nabla n\|_{L^\infty(Q_T)}^2 t)$, where $Q_T = \Omega \times [0, T]$. Then

$$D_t + (u \cdot \nabla)D = v \Delta D + (|\nabla n|^2 - \|\nabla n\|_{L^\infty(Q_T)}^2)D$$

with $D|_{\partial\Omega} = 0$. So from the maximum principle of parabolic equations, we deduce

$$D \equiv 0 \quad \text{in } ((0, T) \times \Omega).$$

Thus we complete the proof of the theorem. \square

References

- [Calderer and Liu 2000] M. C. Calderer and C. Liu, “Liquid crystal flow: dynamic and static configurations”, *SIAM J. Appl. Math.* **60**:6 (2000), 1925–1949 (electronic). MR 1763310 (2001e:76009) Zbl 0956.35104
- [Cho and Kim 2006] Y. Cho and H. Kim, “Existence results for viscous polytropic fluids with vacuum”, *J. Differential Equations* **228**:2 (2006), 377–411. MR 2289539 (2007j:35155) Zbl 1135.35071
- [Cho et al. 2004] Y. Cho, H. J. Choe, and H. Kim, “Unique solvability of the initial boundary value problems for compressible viscous fluids”, *J. Math. Pures Appl. (9)* **83**:2 (2004), 243–275. MR 2038120 (2005a:76133) Zbl 1080.35066
- [Choe and Kim 2003] H. J. Choe and H. Kim, “Strong solutions of the Navier-Stokes equations for isentropic compressible fluids”, *J. Differential Equations* **190**:2 (2003), 504–523. MR 1970039 (2004b:35258) Zbl 1022.35037
- [Ericksen 1962] J. L. Ericksen, “Hydrostatic theory of liquid crystals”, *Arch. Rational Mech. Anal.* **9** (1962), 371–378. MR 0137403 (25 #855) Zbl 0105.23403

- [Ericksen and Kinderlehrer 1987] J. L. Ericksen and D. Kinderlehrer (editors), *Theory and applications of liquid crystals*, vol. 5, The IMA Volumes in Mathematics and its Applications, Springer, New York, 1987. Papers from the IMA workshop held in Minneapolis, Minn., January 21–25, 1985, Edited by J. L. Ericksen and D. Kinderlehrer. MR 900827 (88d:82007) Zbl 0713.76006
- [Feireisl 2004] E. Feireisl, *Dynamics of viscous compressible fluids*, vol. 26, Oxford Lecture Series in Mathematics and its Applications, Oxford University Press, Oxford, 2004. MR 2040667 (2005i:76092) Zbl 1080.76001
- [Feireisl et al. 2001] E. Feireisl, A. Novotný, and H. Petzeltová, “On the existence of globally defined weak solutions to the Navier-Stokes equations”, *J. Math. Fluid Mech.* **3**:4 (2001), 358–392. MR 1867887 (2002k:35253) Zbl 0997.35043
- [Forster et al. 1971] D. Forster, T. Lubensky, P. Martin, J. Swift, and P. Pershan, “Hydrodynamics of liquid crystals”, *Phys. Rev. Lett.* **26**:17 (1971), 1016–1019.
- [Frank 1958] F. C. Frank, “On the theory of liquid crystals”, *Discussions Faraday Soc.* **25** (1958), 19–28.
- [Gennes and Prost 1993] P. G. de Gennes and J. Prost, *The physics of liquid crystals*, 2nd ed., International Series of Monographs on Physics (Oxford, England) **83**, Clarendon, Oxford, 1993.
- [Leslie 1966] F. M. Leslie, “Some constitutive equations for anisotropic fluids”, *Quart. J. Mech. Appl. Math.* **19** (1966), 357–370. MR 0207302 (34 #7118) Zbl 0148.20504
- [Leslie 1968] F. M. Leslie, “Some constitutive equations for liquid crystals”, *Arch. Rational Mech. Anal.* **28**:4 (1968), 265–283. MR 1553506 Zbl 0159.57101
- [Lin 1989] F.-H. Lin, “Nonlinear theory of defects in nematic liquid crystals; phase transition and flow phenomena”, *Comm. Pure Appl. Math.* **42**:6 (1989), 789–814. MR 1003435 (90g:82076) Zbl 0703.35173
- [Lin and Liu 1995] F.-H. Lin and C. Liu, “Nonparabolic dissipative systems modeling the flow of liquid crystals”, *Comm. Pure Appl. Math.* **48**:5 (1995), 501–537. MR 1329830 (96a:35154) Zbl 0842.35084
- [Lin and Liu 1996] F.-H. Lin and C. Liu, “Partial regularity of the dynamic system modeling the flow of liquid crystals”, *Discrete Contin. Dynam. Systems* **2**:1 (1996), 1–22. MR 1367385 (96m:35255) Zbl 0948.35098
- [Lin and Liu 2000] F.-H. Lin and C. Liu, “Existence of solutions for the Ericksen-Leslie system”, *Arch. Ration. Mech. Anal.* **154**:2 (2000), 135–156. MR 1784963 (2003a:76014) Zbl 0963.35158
- [Lin and Liu 2001] F. Lin and C. Liu, “Static and dynamic theories of liquid crystals”, *J. Partial Differential Equations* **14**:4 (2001), 289–330. MR 1883167 (2003b:82063)
- [Lions 1998] P.-L. Lions, *Mathematical topics in fluid mechanics. Vol. 2*, vol. 10, Oxford Lecture Series in Mathematics and its Applications, The Clarendon Press Oxford University Press, New York, 1998. Compressible models, Oxford Science Publications. MR 1637634 (99m:76001) Zbl 0908.76004
- [Liu and Qing 2011] X. Liu and J. Qing, “Globally weak solutions to the flow of compressible liquid crystals system”, preprint, 2011. To appear in *Discret. Contin. Dyn. Syst. A*.
- [Liu and Zhang 2009] X. Liu and Z. Zhang, “Global existence of weak solutions for the incompressible liquid crystals”, *Chinese Ann. Math.* **30**:1 (2009), 1–20.
- [Oseen 1933] C. W. Oseen, “The theory of liquid crystals”, *Trans. Faraday Soc.* **29** (1933), 883–899. Zbl 0008.04203
- [Pasechnik et al. 2009] S. V. Pasechnik, V. G. Chigrinov, and D. V. Shmeliova, *Liquid crystals: viscous and elastic properties*, Wiley-VCH, Weinheim, 2009. Zbl 0999.35078

[Stephen 1970] M. J. Stephen, “Hydrodynamics of liquid crystals”, *Phys. Rev. A* **2**:4 (1970), 1558–1562.

[Xie 1988] Y. Z. Xie, *The physics of liquid crystals*, Scientific Press, Beijing, 1988.

Received May 16, 2011. Revised November 12, 2011.

YU-MING CHU
SCHOOL OF MATHEMATICS AND COMPUTATION SCIENCES
HUNAN CITY UNIVERSITY
YIYANG, 413000
CHINA
chuyuming@hutc.zj.cn

XIAN-GAO LIU
INSTITUTE OF MATHEMATICS
FUDAN UNIVERSITY
SHANGHAI, 200433
CHINA
xgliu@fudan.edu.cn

XIAO LIU
INSTITUTE OF MATHEMATICS
FUDAN UNIVERSITY
SHANGHAI, 200433
CHINA
shaw0820@gmail.com

PRESENTATIONS FOR THE HIGHER-DIMENSIONAL THOMPSON GROUPS nV

JOHANNA HENNIG AND FRANCESCO MATUCCI

M. G. Brin has introduced the higher-dimensional Thompson groups nV that are generalizations to the Thompson group V of self-homeomorphisms of the Cantor set and found a finite set of generators and relations in the case $n = 2$. We show how to generalize his construction to obtain a finite presentation for every positive integer n . As a corollary, we obtain another proof that the groups nV are simple (first proved by Brin).

1. Introduction

The higher-dimensional groups nV were introduced by Brin in [2004; 2005] and generalize Thompson's group V . The group V is a group of self-homeomorphisms of the Cantor set \mathcal{C} that is simple and finitely presented — the standard introduction to V is the paper by Cannon, Floyd and Parry [1996]. The groups nV generalize the group V and act on powers of the Cantor set \mathcal{C}^n . Brin shows in [2004] that the groups V and $2V$ are not isomorphic and shows in [2005] that the group $2V$ is finitely presented. Bleak and Lanoue [2010] have recently shown that two groups mV and nV are isomorphic if and only if $m = n$.

In this paper we give a finite presentation for each of the higher-dimensional Thompson groups nV . The argument extends to the ascending union ωV of the groups nV and returns an infinite presentation of the same flavor. As a corollary, we obtain another proof that the groups nV and ωV are simple. Our arguments follow closely and generalize those of Brin in [2004; 2005] for the group $2V$.

This work arose during a Research Experience for Undergraduates program at Cornell University. The motivation for the project sprang from a commonly held opinion that the bookkeeping required to generalize Brin's presentations to the groups nV would be overwhelming. One would expect from the similarity of the groups' constructions that all arguments for $2V$ would carry over to nV for all n . Standing in the way of this are the cross relations. Thus our paper has two kinds

Partially supported by the NSF grant for Research Experiences for Undergraduates (REU).

MSC2010: 20F05, 20F65.

Keywords: Thompson groups, groups of piecewise-linear homeomorphisms, finiteness properties, finite presentations.

of arguments: those that verify the parts of [Brin 2005] that carry over with no change to nV and those involving the cross relations that have to be modified to hold in nV (see Lemmas 6 and 20 and Remark 13 below).

Following a suggestion of Collin Bleak the authors have also explored an alternative generating set (see Section 8). An interesting project would be to find a set of relators for this alternative generating set in order to use a known procedure that significantly reduces the number of relations, and which has been successfully implemented in a number of papers by Guralnick, Kantor, Kassabov and Lubotzky; see for example [Guralnick et al. 2011].

After a careful reading of Brin’s original paper [2005], it became clear what was needed to generalize his proof, and the current paper borrows heavily from Brin’s. Brin was already aware that many of his arguments would probably extend (and he points out in several places in [2004; 2005] where it is evident that they do). We show how to deal with generators in higher dimensions and what steps are needed to obtain the same type of normalized words that are built for $2V$ in [Brin 2005].

We also mention that Brin asks in [2005] whether or not the group $2V$ has type F_∞ (that is, it has a classifying space that is finite in each dimension). This has recently been answered by Kochloukova, Martinez-Perez and Nucinkis [2010], who have shown that the groups $2V$ and $3V$ have type F_∞ , therefore obtaining a new proof that these groups are finitely presented.

2. The main ingredient and structure of this paper

Many arguments of Brin [2004; 2005] generalize verbatim from $2V$ to nV . The key observation that allows us to restate many results without proof (or with little additional effort) is the following: Many statements of Brin do not depend on dimension 2, except those that need to make use of the “cross relation” (relation (18) in Section 4 below) to rewrite a cut in dimension d followed by a cut in dimension d' as one in dimension d' followed by one in dimension d .

As a result, proofs that need to make use of this new relation require a slight generalization (for example, the normalization of words in the monoid across fully divided dimensions) while those that do not can be obtained directly using Brin’s original proof. In any case, since statements need to be adapted to our context we sketch certain proofs to make it clear that they generalize directly. For example, we will show why Brin’s proof that $2V$ is simple does not use the new relation (18) and therefore it lifts immediately to higher dimensions.

3. The monoid Π_n

In [2004, Section 4.5], Brin defines the monoid Π and $\widehat{2V}$ and observes that one can extend the definition for all n . Elements of Π_n are given by numbered patterns

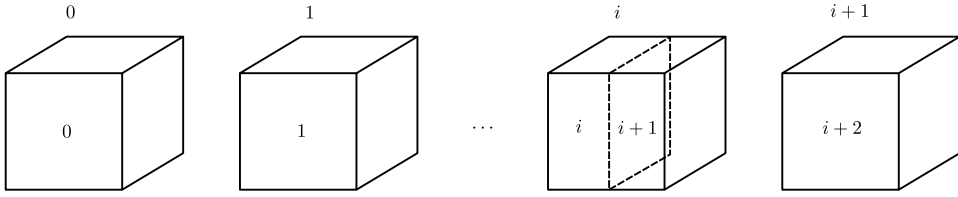


Figure 1. The generator $s_{i,d}$.

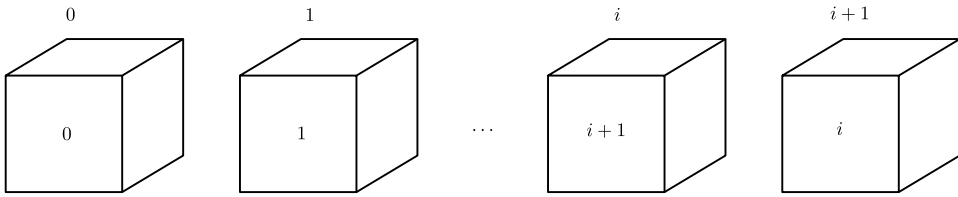


Figure 2. The generator σ_i .

in X , where X is the union of the set $\{S_0, S_1, \dots\}$ of unit n -cubes. Fix $n \in \mathbb{N}$ and fix an ordering on the dimensions d for $1 \leq d \leq n$. The monoid Π_n is generated by the elements $s_{i,d}$ and σ_i , where $s_{i,d}$ denotes the element that cuts the rectangle S_i in half across the d -th dimension (see Figure 1) and σ_i is the transposition that switches the rectangle labeled i with that labeled $i + 1$, as defined for $2V$ (see Figure 2).

After each cut, the numbering shifts as before. The following relations hold in Π_n .

- (M1) $s_{j,d'}s_{i,d} = s_{i,d}s_{j+1,d'}$ for $i < j, 1 \leq d, d' \leq n$,
- (M2) $\sigma_i^2 = 1$ for $i \geq 0$,
- (M3) $\sigma_i\sigma_j = \sigma_j\sigma_i$ for $|i - j| \geq 2$,
- (M4) $\sigma_i\sigma_{i+1}\sigma_i = \sigma_{i+1}\sigma_i\sigma_{i+1}$ for $i \geq 0$,
- (M5a) $\sigma_js_{i,d} = s_{i,d}\sigma_{j+1}$ for $i < j$,
- (M5b) $\sigma_js_{i,d} = s_{j+1,d}\sigma_j\sigma_{j+1}$ for $i = j$,
- (M5c) $\sigma_js_{i,d} = s_{j,d}\sigma_{j+1}\sigma_j$ for $i = j + 1$,
- (M5d) $\sigma_js_{i,d} = s_{i,d}\sigma_j$ for $i > j + 1$,
- (M6) $s_{i,d}s_{i+1,d'}s_{i,d'} = s_{i,d'}s_{i+1,d}s_{i,d}\sigma_{i+1}$ for $i \geq 0, d \neq d'$.

Relations (M5b) and (M5c) are actually equivalent, because σ_i is its own inverse.

Remark 1. The proofs of [Brin 2005, Section 2] that use relations (M1)–(M5d) do not depend on the dimension being 2. For this reason, they generalize immediately

to the case of the monoid Π_n and we do not prove them again. This includes every result up to and including [Brin 2005, Lemma 2.9].

On the other hand, Proposition 2.11 in [Brin 2005] uses the cross relation (M6) and it requires us to choose how we write elements to get some underlying pattern. Brin achieves this type of normalization by writing elements so that vertical cuts appear first, whenever possible. We generalize his argument by describing how to order nodes in forests (which represent cuts in some dimension).

The following definition is given inductively on the subtrees.

Definition 2. Given a forest F , we say that a subtree T of some tree of F is *fully divided* across some dimension d if the root of T is labeled d or if both its left and right subtrees are fully divided across dimension d . We say a forest F is *normalized* if every subtree T is such that if T is fully divided across different the dimensions $d_1 < d_2 < \dots < d_u$, then the root of T is labeled with d_1 , the lowest among all possible dimensions over which T is fully divided.

Given that a word w is a word in the generators $\{s_{i,d}, \sigma_i\}$, we define the *length* $\ell(w)$ of w to be the number of times an element of $\{s_{i,d}\}$ appears in w . It can easily be seen that the length of a word is preserved by relations (M1)–(M6).

We restate some results adapted to our case.

Lemma 3 [Brin 2005, Lemma 2.7]. *If the numbered, labeled forest F comes from a word in $\{s_{i,d} \mid d, i \in \mathbb{N}\}$, then the leaves of F are numbered so that the leaves in F_i have numbers lower than those in F_j whenever $i < j$ and the leaves in each tree of F are numbered in increasing order under the natural left-right ordering of the leaves.*

Lemma 4 [Brin 2005, Lemma 2.8]. *If two words in the generators*

$$\{s_{i,d}, \sigma_i \mid i \in \mathbb{N}, 1 \leq d \leq n\}$$

lead to the same numbered, labeled forest, then they are related by (M1)–(M5d).

Lemma 5 [Brin 2005, Lemma 2.9]. *If F is a numbered, labeled forest with the numbering as in Lemma 3, and if a linear order is given on the interior vertices (and thus of the carets) of F that respects the ancestor relation, then there is a unique word w in $\{s_{i,d} \mid d, i \in \mathbb{N}\}$ leading to F such that the order on the interior vertices of F derived from the order on the entries in w is identical to the given linear order on the interior vertices.*

The next lemma and corollary are used to prove results analogous to [Brin 2005, Lemma 2.10 and Proposition 2.11].

Lemma 6. *Let w be a word in the set $\{s_{i,d}, \sigma_i\}$ and suppose that the underlying pattern P has a fully divided hypercube S_i across dimension d . Then $w \sim w' = s_{i,d}a$ for some word $a \in \langle s_{i,d}, \sigma_i \rangle$.*

Proof. We use induction on $g := \ell(w)$. By using relations (M5a)–(M5d) as in [Brin 2005, Lemma 2.3] we can assume that $w = pq$, where $p \in \langle s_{i,d} \rangle$ and $q \in \langle \sigma_i \rangle$. This does not alter the length of w . If $g = 3$, then $p = p_1 p_2 p_3$. If $p_1 = s_{i,d}$, we are done; otherwise we have two cases: either $p_2 = s_{i+1,d}$ and $p_3 = s_{i,d}$ or $p_2 = s_{i,d}$ and $p_3 = s_{i+2,d}$. Up to using relation (M1), we can assume that $p_2 = s_{i+1,d}$ and $p_3 = s_{i,d}$ which is what we want to apply relation (M6) to p to get $w \sim w' = s_{i,d} s_{i+1,k} s_{i,k} q$.

Now assume the thesis true for all words of length less than g . We consider the word p and look at the labeled unnumbered tree F_i corresponding to S_i with root vertex u and children u_0 and u_1 . Let T_r be the subtree of F_i with root vertex u_r for $r = 0, 1$. We choose an ordering of the vertices of F_i that respects the ancestor relation and such that u corresponds to 1, u_0 corresponds to 2, the other interior nodes of T_0 correspond to the numbers from 3 to $j = \#(\text{interior nodes of } T_0)$ and u_2 corresponds to $j + 1$.

By Lemma 5, the word p is equivalent to

$$p \sim s_{i,k}(s_{i,m}p_0)(s_{f,l}p_1),$$

where $s_{i,m}p_0$ is the subword corresponding to the subtree T_0 and $s_{f,l}p_1$ is the subword corresponding to the subtree T_1 and with $p_0, p_1 \in \langle s_{i,d} \rangle$. We observe that

$$\ell(s_{i,m}p_0) < \ell(p) = g \quad \text{and} \quad \ell(s_{f,l}p_1) < \ell(p) = g$$

and that the underlying squares S_i for $s_{i,m}p_0$ and S_{i+1} for $s_{f,l}p_1$ are fully divided across dimension d . We can thus apply the induction hypothesis and rewrite

$$s_{i,m}p_0 \sim s_{i,d}\tilde{p}_0\tilde{q}_0 \quad \text{and} \quad s_{f,l}p_1 \sim s_{f,d}\tilde{p}_1\tilde{q}_1.$$

We restrict our attention to the subword $s_{i,d}\tilde{p}_0\tilde{q}_0s_{f,d}$. Using the relations (M5a)–(M5d), we can move \tilde{q}_0 to the right of $s_{f,d}$ and obtain

$$s_{i,d}\tilde{p}_0\tilde{q}_0s_{f,d} \sim s_{i,d}\tilde{p}_0s_{g,d}\tilde{q}$$

for some permutation word \tilde{q} . Since the word \tilde{p}_0 acts on the rectangle S_i and $s_{g,d}$ acts on the rectangle S_{i+1} , we can apply Lemma 4 and 5 and put a new order on the nodes so that the node corresponding to $s_{i,d}$ is 1 and $s_{g,d}$ is 2. Thus we have

$$s_{i,d}\tilde{p}_0s_{g,d}\tilde{q} \sim s_{i,d}s_{i+2,d}\tilde{p}\tilde{q}$$

for some \tilde{p} word in the set $\{s_{i,d}\}$. Thus we have $w \sim w'' = s_{i,k}s_{i,d}s_{i+2,d}\tilde{p}\tilde{q}$ and so, by applying the cross relation (M6) to the first three letters of w'' , we get

$$w \sim w'' \sim w' = s_{i,d}s_{i,k}s_{i+2,k}\tilde{p}\tilde{q} = s_{i,d}a. \quad \square$$

We have now proved [Brin 2005, Lemma 2.10], since in order for a tree in a forest to be nonnormalized, one of the rectangles in the pattern corresponding to that tree must be fully divided across two different dimensions.

Lemma 7 [Brin 2005]. *If two different forests correspond to the same pattern in X , then at least one of the two forests is not normalized.*

Remark 8. Lemma 6 is used in our extension of [Brin 2005, Proposition 2.11] so that we can push dimension d under the root. This is explained better in the following corollary.

Corollary 9. *Let w be a word in the generators $\{s_{i,d}, \sigma_i\}$ such that its underlying square S_i is fully divided across dimensions d and ℓ . Then*

$$w \sim w' = s_{i,d}s_{i,\ell}s_{i+2,\ell}a \sim w'' = s_{i,\ell}s_{i,d}s_{i+2,d}b$$

for some suitable words a and b in the generators $\{s_{i,d}, \sigma_i\}$.

Proof. This is achieved by a repeated application of Lemma 6. We apply it to w and obtain $w \sim s_{i,d}a_1$. By construction, the underlying squares S_i and S_{i+1} of a_1 are fully divided across dimension ℓ , so we can apply the previous lemma to a_1 to get $a_1 \sim s_{i,\ell}a_2$ and finally we apply it again to $a_2 \sim s_{i+2,\ell}a$. Hence $w \sim w' = s_{i,\ell}s_{i+2,\ell}a$. To get w'' we apply the cross relation (M6) to the subword $s_{i,\ell}s_{i,d}s_{i+2,d}$. \square

Proposition 10. *A word w is related by (M1)–(M6) to a word corresponding to a normalized, labeled forest.*

Proof. We proceed by induction on the length of w . Let g be the length of w and assume the result holds for all words of length less than g . As before, write $w = pq$, where $p = s_{i_0}s_{i_1} \cdots s_{i_{n-1}}$ (here, the i_j refers to the cube that is being cut; we omit the second index indicating dimension as it is unimportant for now). Write $w = s_{i_0}w'$; since the order of the interior vertices of the forest for p given by the order of the letters in p must respect the ancestor relation, we know that the interior vertex corresponding to s_{i_0} must be a root of some tree T . As w' is a word of length less than g , we may apply our inductive hypothesis and assume that w' can be rewritten via relations (M1)–(M6) to obtain a corresponding normalized forest. The pattern P for w is obtained from the pattern P' for w' by applying the pattern of P' in unit square S_i to the rectangle numbered i in the pattern for s_{i_0} . The forest F for w is obtained from the forest F' for w' by attaching the i -th tree of F' to the i -th leaf of the forest for s_{i_0} . Since F' is normalized, it is seen that F has all interior vertices normalized except possibly for the root vertex of one tree, T .

Let u be the root vertex of T with label k and with children u_1 and u_2 . Let T_1 and T_2 be the subtrees of T whose roots are u_1 and u_2 , respectively. By hypothesis, T_1 and T_2 are already normalized. If T is not normalized already, then T must

be fully divided across the dimension k that u is labeled with, and some other dimension less than k . Let d be the minimal dimension across which T is fully divided. Since T_1 and T_2 are also fully divided across d , by Lemma 6, we may apply relations (M1)–(M6) to the subwords of w corresponding to T_1 and T_2 until u_1 and u_2 are each labeled d . Now by [Brin 2005, Lemma 2.9], we may assume $w = s_{i_0,k} s_{i_0,d} s_{i_0+2,d} w''$, where w'' is the remainder of w . We apply relation (M6) to obtain

$$w = s_{i_0,d} s_{i_0,k} s_{i_0+2,k} \sigma_{i_0} w''.$$

Now, we have normalized the vertex u , and we may now use the inductive hypothesis to renormalize the trees T_1 and T_2 . The result is a normalized forest. \square

The proof of the next result follows the argument of [Brin 2005, Theorem 1], using [Lemma 2.10] and Proposition 10 (to extend [Proposition 2.11]).

Theorem 11. *The monoid Π_n is presented by using the generators $\{s_{i,d}, \sigma_i\}$ and relations (M1)–(M6).*

4. Relations in nV

4.1. Generators for nV . The following generators are defined as in [Brin 2004] and analogous arguments show why they are a generating set for nV .

$$\begin{aligned} X_{i,d} &= (s_{0,1}^{i+1} s_{1,d}, s_{0,1}^{i+2}) && \text{for } i \geq 0, 1 \leq d \leq n, \\ C_{i,d} &= (s_{0,1}^i s_{0,d}, s_{0,1}^{i+1}) && \text{for } i \geq 0, 2 \leq d \leq n, && \text{(baker's maps),} \\ \pi_i &= (s_{0,1}^{i+2} \sigma_1, s_{0,1}^{i+2}) && \text{for } i \geq 0 && (\sigma_i \text{ defined as above),} \\ \bar{\pi}_i &= (s_{0,1}^{i+1} \sigma_0, s_{0,1}^{i+1}) && \text{for } i \geq 0 \end{aligned}$$

4.2. Relations involving cuts and permutations. In the following relations (1)–(7), the reader can assume that $1 \leq d, d' \leq n$ unless otherwise stated.

$$\begin{aligned} (1) \quad & X_{q,d} X_{m,d'} = X_{m,d'} X_{q+1,d} && \text{for } m < q, \\ (2) \quad & \pi_q X_{m,d} = X_{m,d} \pi_{q+1} && \text{for } m < q, \\ (3) \quad & \pi_q X_{q,d} = X_{q+1,d} \pi_q \pi_{q+1} && \text{for } q \geq 0, \\ (4) \quad & \pi_q X_{m,d} = X_{m,d} \pi_q && \text{for } m > q + 1, \\ (5) \quad & \bar{\pi}_q X_{m,d} = X_{m,d} \bar{\pi}_{q+1} && \text{for } m < q, \\ (6) \quad & \bar{\pi}_m X_{m,1} = \pi_m \bar{\pi}_{m+1} && \text{for } m \geq 0, \\ (7) \quad & X_{m,d} X_{m+1,d'} X_{m,d'} = X_{m,d'} X_{m+1,d} X_{m,d} \pi_{m+1} && \text{for } m \geq 0, d \neq d'. \end{aligned}$$

4.3. Relations involving permutations only. We have

$$(8) \quad \pi_q \pi_m = \pi_m \pi_q \quad \text{for } |m - q| > 2,$$

$$(9) \quad \pi_m \pi_{m+1} \pi_m = \pi_{m+1} \pi_m \pi_{m+1} \quad \text{for } m \geq 0,$$

$$(10) \quad \bar{\pi}_q \pi_m = \pi_m \bar{\pi}_q \quad \text{for } q \geq m + 2,$$

$$(11) \quad \pi_m \bar{\pi}_{m+1} \pi_m = \bar{\pi}_{m+1} \pi_m \bar{\pi}_{m+1} \quad \text{for } m \geq 0,$$

$$(12) \quad \pi_m^2 = 1 \quad \text{for } m \geq 0,$$

$$(13) \quad \bar{\pi}_m^2 = 1 \quad \text{for } m \geq 0.$$

4.4. Relations involving baker's maps. In the relations (14)–(18) the reader can assume that $2 \leq d \leq n$ and $1 \leq d' \leq n$ unless otherwise stated.

$$(14) \quad \bar{\pi}_m X_{m,d} = C_{m+1,d} \pi_m \bar{\pi}_{m+1} \quad \text{for } m \geq 0,$$

$$(15) \quad C_{q,d} X_{m,d'} = X_{m,d'} C_{q+1,d} \quad \text{for } m < q,$$

$$(16) \quad C_{m,d} X_{m,1} = X_{m,d} C_{m+2,d} \pi_{m+1} \quad \text{for } m \geq 0,$$

$$(17) \quad \pi_q C_{m,d} = C_{m,d} \pi_q \quad \text{for } m > q + 1,$$

$$(18) \quad C_{m,d} X_{m,d'} C_{m+2,d'} = C_{m,d'} X_{m,d} C_{m+2,d} \pi_{m+1} \quad \text{for } m \geq 0, 1 < d' < d \leq n.$$

Relations (1)–(17) are generalizations of those given in [Brin 2004] and their proofs are completely analogous. The only new family of relations is (18), which we prove using relation (M6) from the monoid:

Proof. We have

$$\begin{aligned} C_{m,d} X_{m,d'} C_{m+2,d'} &= (s_{0,1}^m s_{0,d}, s_{0,1}^{m+1}) (s_{0,1}^{m+1} s_{1,d'}, s_{0,1}^{m+2}) (s_{0,1}^{m+2} s_{0,d'}, s_{0,1}^{m+3}) \\ &= (s_{0,1}^m s_{0,d} s_{1,d'} s_{0,d'}, s_{0,1}^{m+3}) \\ &= (s_{0,1}^m s_{0,d'} s_{1,d} s_{0,d} \sigma_1, s_{0,1}^{m+3}) \\ &= (s_{0,1}^m s_{0,d'}, s_{0,1}^{m+1}) (s_{0,1}^{m+1} s_{1,d}, s_{0,1}^{m+2}) (s_{0,1}^{m+2} s_{0,d}, s_{0,1}^{m+3}) (s_{0,1}^{m+3} \sigma_1, s_{0,1}^{m+3}) \\ &= C_{m,d'} X_{m,d} C_{m+2,d} \pi_{m+1}. \quad \square \end{aligned}$$

Lemma 12 (subscript raising formulas). *We have*

$$C_{r,d} \sim C_{r+1,d} X_{r,d} \pi_{r+1} X_{r,1}^{-1} \quad \text{and} \quad \bar{\pi}_r \sim \pi_r \bar{\pi}_{r+1} X_{r,1}^{-1} \sim X_{r,1} \bar{\pi}_{r+1} \pi_r.$$

The first formula of Lemma 12 follows from relations (15) and (16), while the second is a generalization of the one found in [Brin 2005].

4.5. Secondary relations for nV . These are as follows.

$$\begin{aligned}
X_{q,d}^{-1}X_{r,d} &\sim \begin{cases} X_d X_d^{-1} & \text{if } r \neq q, \\ 1 & \text{if } r = q \end{cases} \quad \text{for } 1 \leq d \leq n, \\
X_{q,d}^{-1}X_{r,d'} &\sim \begin{cases} X_{d'} X_{d'}^{-1} & \text{if } r \neq q, \\ w(X_{d'})\pi w(X_{d'}^{-1}) & \text{if } r = q \end{cases} \quad \text{for } 1 \leq d, d' \leq n, d \neq d', \\
C_{q,d}^{-1}X_{r,d'} &\sim \begin{cases} X_{d'} C_{d'}^{-1} & \text{if } r < q, \\ w(X_1, \pi, X_d^{-1})X_{d'} C_{d'}^{-1} & \text{if } r \geq q \end{cases} \quad \text{for } 2 \leq d \leq n, 1 \leq d' \leq n, \\
X_{r,d'}^{-1}C_{q,d} &\sim \begin{cases} C_d X_{d'}^{-1} & \text{if } r < q, \\ C_d X_{d'}^{-1} w(X_d, \pi, X_1^{-1}) & \text{if } r \geq q \end{cases} \quad \text{for } 2 \leq d \leq n, 1 \leq d' \leq n, \\
\pi_q X_{r,d} &\sim X_d w(\pi) \quad \text{for } 1 \leq d \leq n, \\
\bar{\pi}_q X_{r,1} &\sim \begin{cases} X_1 \bar{\pi} & \text{if } r < q, \\ \pi \bar{\pi} & \text{if } r = q, \\ w(X_1) \bar{\pi} w(\pi) & \text{if } r > q, \end{cases} \\
\bar{\pi}_q X_{r,d} &\sim \begin{cases} X_d \bar{\pi} & \text{if } r < q, \\ C_d \pi \bar{\pi} & \text{if } r = q, \\ w(X_1) X_d \bar{\pi} w(\pi) & \text{if } r > q \end{cases} \quad \text{for } 2 \leq d \leq n, \\
\pi_q C_{r,d} &\sim \begin{cases} C_d \pi & \text{if } r > q + 1, \\ C_d w(X_1^{-1}, \pi, X_d) & \text{if } r \leq q + 1 \end{cases} \quad \text{for } 2 \leq d \leq n, \\
\bar{\pi}_q C_{r,d} &\sim \begin{cases} X_d \bar{\pi} \pi & \text{if } r = q + 1, \\ w(X_1) X_d \bar{\pi} w(\pi) & \text{if } r > q + 1, \\ w(X_d) C_d \pi \bar{\pi} w(\pi, X_1^{-1}) & \text{if } r < q + 1 \end{cases} \quad \text{for } 2 \leq d \leq n, \\
C_{q,d}^{-1}C_{r,d} &\sim \begin{cases} w(X_1^{-1}, \pi, X_d) & \text{if } q < r, \\ 1 & \text{if } q = r, \\ w(X_1, \pi, X_d^{-1}) & \text{if } q > r \end{cases} \quad \text{for } 2 \leq d \leq n, \\
C_{q,d}^{-1}C_{r,d'} &\sim \begin{cases} X_{d'} C_{d'} \pi C_d^{-1} X_d^{-1} w(X_{d'}, \pi, X_1^{-1}) & \text{if } q > r, \\ X_{d'} C_{d'} \pi C_d^{-1} X_d^{-1} & \text{if } q = r, \\ w(X_1, \pi, X_{d'}^{-1}) X_d C_d \pi C_{d'}^{-1} X_{d'}^{-1} & \text{if } q < r \end{cases} \quad \text{for } 1 \leq d' < d \leq n.
\end{aligned}$$

Proof. We only prove the last set of secondary relations as it is the only one that does not immediately descend from the computations in [Brin 2005]. If $q > r$ we can apply the subscript raising formulas repeatedly for j times until $r + j = q$ and

rewrite the product as

$$C_{q,d}^{-1}C_{r,d'} \sim C_{q,d}^{-1}C_{r+1,d'}X_{r,d'}\pi_{r+1}X_{r,1}^{-1} \sim \dots \sim C_{q,d}^{-1}C_{r+j,d'}w(X_{d'}, \pi, X_1^{-1}).$$

We argue similarly if $q < r$. We now have to study the product $C_{q,d}^{-1}C_{q,d'}$. Without loss of generality we assume $d' < d$ and apply relation (18):

$$C_{q,d}^{-1}C_{q,d'} = X_{q,d'}C_{q+2,d'}\pi_{q+1}C_{q+2,d}^{-1}X_{q,d}^{-1},$$

which is what was claimed. Similar relations can be derived if $d' > d$. \square

Remark 13. The last two secondary relations allow us to rewrite a word of type $w(X, C, \pi, C^{-1}, X^{-1})$ in *LMR* form without increasing the number of times C appears, and thereby to generalize the proof of [Brin 2005, Lemma 4.6]; see Lemma 15 below. This observation also lets us generalize [Brin 2005, Lemma 4.7]; see Lemma 16 below. In fact, all our secondary relations are immediate generalizations of those in [Brin 2005]; the last one does not introduce appearances of $\bar{\pi}$ and therefore all the letters in the last secondary relations can be migrated to their needed position by means of the previous secondary relations, without altering the original argument of [Brin 2005, Lemma 4.7]. Therefore even in the case of nV one is able to do the bookkeeping without risk of creating extra letters that cannot be passed safely without recreating them, and hence we obtain an argument that terminates.

5. Presentations for nV

We now show how the relations above enable us to put our group elements into a normal form, starting with words in the generators of nV corresponding to elements from \widehat{nV} .

Lemma 14. *Let w be a word in $\{X_{i,d}, \pi_i, X_{i,d}^{-1} \mid 1 \leq d \leq n, i \in \mathbb{N}\}$. Then $w \sim LMR$, where L and R^{-1} are words in $\{X_{i,d}\}$ and M is a word in $\{\pi_i\}$.*

Proof. There is a homomorphism from \widehat{nV} to nV given by $s_{i,d} \mapsto X_{i,d}$ and $\sigma_i \mapsto \pi_i$. This follows from the correspondence between the relations for \widehat{nV} and nV as given below:

$$\begin{aligned} \text{(M1)} &\rightarrow (1), & \text{(M5a)} &\rightarrow (2), \\ \text{(M2)} &\rightarrow (12), & \text{(M5b), (M5c)} &\rightarrow (3), \\ \text{(M3)} &\rightarrow (8), & \text{(M5d)} &\rightarrow (4), \\ \text{(M4)} &\rightarrow (9), & \text{(M6)} &\rightarrow (7). \end{aligned}$$

Hence, any word w as given above is the image under this homomorphism of a word w' in \widehat{nV} . Since \widehat{nV} is the group of right fractions of the monoid Π_n , we can represent w' as pq^{-1} , where p and q are words in $\{s_{i,d}, \sigma_i \mid 1 \leq d \leq n, i \in \mathbb{N}\}$.

Now, as noted before in the proof of Lemma 6, we can assume p and q are of the form ab , where $a \in \langle s_{i,d} \rangle$ and $b \in \langle \sigma_i \rangle$. Hence, we have written w' as lmr for $l, r^{-1} \in \langle s_{i,d} \rangle$ and $m \in \langle \sigma_i \rangle$ since elements of $\langle \sigma_i \rangle$ are their own inverse. Applying the homomorphism to w' puts w in the desired form. \square

The next two results follow the original proofs of [Brin 2005, Lemmas 4.6 and 4.7] via Remark 13.

Lemma 15. *Let w be of the form $w(X, C, \pi, X^{-1}, C^{-1})$. Then $w \sim LMR$, where L and R^{-1} are words of the form $w(X, C)$ and M is of the form $w(\pi)$. Further the number of appearances of C in L will be no larger than the number of appearances of C in w and the number of appearances of C^{-1} in R will be no larger than the number of appearances of C^{-1} in w .*

Lemma 16. *Let w be a word in the generating set*

$$\{X_{i,d}, C_{i,d'}, \pi_i, \bar{\pi}_i, X_{i,d}^{-1}, C_{i,d'}^{-1} \mid 1 \leq d \leq n, 2 \leq d' \leq n, i \in \mathbb{N}\}.$$

Then $w \sim LMR$, where L and R^{-1} are words of the form $w(X, C)$ and M is of the form $w(\pi, \bar{\pi})$.

Lemma 17. *Let w be a word in the generating set*

$$\{X_{i,d}, C_{i,d'}, \pi_i, \bar{\pi}_i, X_{i,d}^{-1}, C_{i,d'}^{-1} \mid 1 \leq d \leq n, 2 \leq d' \leq n, i \in \mathbb{N}\}.$$

Then $w \sim LMR$, where

- $L = C_{i_0, d_0} C_{i_1, d_1} \dots C_{i_g, d_g} q$ with $i_0 < i_1 < \dots < i_g$ for $g \geq -1$ and q is a word in the set $\{X_{i,d} \mid 1 \leq d \leq n, i \in \mathbb{N}\}$
- $R^{-1} = C_{j_0, d'_0} C_{j_1, d'_1} \dots C_{j_m, d'_m} q'$ with $j_0 < j_1 < \dots < j_m$ for $m \geq -1$ and q' is a word in the set $\{X_{i,d} \mid 1 \leq d \leq n, i \in \mathbb{N}\}$
- M is a word in the set $\{\pi_i, \bar{\pi}_i \mid i \in \mathbb{N}\}$

Proof. By using the secondary relations, we can assume that $w \sim LMR$, where L and R^{-1} are words in $\{X_{i,d}, C_{i,d}\}$ and M is a word in $\{\pi_i, \bar{\pi}_i\}$ by analogous arguments used in [Brin 2005, Lemmas 4.6 and 4.7]. We then improve L using the subscript raising formula for the $C_{i,d}$ and relation (15) as in the proof of [ibid., Lemma 4.8]. To adapt the quoted lemmas from [Brin 2005] we need to use Remark 13 to make sure that appearances of C and $\bar{\pi}$ do not increase. \square

We define the notions of *primary* and *secondary tree* and of *trunk* exactly the same way that Brin does [2005]. The primary tree is the tree corresponding to the word t in Lemma 18 and any extension to the left is a secondary tree for L . The following extends [Brin 2005, Lemma 4.15] adapted to our case. The proof is completely analogous.

Lemma 18. *Let*

$$L = C_{i_0, d_0} C_{i_1, d_1} \cdots C_{i_g, d_g} X_{i_{n+1}, d_{n+1}} \cdots X_{i_{l-1}, d_{l-1}},$$

where $i_0 < i_1 < \cdots < i_g$, where $2 \leq d_k \leq n$ for $k \in \{0, \dots, g\}$ and $1 \leq d_k \leq n$ for $k \in \{g+1, \dots, l-1\}$. Let m equal the maximum of

$$\{i_j + g + 2 - j \mid g+1 \leq j \leq l-1\} \cup \{i_g + 1\}.$$

Then L can be represented as $L = (t, s_{0,1}^k)$, where t is a word in $\{s_{i,d}\}$ and k is the length of t , so that $k = m + l - g$, and so that the tree T for t is the primary tree for L and is described as follows. The tree T consists of a trunk Λ with a finite forest F attached. The trunk Λ has m carets and $m+1$ leaves numbered 0 through m in the right-left order. If the carets in Λ are numbered from 0 starting at the top, then the label of the i -th caret is d_k if $i = i_k$ for k in $\{0, 1, \dots, g\}$ and 1 otherwise.

The following two lemmas are used in proving Remark 13, which allows us to assume the trees corresponding to our group elements are in normal form.

Lemma 19. *Let*

$$L = C_{i_0, d_0} C_{i_1, d_1} \cdots C_{i_g, d_g} u \quad \text{and} \quad L' = C_{k_0, d'_0} C_{k_1, d'_1} \cdots C_{k_g, d'_g} u',$$

where $i_0 < i_1 < \cdots < i_g$, where $k_0 < k_1 < \cdots < k_g$, where u is a word in the set $\{X_{i,d} \mid 1 \leq d \leq n, i \in \mathbb{N}\}$, and where u' is a word in the set $\{X_{i,d}, \pi_i \mid 1 \leq d \leq n, i \in \mathbb{N}\}$. Assume that L is expressible as $(t, s_{0,1}^p)$ as an element of $n\widehat{V}$ with t a word in $\{s_{i,d}\}$ and p the length of t . Let m be the number of carets of the trunk of the tree T corresponding to t and assume that $m \geq k_g + 1$.

If $L \sim L'$, then there is a word u'' in $\{X_{i,d}\}$, and there is a word z in $\{\pi_i \mid i \leq p-2\}$ such that setting $L_1 = C_{k_0, d'_0} C_{k_1, d'_1} \cdots C_{k_g, d'_g} u''$ and $L_2 = L_1 z$ gives that $L \sim L_2$ and L_1 is expressible as $(t', s_{0,1}^p)$ with t' a word in $\{s_{i,d}\}$ of length p , so that the tree T' for t' is normalized except possibly at interior vertices in the trunk of the tree, and so that the trunk of T' has m carets.

Proof. The homomorphism $n\widehat{V} \rightarrow nV$ given by $s_{i,d} \mapsto X_{i,d}$ and $\sigma_i \mapsto \pi_i$ allows us to write $u' \sim u''z'$ with u'' a word in $\{X_{i,d}\}$ and z' a word in $\{\pi_i \mid i \in \mathbb{N}\}$ such that the forest F for u'' is normalized. The rest of the proof goes through as before, but we describe the slight modifications needed for our case. We write $L = (ts_{0,1}^k, s_{0,1}^{p+k}) = (\hat{t}s_{1,0}^r x, s_{1,0}^{q+r}) = L_2$ as elements in $n\widehat{V}$, where x is a word in $\{\sigma_i\}$ and $p+k = q+r$. As before, we can conclude that the unnumbered patterns for $ts_{0,1}^k$ and $\hat{t}s_{1,0}^r$ are identical.

In the tree for $ts_{0,1}^k$, let the left edge vertices be a_0, a_1, \dots, a_b reading from the top, so that a_0 is the root of the tree. Since we assume the trunk of the tree has m carets, we know $b = m+k$ and for $m \leq i < b$, the label for a_i is 1. Similarly, in the tree for $\hat{t}s_{1,0}^r$, let the left edge vertices be a'_0, a'_1, \dots, a'_b reading from the top. Note

that remark (*) in the proof of [Brin 2005, Theorem 4.21] (which we are about to restate) remains true in our general case, by giving a new definition: For each left edge vertex a_i , define the n -tuple (x_1^i, \dots, x_n^i) , where x_k^i equals the number of left edge vertices above a_i with label k . (Note we are using i to denote an index, not an exponent). It follows that $x_1^i + x_2^i + \dots + x_n^i$ is the total number of left edge vertices above a_i . Then we can say,

- (*) The rectangle corresponding to a left edge vertex a_i depends only on the n -tuple (x_1^i, \dots, x_n^i) .

In other words, for the rectangle labeled 0 in any pattern, the order of the different cuts does not matter. This is because the rectangle labeled 0 must contain the origin and its size in each dimension k will be $2^{-x_k^i}$. Hence, the analogous statement for our case follows, and we conclude that the n -rectangle R corresponding to a_m is identical to the n -rectangle R' corresponding to a'_m . Since R is divided k times across dimension 1, so is R' , and hence the tree below a'_m must consist of an extension to the left by k carets all labeled 1, and we can conclude that $r \geq k$. The rest of the proof follows exactly as before. \square

Here, we define a notion of *complexity* to measure progress in the following lemma and proposition towards normalizing trees. If T is a labeled tree, we let a_0, a_1, \dots, a_m be the interior, left edge vertices of T reading from top to bottom so that a_0 is the root. Let $b_0 b_1 \dots b_m$ be a word in $\{1, 2, \dots, n\}$ where $b_i = k$ if a_i is labeled k for $0 \leq i \leq m$. We say $b_0 b_1 \dots b_m$ is the complexity of T . We impose the length-lex ordering on such words, that is, if w_1 and w_2 are two such words, then we say $w_1 < w_2$ if w_1 is shorter than w_2 or if $w_1 = b_0^1 \dots b_m^1$ and $w_2 = b_0^2 \dots b_m^2$ are two such words of the same length, then $w_1 < w_2$ if when we take $j \in \{0, \dots, m\}$ minimal where $b_j^1 \neq b_j^2$, we have $b_j^1 < b_j^2$.

Lemma 20. *Let $L = C_{i_0, d_0} C_{i_1, d_1} \dots C_{i_g, d_g} u$, where $i_0 < i_1 < \dots < i_g$ and u is a word in the set $\{X_{i, d}\}$. Assume that the primary tree T for L is normalized except at one or more vertices in the trunk of T . Let m be the number of carets in the trunk of T . Then $L \sim L' = C_{k_0, c_0} C_{k_1, c_1} \dots C_{k_g, c_g} u'$, where $k_0 < k_1 < \dots < k_g$ and u' is a word in the set $\{X_{i, d}, \pi_s\}$, so that $m \geq k_g + 1$, and so that the complexity of the primary tree T' of L' is strictly less than the complexity of T .*

Proof. We want to use the relations to push a suitable instance of an $X_{u, v}$ in the word L as far as possible to the left to be able to apply a cross relation. This operation normalizes a suitable vertex and decreases the complexity of the primary tree T .

Let Λ be the trunk of T . The interior vertices of Λ are the interior, left edge vertices of T and let these be a_0, a_1, \dots, a_{m-1} . Let r be the highest value with $0 \leq r < m$ for which a_r is not normalized. This is the lowest nonnormalized

interior vertex of Λ , and since a_r is not normalized it is labeled $\ell \neq 1$ and must correspond to some $C_{i_j, \ell}$. From Lemma 18, we have $i_j = r$.

Since it is not normalized, a_r must correspond to some hypercube S_{i_j} that is fully divided across dimension ℓ and some other dimension d , with $1 \leq d < \ell$.

By rewriting L as $(t, s_{0,1}^k)$ (which we can do by Lemma 18) and applying Corollary 9 to t , we can assume that the children of a_r , v_1 and v_2 , are both labeled d . We divide our work in two cases, $d = 1$ and $d > 1$. We observe that the case $d = 1$ is entirely analogous to the proof of [Brin 2005, Theorem 4.22] while the case $d > 1$ is slightly different.

Case 1: $d = 1$. In this case, the left child v_1 , which is in the trunk Λ , is labeled 1. In the case that $j < n$ we observe that $i_{j+1} > r + 1 = i_j + 1$, since the interior vertex of the trunk corresponding to $C_{i_{j+1}, d_{j+1}}$ is not labeled 1 (otherwise, $a_r = a_{i_j}$ would not be the lowest nonnormalized interior vertex). Since the right child v_2 is an interior vertex not on the trunk, there must be a letter $X_{q,1}$ corresponding to it. By Lemma 5 we can assume that $X_{q,1}$ occurs as the first letter of u , that is, $u = X_{q,1}u''$. Hence

$$L = C_{i_0} \cdots C_{i_{j-1}} \underline{C_{i_j, \ell}} C_{i_{j+1}} \cdots C_{i_g} \underline{X_{q,1}} u'',$$

where we have omitted all the dimension subscripts of the baker's maps $C_{i,d}$ (except for one map) since they are not important for the argument. The subword $C_{i_0} \cdots C_{i_j, \ell} \cdots C_{i_g} X_{q,1}$ is a trunk with a single caret labeled 1 attached at the caret i_j of the trunk on its right child. By a careful observation of the right-left ordering it is evident that $q = i_j$. By using relation (15) repeatedly on L we can move $X_{q,1} = X_{i_j,1}$ to the left and rewrite the word L as

$$C_{i_0} \cdots C_{i_{j-1}} \underline{C_{i_j, \ell} X_{i_j,1}} C_{i_{j+1}+1} \cdots C_{i_g+1} u'',$$

since $i_0 < i_1 < \cdots < i_g$ and $i_{j+1} > i_j + 1$. Combining relations (15) and (16) on the product $C_{i_j, \ell} X_{i_j,1}$, we rewrite L as

$$C_{i_0} \cdots C_{i_{j-1}} \underline{C_{i_j+1, \ell} X_{i_j, \ell} \pi_{i_j+1}} C_{i_{j+1}+1} \cdots C_{i_g+1} u''.$$

Now we apply (17) to commute π_{i_j+1} back to the right without affecting the indices of the baker's maps. This is possible since $i_{j+1} > i_j + 1$ and therefore $i_{j+1} + 1 > i_j + 2$. Now we apply (15) repeatedly to the word

$$C_{i_0} \cdots C_{i_{j-1}} \underline{C_{i_j+1, \ell} X_{i_j, \ell}} C_{i_{j+1}+1} \cdots C_{i_g+1} \underline{\pi_{i_j+1}} u''$$

to bring $X_{i_j, \ell}$ back to the right, decreasing the indices of the baker's maps by 1

$$C_{i_0} \cdots C_{i_{j-1}} \underline{C_{i_j+1, \ell}} C_{i_{j+1}} \cdots C_{i_g} \underline{X_{i_j, \ell} \pi_{i_j+1}} u''.$$

By setting $u' = X_{i_j, \ell} \pi_{i_j+1} u''$ in the previous equation and relabeling the indices with the k_i , we obtain the word $L' = C_{k_0, c_0} C_{k_1, c_1} \cdots C_{k_g, c_g} u'$ whose primary tree T' is the same as T up until the vertex a_r , which is now labeled $d = 1$ instead of ℓ . Thus, $L \sim L' = C_{k_0, c_0} C_{k_1, c_1} \cdots C_{k_g, c_g} u'$ and the complexity of the primary tree T' of L' is strictly less than the complexity of T .

The only thing we still need to prove in this case is that $m \geq k_g + 1$. However, it has been observed above that $i_j = r < m - 1$ so $i_j + 2 \leq m$. This gives the result in the case that $j = n$. If $j < n$, then $k_g = i_g$ and $m \geq i_g + 1$ by Lemma 18.

Case 2: $1 < d < \ell$. We observe that a_r corresponds to $C_{i_j, \ell}$ and that v_1 corresponds to $C_{i_k, d}$. By Lemma 18, we have $r + 1 = i_k$, which implies $i_k = i_j + 1 = i_{j+1}$. In fact, if $i_j + 1 < i_{j+1}$, there would be a vertex labeled 1 on the trunk between the vertices i_j and i_{j+1} (and this is impossible since $d > 1$). Let $X_{i_j, d}$ correspond to the right child v_2 . Arguing as in the case $d = 1$ we have

$$L = C_{i_0} \cdots C_{i_{j-1}} \underline{C_{i_j, \ell} C_{i_j+1, d} C_{i_{j+2}}} \cdots C_{i_g} X_{q, d} u''.$$

We apply relation (15) as before to move $X_{q, d} = X_{i_j, d}$ to the left while increasing the subscript of each baker's map by 1:

$$C_{i_0} \cdots C_{i_{j-1}} \underline{C_{i_j, \ell} X_{i_j, d} C_{i_j+2, d} C_{i_{j+2}+1}} \cdots C_{i_{g+1}} u''.$$

By using the cross relation (18) on the underlined portion, we read it as

$$C_{i_0} \cdots C_{i_{j-1}} \underline{C_{i_j, d} X_{i_j, \ell} C_{i_j+2, \ell} \pi_{i_j+1} C_{i_{j+2}+1}} \cdots C_{i_{g+1}} u''.$$

Since $i_{j+2} > i_{j+1}$, then $i_{j+2} + 1 > i_{j+1} + 1$; hence π_{i_j+1} and the baker's maps to its right commute, so the word becomes

$$C_{i_0} \cdots \underline{C_{i_j, d} X_{i_j, \ell} C_{i_j+2, \ell} C_{i_{j+2}+1}} \cdots C_{i_{g+1}} \underline{\pi_{i_j+1}} u''.$$

We apply (15) repeatedly and move $X_{i_j, \ell}$ back to the right to obtain

$$L \sim C_{i_0} \cdots \underline{C_{i_j, d} C_{i_j+1, \ell} C_{i_{j+2}}} \cdots C_{i_g} X_{i_j, \ell} \pi_{i_j+1} u'',$$

where the product $C_{i_j, d} C_{i_j+2, \ell}$ has been underlined to stress that the new trunk has the vertices labeled d and ℓ , which are now switched. Thus the complexity of the tree has been lowered. In this second case, the new sequence $k_0 < \cdots < k_g$ is exactly equal to the initial one $i_0 < \cdots < i_g$. By the definition of m (given in Lemma 18) applied on the initial word L , we have $m \geq i_g + 1$ and so, since $k_g = i_g$, we are done. \square

Remark 21. As observed in the proof above, the case $d = 1$ is equivalent to [Brin 2005, Theorem 4.22], though the proof therein leads to a condition that is equivalent to lowering the complexity. When the index in some $C_{i_j, d}$ goes up by 1, this

corresponds to switching the vertices with labels d and 1 in the primary tree and thus lowering the complexity by making more vertices normalized.

Proposition 22. *Let w be a word in the generating set*

$$\{X_{i,d}, C_{i,d'}, \pi_i, \bar{\pi}_i, X_{i,d}^{-1}, C_{i,d'}^{-1} \mid 1 \leq d \leq n, 2 \leq d' \leq n, i \in \mathbb{N}\}.$$

Then $w \sim LMR$ as in Lemma 17 and when expressed as elements of \widehat{nV} we have

$$L = ts_{0,1}^{-p}, \quad R^{-1} = ys_{0,1}^{-p}, \quad M = s_{0,1}^p us_{0,1}^{-p},$$

where t, y are words in $\{s_{i,d} \mid 1 \leq d \leq n, i \in \mathbb{N}\}$, u is a word in $\{\sigma_j \mid 0 \leq j \leq p-1\}$, and the lengths of t and y are both p . Further, we may assume the trees for t and y are normalized, and if u can be reduced to the trivial word using relations (2)–(4), then M can be reduced to the trivial word using relations (13)–(17).

Proof. The proof of the first conclusion is exactly the same as that of [Brin 2010, Lemma 4.19]. In order to assume the trees for t and y are normalized, we alternate applying Lemmas 19 and 20. We have L expressed as $(t, s_{0,1}^p)$, where p is the length of t and the number of carets in the trunk of the tree T for t is m . Setting $L = L'$ certainly gives that $L \sim L'$ and $m \geq k_g + 1$ by Lemma 18, so we have satisfied the hypotheses of Lemma 19. Therefore, $L \sim L_1 z$ where L_1 expressed as $(t', s_{0,1}^p)$, where the trunk of the tree T' for t' has m carets. Since we set $L = L'$, we see that the trunks of T and T' are identical and the only way in which the two trees differ is that T' is normalized off the trunk. Since z is a word in $\{\pi_i\}$, z can be absorbed into M without disrupting the assumptions on M , namely, M can still be written in the form $M = s_{0,1}^p us_{0,1}^{-p}$ as above. We now replace L with L_1 and proceed to use Lemma 20.

Since the tree for L is now normalized off the trunk, we satisfy the hypotheses of Lemma 20 and write $L \sim L'$, where the tree for L' has complexity lower than the tree for L and $m \geq k_g + 1$. Hence, we can now apply Lemma 19 again and obtain $L \sim L_1 z$ and let z be absorbed into M . We apply this process over and over, decreasing the complexity of the tree associated to L each time. Since there are only finitely many linearly ordered complexities, eventually this process will terminate, at which point the tree for L will be normalized. We can apply the same procedure to the inverse of LMR to normalize the tree for R . The last statement regarding M follows immediately from [Brin 2005, Lemma 4.18]. \square

Theorem 23. *Let w be a word in the generating set*

$$\{X_{i,d}, C_{i,d'}, \pi_i, \bar{\pi}_i, X_{i,d}^{-1}, C_{i,d'}^{-1} \mid 1 \leq d \leq n, 2 \leq d' \leq n, i \in \mathbb{N}\}$$

that represents the trivial element of nV . Then $w \sim 1$ using the relations in (1)–(18). Hence, we have a presentation for nV .

Proof. Using Proposition 22, we can assume

$$w \sim LMR = (ts_{0,1}^{-p})(s_{0,1}^p u s_{0,1}^{-p})(s_{0,1}^p y^{-1}), = tuy^{-1}$$

where t and y are words in $\{s_{i,d} \mid 1 \leq d \leq n, i \in \mathbb{N}\}$, u is a word in $\{\sigma_j \mid 0 \leq j \leq p-1\}$, and the trees associated to t and y are normalized. By assumption, $tuy^{-1} = (tu, y)$ is the trivial element of $n\widehat{V}$ and so tu and y represent the same numbered patterns in Π_n . Furthermore, t and y must give the same unnumbered pattern, while u enacts a permutation on the numbering. Since the forests for t and y are normalized and give the same pattern, the forests are identical with the same labeling by Lemma 7. The numbering on the leaves for both forests follows the left-right ordering; hence t and y give the same numbered patterns, which implies that u enacts the trivial permutation and $M \sim 1$ by Proposition 22.

We now wish to show that $L \sim R^{-1}$. By Lemma 17, we have

$$L = C_{i_0, d_0} C_{i_1, d_1} \cdots C_{i_g, d_g} q \quad \text{and} \quad R^{-1} = C_{j_0, d'_0} C_{j_1, d'_1} \cdots C_{j_m, d'_m} q'$$

Since we know that the trunks of the trees corresponding to L and R^{-1} are identical with the same labeling, the sequences (i_0, i_1, \dots, i_g) and (j_0, j_1, \dots, j_m) are identical and $d_k = d'_k$ for each $k \in \{0, 1, \dots, n = m\}$. Hence, the subwords $C_{i_0, d_0} C_{i_1, d_1} \cdots C_{i_g, d_g}$ and $C_{j_0, d'_0} C_{j_1, d'_1} \cdots C_{j_m, d'_m}$ are the same and it remains to show that $q \sim q'$. This follows from Lemma 4 and the homomorphism from $n\widehat{V}$ to nV as before. \square

6. Finite presentations

6.1. Finite presentation for $n\widehat{V}$. We now give a finite presentation for $n\widehat{V}$, using arguments analogous to those found in [Brin 2005] to show that the full set of relations is the result of only finitely many of them.

Theorem 24. *The group $n\widehat{V}$ is presented by the $2n + 2$ generators $\{s_{i,d}, \sigma_i \mid i \in \{0, 1\}, 1 \leq d \leq n\}$ and the $5n^2 + 7n + 6$ relations given below:*

- (M1) $s_{1,1}^{-1} s_{1+k, d'} s_{1,1} = s_{2+k, d'}$ for $k = 1, 2,$
 $s_{i,d}^{-1} s_{i+k, d'} s_{i,d} = s_{i+k+1, d'}$ for $i = 0, 1, k = 1, 2, 2 \leq d \leq n,$
- (M2) $\sigma_i^2 = 1$ for $i = 0, 1,$
- (M3) $\sigma_i \sigma_{i+k} = \sigma_{i+k} \sigma_i$ for $i = 0, 1, k = 2, 3,$
- (M4) $\sigma_i \sigma_{i+1} \sigma_i = \sigma_{i+1} \sigma_i \sigma_{i+1}$ for $i = 0, 1,$
- (M5a) $\sigma_{k+1} s_{1,1} = s_{1,1} \sigma_{k+2}$ for $k = 1, 2,$
 $\sigma_{i+k} s_{i,d} = s_{i,d} \sigma_{i+k+1}$ for $i = 0, 1, k = 1, 2, 2 \leq d \leq n,$
- (M5b)/(M5c) $\sigma_i s_{i,d} = s_{i+1, d} \sigma_i \sigma_{i+1}$ for $i = 0, 1,$

$$(M5d) \quad \sigma_i s_{i+k,d} = s_{i+k,d} \sigma_i \quad \text{for } i = 0, 1, k = 2, 3,$$

$$(M6) \quad s_{i,d} s_{i+1,d'} s_{i,d'} = s_{i,d'} s_{i+1,d} s_{i,d} \sigma_{i+1} \quad \text{for } i = 0, 1, d \neq d'.$$

Proof. First, recall our generating set is $\{s_{i,d}, \sigma_i \mid i \in \mathbb{N}, 1 \leq d \leq n\}$. When $i < j$, relations (M1) and (M5a) give $s_{i,1}^{-1} x_j s_{i,1} = x_{j+1}$, where $x_j = s_{j,d}$ (for some d) or σ_j . Hence, we can use

$$s_{i,d} = s_{0,1}^{1-i} s_{1,d} s_{0,1}^{i-1} \quad \text{and} \quad \sigma_i = s_{0,1}^{1-i} \sigma_1 s_{0,1}^{i-1}$$

as definitions for $i \geq 2$. Therefore, \widehat{nV} is generated by

$$\{s_{i,d}, \sigma_i \mid i \in \{0, 1\}, 1 \leq d \leq n\},$$

which gives a generating set of size $2n + 2$ for each n .

We treat relations (M1)–(M6) as they are treated in [Brin 2005]. Relations involving only one parameter, such as (M2), (M4), and (M6), are obtained for $i \geq 2$ by setting $i = 1$ and conjugating by powers of $s_{0,1}$; therefore the only necessary relations to include are those having $i = 0$ and $i = 1$. As before, (M2) and (M4) follow from $\sigma_0^2 = 1$, $\sigma_1^2 = 1$, $\sigma_0 \sigma_1 \sigma_0 = \sigma_1 \sigma_0 \sigma_1$, and $\sigma_1 \sigma_2 \sigma_1 = \sigma_2 \sigma_1 \sigma_2$, or 4 relations for each n . Relation (M6) follows from 2 relations for each pair of distinct dimensions, giving $2 \binom{n}{2} = n(n-1)$ relations for each n .

Relation (M3) is treated the same way as in [Brin 2005] for each n . Hence, for all i and j , (M3) follows from the 4 relations $\sigma_0 \sigma_2 = \sigma_2 \sigma_0$, $\sigma_0 \sigma_3 = \sigma_3 \sigma_0$, $\sigma_1 \sigma_3 = \sigma_3 \sigma_1$, $\sigma_1 \sigma_4 = \sigma_4 \sigma_1$.

For relation (M1), which can be rewritten as $s_{i,d}^{-1} s_{i+k,d'} s_{i,d} = s_{i+k+1,d'}$ for $k > 0$, we have two cases: the case where $d = 1$ and the case where $d \neq 1$. If $d = 1$, then the case $i = 0$ follows by definition, and by the same induction argument used in [Brin 2005] implies that the relation for all i and k follows from the cases where $i = 1$ and $k = 1, 2$; hence we need only 2 relations per dimension. If $d \neq 1$, we do not get the case $i = 0$ by definition and we must include $i = 0, 1$ and $k = 1, 2$, that is, 4 relations per each pair of dimensions. There are $n - 1$ choices for d , as $d \neq 1$, and n choices for d' , so this case yields $4n(n-1)$ relations. Hence, in total (M1) can be obtained for all i and k by $2n + 4n(n-1) = 4n^2 - 2n$ relations.

For relation (M5b), $\sigma_i s_{i,d} = s_{i+1,d} \sigma_i \sigma_{i+1}$, there is only a single parameter to deal with; hence the relation for $i \geq 2$ can be obtained from the cases where $i = 0, 1$ by conjugating by $s_{0,1}$ as before. Relation (M5c) is actually equivalent to (M5b); hence for each n we only need $2n$ relations for (M5b) and (M5c). We treat (M5a) $\sigma_{i+k} s_{i,d} = s_{i,d} \sigma_{i+k+1}$ for $k > 0$ the same way as for (M1), hence 2 relations are required for $d = 1$ and 4 for $d \neq 1$ for a total of $4n - 2$ relations. And lastly, (M5d) $\sigma_i s_{i+k,d} = s_{i+k,d} \sigma_i$ can be obtained in the same way as the second case of (M1) where the relation for all i, k is obtained by $i = 0, 1, k = 2, 3$, that is, $4n$ relations. \square

6.2. Finite presentation for nV .

Theorem 25. *The group nV is presented by the $2n + 4$ generators*

$$\{X_{i,d}, \pi_i, \bar{\pi}_i \mid i \in \{0, 1\}, 1 \leq d \leq n\},$$

the $5n^2 + 7n + 6$ relations obtained from the homomorphism $\widehat{nV} \rightarrow nV$, and the additional $5n^2 + 3n + 4$ relations given below, for a total of $10n^2 + 10n + 10$ relations.

- (5) $\bar{\pi}_{k+1}X_{1,1} = X_{1,1}\bar{\pi}_{k+2}$ *for $k = 1, 2$,*
 $\bar{\pi}_{m+k}X_{m,d} = X_{m,d}\bar{\pi}_{m+k+1}$ *for $m = 0, 1, k = 1, 2,$*
 $2 \leq d \leq n,$
- (10) $\bar{\pi}_{m+k}\pi_m = \pi_m\bar{\pi}_{m+k}$ *for $m = 0, 1, k = 2, 3,$*
- (11) $\pi_m\bar{\pi}_{m+1}\pi_m = \bar{\pi}_{m+1}\pi_m\bar{\pi}_{m+1}$ *for $m = 0, 1$*
- (13) $\bar{\pi}_m^2 = 1$ *for $m = 0, 1,$*
- (6) $\bar{\pi}_mX_{m,1} = \pi_m\bar{\pi}_{m+1}$ *for $m = 0, 1,$*
- (14) $\bar{\pi}_mX_{m,d} = C_{m+1,d}\pi_m\bar{\pi}_{m+1}$ *for $m = 0, 1, d \neq 1,$*
- (15) $C_{k+1,d}X_{1,1} = X_{1,1}C_{k+2,d}$ *for $k = 1, 2,$*
 $C_{m+k,d}X_{m,d'} = X_{m,d'}C_{m+k+1,d}$ *for $m = 0, 1, k = 1, 2,$*
 $2 \leq d, d' \leq n,$
- (16) $C_{m,d}X_{m,1} = X_{m,d}C_{m+2,d}\pi_{m+1}$ *for $m = 0, 1, 2 \leq d \leq n,$*
- (17) $\pi_mC_{m+k,d} = C_{m+k,d}\pi_m$ *for $m = 0, 1, k = 2, 3,$*
- (18) $C_{m,d}X_{m,d'}C_{m+2,d'} = C_{m,d'}X_{m,d}C_{m+2,d}\pi_{m+1}$ *for $m = 0, 1,$*
 $1 < d' < d \leq n,$

Proof. We can use the relations in nV to write, for $i \geq 2$ and $1 \leq d \leq n$,

$$X_{i,d} = X_{0,1}^{1-i}X_{1,d}X_{0,1}^{i-1}, \quad \pi_i = X_{0,1}^{1-i}\pi_1X_{0,1}^{i-1}, \quad \bar{\pi}_i = X_{0,1}^{1-i}\bar{\pi}_1X_{0,1}^{i-1}.$$

We can also use the relations for nV as in [Brin 2004, Proposition 6.2] to write

$$C_{m,d} = (\bar{\pi}_mX_{m,d}\bar{\pi}_{m+1}\pi_m)(X_{m,d}\pi_{m+1}X_{m,1}^{-1})$$

for $m \geq 0$ and $2 \leq d \leq n$, which we use as a definition. Hence, the $C_{m,d}$ are not needed to generate nV .

The homomorphism $\widehat{nV} \rightarrow nV$ given by $s_{i,d} \mapsto X_{i,d}$ and $\sigma_i \mapsto \pi_i$ implies that the work done for the relations for \widehat{nV} carries over to relations (1)–(4), (7)–(9), and (12) (see Lemma 14). Relations (10), (11), (13) and (6) are exactly the same as those from $2V$ and can be treated as in [Brin 2005], contributing a total of 10 relations to our finite set.

Relation (5) can be treated in a manner similar to (M1) from \widehat{nV} , where 2 relations are needed for dimension 1 and 4 for all others, contributing a total of $4(n-1) + 2$ relations. Relations (14) and (16) include only one parameter and hence can be obtained from the cases where $i = 0, 1$ as before, contributing $2(n-1)$ relations apiece. And (17) requires 4 relations for each $d \neq 1$, hence adding an additional $4(n-1)$ relations.

For relation (15), we have two cases: For $d' = 1$, all cases follow from when $i = 0, 1$, giving us $2(n-1)$ relations since $2 \leq d \leq n$. For $d' \neq 1$, four relations are required for each pair $d, d' \in \{2, \dots, n\}$, contributing $4(n-1)(n-1)$ relations. Lastly, since (18) involves only one parameter in the first component, we only need 2 relations for each $1 < d' < d \leq n$, the number of pairs being $(n-1)(n-2)/2$. \square

Remark 26. Since ωV is an ascending union of the nV , a word

$$w \in \{X_{i,d}, \pi_i, \bar{\pi}_i \mid i \in \{0, 1\}, d \in \mathbb{N}\}$$

such that $w =_{\omega V} 1$ must be contained in some nV (for some $n \in \mathbb{N}$) and so we can use the same ideas and the relations inside nV to transform w into the empty word. Therefore, the following result is an immediate consequence of Theorem 25.

Corollary 27. *The group ωV is generated by the set $\{X_{i,d}, \pi_i, \bar{\pi}_i \mid i \in \{0, 1\}, d \in \mathbb{N}\}$ and satisfies the family of relations in Theorem 25 with the only exception that the parameters $d, d' \in \mathbb{N}$.*

7. Simplicity of nV and ωV

Brin [2010] proved that the groups nV and ωV are simple by showing that the baker's map is a product of transpositions and following the outline of an existing proof that V is simple.

We prove again Brin's simplicity result verify that Brin's original proof that $2V$ is simple [2004, Theorem 7.2] generalizes using the generators and the relations that have been found.

Theorem 28. *The groups nV equal their commutator subgroups for $n \leq \omega$.*

Proof. The goal is to show that the generators $X_{m,i}$, π_m and $\bar{\pi}_m$ are products of commutators. We write $f \simeq g$ to mean that $f = g$ modulo the commutator subgroup. The arguments below are independent of the dimension i .

From relation (1) we see that $X_{q,i}^{-1} X_{0,1}^{-1} X_{q,i} X_{0,1} = X_{q,i}^{-1} X_{q+1,i}$ for $q \geq 1$ and so $X_{q+1,i} \simeq X_{q,i}$. Therefore $X_{q,i} \simeq X_{1,i}$, for $q \geq 1$. Using relation (2) and arguing similarly, we see that $\pi_q \simeq \pi_1$ for $q \geq 1$.

From relation (3) we see that $\pi_0 X_{0,i} \pi_0^{-1} X_{0,i}^{-1} = X_{1,i} \pi_1 X_{0,i}^{-1}$ so that $X_{0,i} \simeq X_{1,i} \pi_1$. Also, by relation (3), $X_{2,i} \simeq X_{1,i}$, and the fact that $\pi_2 \simeq \pi_1$, we see $\pi_1 X_{1,i} = X_{2,i} \pi_1 \pi_2 \simeq X_{1,i} \pi_1 \pi_1 = X_{1,i}$. Therefore $\pi_1 \simeq 1$ and so $X_{0,i} \simeq X_{1,i}$.

Relation (9) and $\pi_1 \simeq 1$ give $\pi_0^2 \simeq \pi_0\pi_1\pi_0 = \pi_1\pi_0\pi_1 \simeq \pi_0$, which implies $\pi_0 \simeq 1$.

By relation (6) and the fact that $\pi_1 \simeq 1$ and $\bar{\pi}_1 \simeq \bar{\pi}_0$, we get $\bar{\pi}_1 X_{1,1} = \pi_1 \bar{\pi}_2 \simeq \bar{\pi}_1$. Hence $X_{0,1} \simeq X_{1,1} \simeq 1$.

Now, relation (6) and $X_{0,1} \simeq 1$ give that $\bar{\pi}_0 \simeq \bar{\pi}_0 X_{0,1} = \bar{\pi}_1$. Relation (11) and $\pi_0 \simeq 1$ lead to $\bar{\pi}_1 \simeq \pi_0 \bar{\pi}_1 \pi_0 = \bar{\pi}_1 \pi_0 \bar{\pi}_1 \simeq \bar{\pi}_1^2$. Therefore $\bar{\pi}_0 \simeq \bar{\pi}_1 \simeq 1$.

Finally, by relation (7) and $X_{0,1} \simeq X_{1,1} \simeq 1 \simeq \pi_1$ we get

$$X_{1,i} X_{0,i} \simeq X_{0,1} X_{1,i} X_{0,i} = X_{0,i} X_{1,1} X_{0,1} \pi_1 \simeq X_{0,i},$$

which implies $X_{0,i} \simeq X_{1,i} \simeq 1$. We have thus proved that all the generators of nV are in the commutator subgroup. The case of ωV is identical: Each generator lies in some nV and can be written as a product of commutators within that subgroup. \square

From [Brin 2004, Section 3.1] (which generalizes to nV and ωV as observed by Brin [2005; 2010]) the commutator subgroup of nV and ωV are simple; therefore Theorem 28 implies the following result.

Theorem 29. *The groups nV are simple for $n \leq \omega$.*

8. An alternative generating set

For any $n \in \mathbb{N}$, we have $(n-1)V \times V \leq nV$. It can be shown that another generating set for nV is given by taking a generating set for $(n-1)V \times V$ and adding an involution that swaps two disjoint subcubes of $[0, 1]^n$, one of which has the origin as one of its vertices and the other of which contains the vertex $(1, \dots, 1)$. This second generating set has the advantage of taking the generators of $(n-1)V$ and adding only the generators of V plus another one. This leads to a smaller generating set, which was suggested to us by Collin Bleak. It seems feasible that a good set of relations exist for this alternative generating set.

Acknowledgments

We thank Robert Strichartz and the National Science Foundation for their support during the REU. We thank Collin Bleak and Martin Kassabov for several helpful conversations and Matt Brin for helpful comments and for pointing out that his argument for the simplicity of $2V$ lifts immediately to nV using the presentations that we find. We thank Matt Brin, Collin Bleak, Dessislava Kochloukova, Daniel Lanoue, Conchita Martinez-Perez and Brita Nucinkis for kindly citing this work while it was still in preparation. We thank Roman Kogan for advice on how to create diagrams using Inkscape.

References

[Bleak and Lanoue 2010] C. Bleak and D. Lanoue, “A family of non-isomorphism results”, *Geom. Dedicata* **146** (2010), 21–26. MR 2011d:20054 Zbl 1213.20029

- [Brin 2004] M. G. Brin, “Higher dimensional Thompson groups”, *Geom. Dedicata* **108** (2004), 163–192. MR 2005m:20008 Zbl 1136.20025
- [Brin 2005] M. G. Brin, “Presentations of higher dimensional Thompson groups”, *J. Algebra* **284**:2 (2005), 520–558. MR 2007e:20062 Zbl 1135.20022
- [Brin 2010] M. G. Brin, “On the baker’s map and the simplicity of the higher dimensional Thompson groups nV ”, *Publ. Mat.* **54**:2 (2010), 433–439. MR 2011g:20038 Zbl 05770007
- [Cannon et al. 1996] J. W. Cannon, W. J. Floyd, and W. R. Parry, “Introductory notes on Richard Thompson’s groups”, *Enseign. Math.* (2) **42**:3-4 (1996), 215–256. MR 98g:20058 Zbl 0880.20027
- [Guralnick et al. 2011] R. M. Guralnick, W. M. Kantor, M. Kassabov, and A. Lubotzky, “Presentations of finite simple groups: a computational approach”, *J. Eur. Math. Soc.* **13**:2 (2011), 391–458. MR 2011m:20035 Zbl 05842815
- [Kochloukova et al. 2010] D. H. Kochloukova, C. Martinez-Perez, and B. E. A. Nucinkis, “Cohomological finiteness properties of the Brin–Thompson–Higman groups $2V$ and $3V$ ”, preprint, 2010. arXiv 1009.4600

Received May 18, 2011. Revised February 13, 2012.

JOHANNA HENNIG
DEPARTMENT OF MATHEMATICS
UNIVERSITY OF CALIFORNIA, SAN DIEGO
9500 GILMAN DRIVE
LA JOLLA, CA 92093
UNITED STATES
jhennig@math.ucsd.edu

FRANCESCO MATUCCI
DEPARTMENT OF MATHEMATICS
UNIVERSITY OF VIRGINIA
325 KERCHOF HALL
CHARLOTTESVILLE, VA 22904
UNITED STATES
fm6w@virginia.edu

RESONANT SOLUTIONS AND TURNING POINTS IN AN ELLIPTIC PROBLEM WITH OSCILLATORY BOUNDARY CONDITIONS

ALFONSO CASTRO AND ROSA PARDO

We consider the elliptic equation $-\Delta u + u = 0$ with nonlinear boundary conditions $\partial u / \partial n = \lambda u + g(\lambda, x, u)$, where the nonlinear term g is oscillatory and satisfies $g(\lambda, x, s)/s \rightarrow 0$ as $|s| \rightarrow 0$. We provide sufficient conditions on g for the existence of sequences of resonant solutions and turning points accumulating to zero.

1. Introduction

This work complements the study initiated in [Arrieta et al. 2010] and [Castro and Pardo 2011] on the positive solutions to the following boundary-value problem

$$(1-1) \quad \begin{cases} -\Delta u + u = 0 & \text{in } \Omega, \\ \frac{\partial u}{\partial n} = \lambda u + g(\lambda, x, u) & \text{on } \partial\Omega, \end{cases}$$

where $\Omega \subset \mathbb{R}^N$ is a bounded and sufficiently smooth domain, $N \geq 2$, λ is a real parameter, $g(\lambda, x, s) = o(s)$ as $s \rightarrow 0$ and g is oscillatory. A typical example of such a g is

$$(1-2) \quad g(x, s) := s^\alpha \left(\sin \left| \frac{s}{\Phi_1(x)} \right|^\beta + C \right) \quad \text{with } \alpha + \beta > 1, \quad \beta < 0,$$

where Φ_1 stands for the first eigenfunction of the Steklov eigenvalue problem

$$(1-3) \quad \begin{cases} -\Delta \Phi + \Phi = 0 & \text{in } \Omega, \\ \frac{\partial \Phi}{\partial n} = \sigma \Phi & \text{on } \partial\Omega. \end{cases}$$

Pardo is supported by the Spanish Ministerio de Ciencia e Innovación (MICINN) under Project MTM2009-07540, by UCM-BSCH, Spain, GR58/08, Grupo 920894 and also by the Programa Becas Complutense del Amo. This work was carried out during Pardo's sabbatical visit to the Department of Mathematics, Harvey Mudd College, whose hospitality she thanks.

MSC2010: primary 35B32, 35B34, 35B35, 35J25, 58J55; secondary 35J60, 35J65.

Keywords: turning points, resonance, stability, instability, multiplicity, Steklov eigenvalues, bifurcation, sublinear oscillating boundary conditions.

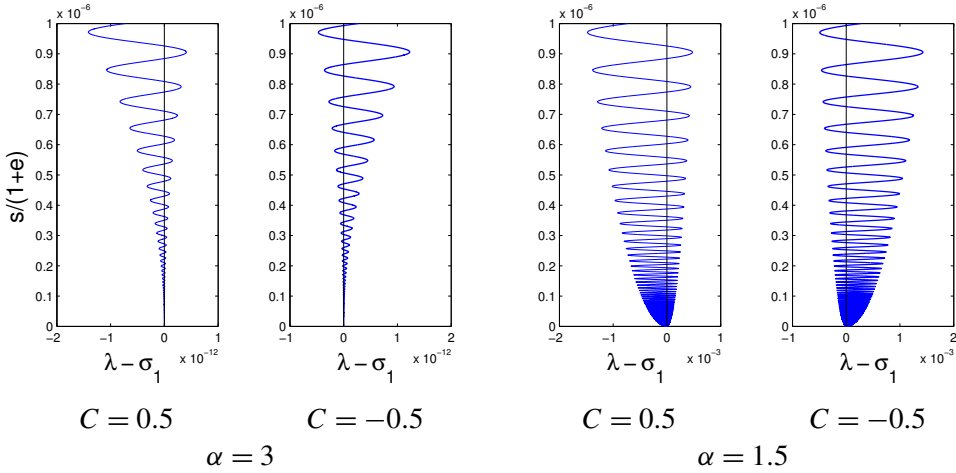


Figure 1. Bifurcation diagram of subcritical and supercritical solutions, containing infinitely many turning points and infinitely many resonant solutions. In all cases, $\beta = -0.35$.

The first eigenvalue σ_1 is simple and, due to Hopf's lemma, we may assume its eigenfunction Φ_1 to be strictly positive in $\overline{\Omega}$ and we take $\|\Phi_1\|_{L^\infty(\partial\Omega)} = 1$.

The case $\alpha + \beta < 1$, $\beta > 0$ was treated in [Arrieta et al. 2010; Castro and Pardo 2011]. Here we focus on the case $\alpha + \beta > 1$, $\beta < 0$, inside of the complementary range. The case with $\alpha < 1$ corresponds to a *bifurcation from infinity* phenomenon; see [Arrieta et al. 2007; 2009; 2010; Castro and Pardo 2011; Rabinowitz 1973]. In contrast, the case with $\alpha > 1$ corresponds to a *bifurcation from zero* phenomenon; see [Arrieta et al. 2007; Crandall and Rabinowitz 1971; Rabinowitz 1971].

The oscillatory situation is in principle more complex than the monotone one, since order techniques such as sub- and supersolutions are not applicable.

One novelty in problem (1-1) is that the parameter appears explicitly in the boundary condition. With respect to this parameter, we perform an analysis of the local bifurcation diagram of nonnegative solutions to (1-1), which turns out to be different from the case $\alpha < 1$ (see Figure 1 for $\alpha > 1$ and Figure 2 for $\alpha < 1$).

Throughout this paper we make the following assumptions:

- (H1) $g : \mathbb{R} \times \partial\Omega \times \mathbb{R} \rightarrow \mathbb{R}$ is a Carathéodory function (i.e. $g = g(\lambda, x, s)$ is measurable in $x \in \Omega$, and continuous with respect to $(\lambda, s) \in \mathbb{R} \times \mathbb{R}$). Moreover, there exist $G_1 \in L^r(\partial\Omega)$ with $r > N - 1$ and continuous functions $\Lambda : \mathbb{R} \rightarrow \mathbb{R}^+$, and $U : \mathbb{R} \rightarrow \mathbb{R}^+$, satisfying

$$\begin{cases} |g(\lambda, x, s)| \leq \Lambda(\lambda)G_1(x)U(s) & \text{for all } (\lambda, x, s) \in \mathbb{R} \times \partial\Omega \times \mathbb{R}, \\ \limsup_{|s| \rightarrow 0} \frac{U(s)}{|s|^\alpha} < +\infty & \text{for some } \alpha > 1. \end{cases}$$

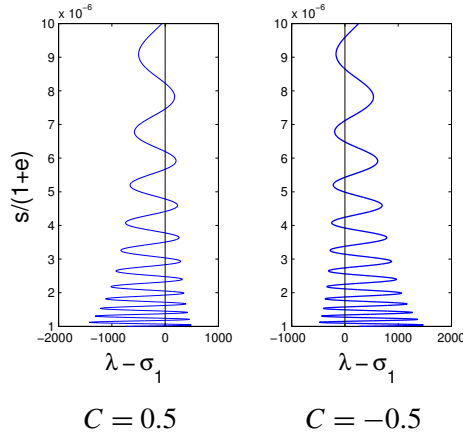


Figure 2. Bifurcation diagram in the case $\alpha = 0.5$, $\beta = -0.35$.

(H2) The partial derivative $g_s(\lambda, \cdot, \cdot)$ (where $g_s := \partial g / \partial s$) belongs to $C(\partial\Omega \times \mathbb{R})$; moreover, $g_s(\cdot, \cdot, 0) = 0$ and there exist $F_1 \in L^r(\partial\Omega)$, with $r > N - 1$ and $\rho > 1$ such that

$$\frac{|g(\lambda, x, s) - sg_s(\lambda, x, s)|}{|s|^\rho} \leq F_1(x) \quad \text{as } \lambda \rightarrow \sigma_1,$$

for $x \in \partial\Omega$ and $s \leq \epsilon$ small enough.

Throughout this paper, by solutions to (1-1) we mean elements $u \in H^1(\Omega)$ such that

$$(1-4) \quad \int_{\Omega} (\nabla u \cdot \nabla v + uv) \, dx = \lambda \int_{\partial\Omega} uv \, d\sigma + \int_{\partial\Omega} g(\lambda, x, u)v \, d\sigma \quad \text{for all } v \in H^1(\Omega).$$

As proven in [Arrieta et al. 2007, Proposition 2.3], all such solutions are in the Holder space $C^\beta(\bar{\Omega})$ for some $\beta > 0$. Moreover, there exists a connected set of positive solutions of (1-1) known as a *branch bifurcating from zero*; see [Arrieta et al. 2007, Theorem 8.1]. We denote it by $\mathcal{C}^+ \subset \mathbb{R} \times C(\bar{\Omega})$, and recall that for $(\lambda, u_\lambda) \in \mathcal{C}^+$

$$u = s\Phi_1 + w, \quad \text{with } w = o(|s|) \quad \text{and} \quad |\sigma_1 - \lambda| = o(1) \quad \text{as } |s| \rightarrow 0.$$

Definition 1.1. A solution (λ^*, u^*) of (1-1) in the branch of solutions $\mathcal{C}^+ \subset \mathbb{R} \times C(\bar{\Omega})$ is called a *turning point* if there is a neighborhood W of (λ^*, u^*) in $\mathbb{R} \times C(\bar{\Omega})$ such that, either $W \cap \mathcal{C}^+ \subset [\lambda^*, \infty) \times C(\bar{\Omega})$ or $W \cap \mathcal{C}^+ \subset (-\infty, \lambda^*] \times C(\bar{\Omega})$.

Our goal is to give conditions on the nonlinear oscillatory term g that guarantee the existence of sequences accumulating to zero of *subcritical* solutions (i.e., for

values of the parameter $\lambda < \sigma_1$), *supercritical* solutions (i.e., for $\lambda > \sigma_1$), *resonant* solutions (i.e., for $\lambda = \sigma_1$), and turning points.

Our main result, Theorem 1.3 below, is exemplified by the case in which g is given by (1-2). In fact we have:

Theorem 1.2. *Assume that g is given by (1-2) with $\beta < 0$. If*

$$|C| < 1 \quad \text{and} \quad \alpha + \beta > 1,$$

then in any neighborhood of the bifurcation point $(\sigma_1, 0)$ in $\mathbb{R} \times C(\overline{\Omega})$, the branch \mathcal{C}^+ of positive solutions of (1-1) contains a sequence of subcritical solutions, a sequence of supercritical solutions, a sequence of turning points, and a sequence of resonant solutions.

The proof of this follows directly from the next theorem.

Theorem 1.3. *Assume the nonlinearity g satisfies hypotheses (H1) and (H2). Assume also that*

$$(1-5) \quad \left| \frac{g(\lambda, x, s) - g(\sigma_1, x, s)}{|s|^\alpha} \right| \rightarrow 0 \quad \text{as} \quad \lambda \rightarrow \sigma_1, s \rightarrow 0$$

pointwise in x .

Let $G : \mathbb{R} \times C(\overline{\Omega}) \rightarrow \mathbb{R}$ be defined by

$$(1-6) \quad G(\lambda, u) := \int_{\partial\Omega} \frac{ug(\lambda, \cdot, u)}{|u|^{1+\alpha}} \Phi_1^{1+\alpha}.$$

If there exist sequences $\{s_n\}, \{s'_n\}$ converging to 0^+ , such that

$$(1-7) \quad \lim_{n \rightarrow +\infty} G(\sigma_1, s'_n \Phi_1) < 0 < \lim_{n \rightarrow +\infty} G(\sigma_1, s_n \Phi_1),$$

then:

(i) *For sufficiently large $n \gg 1$, if (λ, u) is a solution of (1-1) with*

$$P(u) := \frac{\int_{\partial\Omega} u \Phi_1}{\int_{\partial\Omega} \Phi_1^2} = s_n,$$

then (λ, u) is subcritical. Similarly, if $P(u) = s'_n$ it is supercritical. Consequently, there exist two sequences of solutions of (1-1), $\{(\lambda_n, u_n)\}$ and $\{(\lambda'_n, u'_n)\}$ converging to $(\sigma_1, 0)$ as $n \rightarrow \infty$, one of them subcritical, $\lambda_n < \sigma_1$, and the other supercritical, $\lambda'_n > \sigma_1$.

(ii) *There is a sequence converging to zero of turning points $\{(\lambda_n^*, u_n^*)\}$ such that*

$$\lambda_n^* \rightarrow \sigma_1 \quad \text{and} \quad \|u_n^*\|_{L^\infty(\partial\Omega)} \rightarrow 0 \quad \text{as} \quad n \rightarrow \infty.$$

In fact, we can always choose two subsequences of turning points, one of them subcritical, $\lambda_{2n+1}^ < \sigma_1$, and the other supercritical, $\lambda_{2n}^* > \sigma_1$.*

(iii) *There is a sequence converging to zero of resonant solutions; i.e., there are infinitely many solutions $\{(\sigma_1, \tilde{u}_n)\}$ of (1-1) with $\|\tilde{u}_n\|_{L^\infty(\partial\Omega)} \rightarrow 0$.*

The behavior of positive solutions to (1-1) bifurcating from $(\sigma_1, 0)$ described in Theorems 1.2 and 1.3 is similar to that of the solutions bifurcating from (σ_1, ∞) for the sublinear problem; see [Arrieta et al. 2010] for details.

The complex nature of the nonlinearity in (1-2), makes an exhaustive analysis of the global bifurcation diagram outside the scope of this work.

In [Korman 2008] the author considers in the case $\alpha = 1$ $\beta = 1$. He assumes either $N = 1$ or Ω to be a ball and the nonlinearity to be bounded by a constant small enough. He obtains what he calls an oscillatory bifurcation. We refer the reader to [García-Melián et al. 2009] for related problems with nonlinear boundary conditions.

Organization of the paper. Section 2 contains the proof of our main result, giving sufficient conditions for having subcritical, supercritical, and resonant solutions. Section 3 presents two examples; explicit resonant solutions for the oscillatory nonlinearity (1-2) and the one-dimensional case.

2. Subcritical, supercritical and resonant solutions

In this section we give sufficient conditions for the existence of a branch of solutions to (1-1) bifurcating from zero which is neither *subcritical* ($\lambda < \sigma_1$), nor *supercritical*, ($\lambda < \sigma_1$). From this, we conclude the existence of infinitely many *turning points*, see Definition 1.1, and an infinite number of solutions for the resonant problem, i.e. for $\lambda = \sigma_1$. This is achieved in Theorem 1.3

At this step, we analyze when the parameter may cross the first Steklov eigenvalue. To do that, we look at the asymptotic growth rate of the nonlinear term

$$(2-1) \quad \underline{\mathbf{G}}_{0^+} := \int_{\partial\Omega} \liminf_{(\lambda, s) \rightarrow (\sigma_1, 0)} \frac{sg(\lambda, \cdot, s)}{|s|^{1+\alpha}} \Phi_1^{1+\alpha}$$

for $\alpha > 1$. Changing \liminf to \limsup we define the number $\overline{\mathbf{G}}_{0^+}$. If $\underline{\mathbf{G}}_{0^+} > 0$ then \mathcal{C}^+ is subcritical, and if $\overline{\mathbf{G}}_{0^+} < 0$ then \mathcal{C}^+ is supercritical in a neighborhood of $(\sigma_1, 0)$ See [Arrieta et al. 2009, Theorems 3.4 and 3.5] for the bifurcation from infinity case. In this paper we consider nonlinearities for which

$$\underline{\mathbf{G}}_{0^+} < 0 < \overline{\mathbf{G}}_{0^+}.$$

We shall argue as in [Arrieta et al. 2010] for the bifurcation from infinity case. To determine whether a sequence of solutions (λ_n, u_n) is subcritical or supercritical,

one must check the sign of

$$(2-2) \quad \liminf_{n \rightarrow \infty} G(\lambda_n, u_n) \quad \text{and} \quad \limsup_{n \rightarrow \infty} G(\lambda_n, u_n),$$

where G is defined by (1-6). This is done in Lemma 2.3.

In Proposition 2.2, it is proved that when g is such that

$$|g(\lambda, x, s)| = O(|s|^\alpha) \quad \text{as } |s| \rightarrow 0 \text{ for some } \alpha > 1,$$

then the solutions in \mathcal{C}^\pm can be described as

$$u_n = s_n \Phi_1 + w_n, \quad \text{where} \quad \int_{\partial\Omega} w_n \Phi_1 = 0 \quad \text{and} \quad w_n = O(|s_n|^\alpha) \text{ as } n \rightarrow 0.$$

We unveil the signs of the expressions in (2-2) by just looking at the signs of the expressions in (2-2) at $\lambda_n = \sigma_1$ and $u_n = s_n \Phi_1$. This is achieved in Lemma 2.4.

For this we first consider a family of linear Steklov problems with a variable nonhomogeneous term at the boundary h depending on the parameter λ

$$(2-3) \quad \begin{cases} -\Delta u + u = 0 & \text{in } \Omega, \\ \frac{\partial u}{\partial n} = \lambda u + h(\lambda, x) & \text{on } \partial\Omega, \end{cases}$$

where $h(\lambda, \cdot) \in L^r(\partial\Omega)$, $r > N - 1$ and $\lambda \in (-\infty, \sigma_2)$.

We use the decomposition

$$L^r(\partial\Omega) = \text{span}[\Phi_1] \oplus \text{span}[\Phi_1]^\perp,$$

where

$$\text{span}[\Phi_1]^\perp := \left\{ u \in L^r(\partial\Omega) : \int_{\partial\Omega} u \Phi_1 = 0 \right\}.$$

For $h(\lambda, \cdot) \in L^r(\partial\Omega)$, with $r > N - 1$, we write

$$(2-4) \quad h(\lambda, \cdot) = a_1(\lambda) \Phi_1 + h_1(\lambda, \cdot),$$

with

$$a_1(\lambda) = \frac{\int_{\partial\Omega} h(\lambda, \cdot) \Phi_1}{\int_{\partial\Omega} \Phi_1^2}, \quad \int_{\partial\Omega} h_1(\lambda, \cdot) \Phi_1 = 0.$$

For $\lambda \neq \sigma_1$ the solution $u = u(\lambda)$ of (2-3) has a unique decomposition

$$(2-5) \quad u = \frac{a_1(\lambda)}{\sigma_1 - \lambda} \Phi_1 + w, \quad \text{where} \quad \int_{\partial\Omega} w \Phi_1 = 0,$$

and $w = w(\lambda) \in \text{span}[\Phi_1]^\perp$ solves the problem

$$(2-6) \quad \begin{cases} -\Delta w + w = 0 & \text{in } \Omega, \\ \frac{\partial w}{\partial n} = \lambda w + h_1(\lambda, x) & \text{on } \partial\Omega. \end{cases}$$

Note that in (2-6) $w(\lambda) \in \text{span}[\Phi_1]^\perp$ is also well defined for $\lambda = \sigma_1$. Moreover:

Lemma 2.1. *For each compact set $K \subset (-\infty, \sigma_2) \subset \mathbb{R}$ there exists a constant $C = C(K)$, independent of λ , such that*

$$\|w(\lambda)\|_{L^\infty(\partial\Omega)} \leq C \|h_1(\lambda, \cdot)\|_{L^r(\partial\Omega)} \quad \text{for any } \lambda \in K,$$

where $w \in \text{span}[\Phi_1]^\perp$ is the solution of (2-6) and $h_1 \in \text{span}[\Phi_1]^\perp$ is defined in (2-4).

Proof. See Lemma 3.1 of [Arrieta et al. 2010]. □

Now we turn our attention to the nonlinear problem (1-1). Recall that for solutions (λ, u) close to the bifurcation point $(\sigma_1, 0)$ we have

$$(2-7) \quad u = s\Phi_1 + w, \quad \text{where } w \in \text{span}[\Phi_1]^\perp$$

satisfies

$$(2-8) \quad w = o(s) \quad \text{as } s \rightarrow 0.$$

We define

$$(2-9) \quad P(u) := \frac{\int_{\partial\Omega} u(\cdot) \Phi_1}{\int_{\partial\Omega} \Phi_1^2}.$$

Next, we give sufficient conditions on the nonlinear term g in (1-1), for $w = O(|s|^\alpha)$ as $s \rightarrow 0$; compare (2-8). We restrict ourselves below to the branch of positive solutions; a completely analogous result holds for the branch of negative solutions. The next result is essentially Proposition 3.2 in [Arrieta et al. 2010] rewritten for $s \rightarrow 0$; we include the proof for completeness.

Proposition 2.2. *Assume g satisfies hypotheses (H1) and (H2). There exists an open set $\mathbb{O} \subset \mathbb{R} \times C(\overline{\Omega})$ of the form $\mathbb{O} = \{(\lambda, u) : |\lambda - \sigma_1| < \delta_0, \|u\|_{L^\infty(\Omega)} < s_0\}$, for some δ_0 and s_0 , satisfying these conditions:*

- (i) *There exists a constant C_1 independent of λ such that, if $(\lambda, u) \in \mathbb{C}^+ \cap \mathbb{O}$ and $(\lambda, u) \neq (\sigma_1, 0)$ then $u = s\Phi_1 + w$, where $s > 0$, $w \in \text{span}[\Phi_1]^\perp$ and*

$$\|w\|_{L^\infty(\partial\Omega)} \leq C_1 \|G_1\|_{L^r(\partial\Omega)} |s|^\alpha \quad \text{as } |s| \rightarrow 0.$$

- (ii) *There exists a constant $S_0 > 0$ such that for all $|s| \leq S_0$ there exists (λ, u) in $\mathbb{C}^+ \cap \mathbb{O}$ satisfying $u = s\Phi_1 + w$, with $w \in \text{span}[\Phi_1]^\perp$.*

- (iii) *Moreover, for any $(\lambda, u) \in \mathbb{C}^+ \cap \mathbb{O}$, $u = s\Phi_1 + w$, with $w \in \text{span}[\Phi_1]^\perp$,*

$$|\sigma_1 - \lambda| \leq C_2 |s|^{\alpha-1} \quad \text{as } |s| \rightarrow 0,$$

with C_2 independent of λ ; in fact,

$$C_2 = \frac{2\|G_1\|_{L^1(\partial\Omega)}}{\int_{\partial\Omega} \Phi_1^2}.$$

Proof. From (2-7) and (2-8), we have $\Phi_1 + w/s \rightarrow \Phi_1$ as $s \rightarrow 0$ in $L^\infty(\partial\Omega)$. Together with (H1) and Lemma 2.1, this implies that $\|w\|_{L^\infty(\partial\Omega)} \leq C|s|^\alpha$ as $s \rightarrow 0$. This proves part (i).

To prove part (ii) note that $\mathcal{C}^+ \cap \mathcal{C}$ is connected. Hence, using the decomposition in (2-7), we have $u = s\Phi_1 + w$ with $w \in \text{span}[\Phi_1]^\perp$. Since the projection P is continuous, by (2-9), the set

$$\{s \in \mathbb{R} : (1-1) \text{ has a solution of the form } u = s\Phi_1 + w \text{ and } w \in [\text{span}[\Phi_1]^\perp]\}$$

contains an interval in \mathbb{R} containing zero.

To prove part (iii) we observe that if (λ, u) is a solution of (1-1), $u = s\Phi_1 + w$, with $w \in \text{span}[\Phi_1]^\perp$, multiplying the equation by the first Steklov eigenfunction $\Phi_1 > 0$ and integrating by parts we obtain,

$$(\sigma_1 - \lambda)s \int_{\partial\Omega} \Phi_1^2 = \int_{\partial\Omega} g(\lambda, x, s\Phi_1 + w)\Phi_1.$$

Taking into account that

$$\frac{|g(\lambda, x, s\Phi_1 + w)|}{|s|} = \frac{|g(\lambda, x, s\Phi_1 + w)|}{|s\Phi_1 + w|} \left| \Phi_1 + \frac{w}{s} \right| \rightarrow 0 \quad \text{as } s \rightarrow 0$$

we get $\lambda \rightarrow \sigma_1$ as $s \rightarrow 0$.

Moreover, from (H1), we obtain that

$$\begin{aligned} |g(\lambda, x, s\Phi_1 + w)| &= |s|^\alpha \frac{|g(\lambda, x, s\Phi_1 + w)|}{|s\Phi_1 + w|^\alpha} \left| \Phi_1 + \frac{w}{s} \right|^\alpha \\ &\leq C|s|^\alpha G_1(x) \left| \Phi_1 + \frac{w}{s} \right|^\alpha, \end{aligned}$$

and therefore

$$|\sigma_1 - \lambda| \leq C \frac{|s|^{\alpha-1}}{\int_{\partial\Omega} \Phi_1^2} \int_{\partial\Omega} G_1(x) \left| \Phi_1 + \frac{w}{s} \right|^\alpha \Phi_1 \leq C \|G_1\|_{L^r(\partial\Omega)} |s|^{\alpha-1},$$

which ends the proof. \square

Our next result is essentially Lemma 3.1 in [Arrieta et al. 2009] rewritten for $s \rightarrow 0$. It allows us to estimate $\sigma_1 - \lambda_n$ as λ_n converges σ_1 .

Lemma 2.3. *Assume the nonlinearity g satisfies hypotheses (H1) and (H2). Let (λ_n, u_n) be a sequence of solutions of (1-1) with $\lambda_n \rightarrow \sigma_1$ and $\|u_n\|_{L^\infty(\partial\Omega)} \rightarrow 0$.*

If $u_n > 0$ then

$$\begin{aligned}
 (2-10) \quad \frac{\underline{G}_{0+}}{\int_{\partial\Omega} \Phi_1^2} &\leq \frac{1}{\int_{\partial\Omega} \Phi_1^2} \liminf_{n \rightarrow \infty} G(\lambda_n, u_n) \\
 &\leq \liminf_{n \rightarrow \infty} \frac{\sigma_1 - \lambda_n}{\|u_n\|_{L^\infty(\partial\Omega)}^{\alpha-1}} \leq \limsup_{n \rightarrow \infty} \frac{\sigma_1 - \lambda_n}{\|u_n\|_{L^\infty(\partial\Omega)}^{\alpha-1}} \\
 &\leq \frac{1}{\int_{\partial\Omega} \Phi_1^2} \limsup_{n \rightarrow \infty} G(\lambda_n, u_n) \leq \frac{\overline{G}_{0+}}{\int_{\partial\Omega} \Phi_1^2}.
 \end{aligned}$$

A similar statement is obtained for the case $u_n < 0$, just replacing \underline{G}_{0+} by \underline{G}_{0-} and \overline{G}_{0+} by \overline{G}_{0-} .

Proof. We show that $u_n > 0$; the other case has a similar proof. Consider a family of solutions u_n of (1-1) for $\lambda = \lambda_n$ with $\lambda_n \rightarrow \sigma_1$ and $0 < u_n \rightarrow 0$. Multiplying (1-1) by Φ_1 and integrating by parts, we get

$$(2-11) \quad (\sigma_1 - \lambda_n) \int_{\partial\Omega} u_n \Phi_1 = \int_{\partial\Omega} g(\lambda_n, x, u_n) \Phi_1.$$

But

$$\int_{\partial\Omega} g(\lambda_n, x, u_n) \Phi_1 = \|u_n\|_{L^\infty(\partial\Omega)}^\alpha \int_{\partial\Omega} \frac{g(\lambda_n, x, u_n)}{u_n^\alpha} \left(\frac{u_n}{\|u_n\|_{L^\infty(\partial\Omega)}} \right)^\alpha \Phi_1.$$

Taking into account the definition of $G(\lambda, u)$ in (1-6), we can write

$$\begin{aligned}
 \int_{\partial\Omega} \frac{g(\lambda_n, x, u_n)}{u_n^\alpha} \left(\frac{u_n}{\|u_n\|_{L^\infty(\partial\Omega)}} \right)^\alpha \Phi_1 \\
 = \int_{\partial\Omega} \frac{g(\lambda_n, x, u_n)}{u_n^\alpha} \left[\left(\frac{u_n}{\|u_n\|_{L^\infty(\partial\Omega)}} \right)^\alpha - \Phi_1^\alpha \right] \Phi_1 + G(\lambda_n, u_n).
 \end{aligned}$$

Moreover,

$$\int_{\partial\Omega} \frac{g(\lambda_n, x, u_n)}{u_n^\alpha} \left[\left(\frac{u_n}{\|u_n\|_{L^\infty(\partial\Omega)}} \right)^\alpha - \Phi_1^\alpha \right] \Phi_1 \rightarrow 0 \quad \text{as } n \rightarrow \infty,$$

because $u_n/\|u_n\|_{L^\infty(\partial\Omega)} \rightarrow \Phi_1$ uniformly in $\partial\Omega$.

But, firstly from the above, secondly from Fatou's lemma, and thirdly from definition of \underline{G}_{0+} ,

$$\begin{aligned}
 (2-12) \quad \liminf_{n \rightarrow \infty} \int_{\partial\Omega} \frac{g(\lambda_n, x, u_n)}{u_n^\alpha} \left(\frac{u_n}{\|u_n\|_{L^\infty(\partial\Omega)}} \right)^\alpha \Phi_1 \\
 \geq \liminf_{n \rightarrow \infty} G(\lambda_n, u_n) \geq \int_{\partial\Omega} \liminf_{n \rightarrow \infty} \frac{g(\lambda_n, x, u_n)}{u_n^\alpha} \Phi_1^{1+\alpha} \\
 \geq \underline{G}_{0+}.
 \end{aligned}$$

Dividing both sides of (2-11) by $\|u_n\|_{L^\infty(\partial\Omega)}^\alpha$ and passing to the limit we obtain the first two inequalities in the chain (2-10). The third inequality in the chain is trivial and the last two are obtained in a similar manner. \square

Let $\{s_n\}$ and $\{s'_n\}$ satisfy

$$(2-13) \quad -\infty < \lim_{n \rightarrow +\infty} G(\sigma_1, s'_n \Phi_1) < 0 < \lim_{n \rightarrow +\infty} G(\sigma_1, s_n \Phi_1) < \infty.$$

In order to prove Theorem 1.3, we show that the signs in (2-2) can be deduced from those of (2-13). This is stated in the following result, which is a slight variation of [Arrieta et al. 2010, Lemma 3.3].

Lemma 2.4. *Assume that g satisfies hypotheses (H1), (H2), and (1-5).*

If $(\lambda_n, s_n) \rightarrow (\sigma_1, 0)$ and there exists a constant C such that $\|w_n\|_{L^\infty(\partial\Omega)} \leq C|s_n|^\alpha$ for all $n \rightarrow 0$, then

$$\liminf_{n \rightarrow +\infty} G(\lambda_n, s_n \Phi_1 + w_n) \geq \liminf_{n \rightarrow +\infty} G(\sigma_1, s_n \Phi_1),$$

where G is given by (1-6). Similarly,

$$\limsup_{n \rightarrow +\infty} G(\lambda_n, s_n \Phi_1 + w_n) \leq \limsup_{n \rightarrow +\infty} G(\sigma_1, s_n \Phi_1).$$

Proof. Throughout this proof, C denotes several constants depending only on (Ω, g) . Given $\varepsilon > 0$, assume that $|(\lambda_n, s_n) - (\sigma_1, 0)| < \varepsilon$.

By the mean value theorem we have

$$(2-14) \quad \begin{aligned} g(\lambda_n, x, s_n \Phi_1 + w_n) - g(\lambda_n, x, s_n \Phi_1) &= w_n \int_0^1 g_s(\lambda_n, \cdot, s_n \Phi_1 + \tau w_n) d\tau \\ &\leq \|w_n\|_{L^\infty(\partial\Omega)} \sup_{\tau \in [0,1]} \|g_s(\lambda_n, \cdot, s_n \Phi_1 + \tau w_n)\|_{L^\infty(\partial\Omega)}. \end{aligned}$$

Therefore

$$(2-15) \quad \begin{aligned} \int_{\partial\Omega} \left| g(\lambda_n, x, s_n \Phi_1 + w_n) - g(\lambda_n, x, s_n \Phi_1) \right| \Phi_1 dx &\leq \|w_n\|_{L^\infty(\partial\Omega)} \int_{\partial\Omega} \sup_{\tau \in [0,1]} \|g_s(\lambda_n, \cdot, s_n \Phi_1 + \tau w_n)\|_{L^\infty(\partial\Omega)} \\ &\leq |\partial\Omega| \|w_n\|_{L^\infty(\partial\Omega)} \sup_{\tau \in [0,1]} \|g_s(\lambda_n, \cdot, s_n \Phi_1 + \tau w_n)\|_{L^\infty(\partial\Omega)}. \end{aligned}$$

By hypotheses (H1) and (H2), for all $x \in \partial\Omega$,

$$(2-16) \quad \begin{aligned} \frac{|g_s(\lambda_n, x, s)|}{|s|^{\gamma-1}} &\leq |s|^{\rho-\gamma} F_1(x) + C|s|^{\alpha-\gamma} G_1(x) \max\{\Lambda(\lambda_n), n \geq 1\} =: D_1(x), \end{aligned}$$

for n large, and $\gamma = \min\{\rho, \alpha\} > 1$. Hence, $D_1 \in L^r(\partial\Omega)$ with $r > N - 1$ and

$$(2-17) \quad \sup_{|s| \leq 1/n} |g_s(\lambda_n, x, s)| \leq D_1(x) \left(\frac{1}{n}\right)^{\gamma-1}, \quad \text{with } \gamma > 1.$$

Since $\|w_n\|_{L^\infty(\partial\Omega)} = O(|s_n|^\alpha)$, we obtain from (2-15) and (2-17)

$$(2-18) \quad \int_{\partial\Omega} \frac{|g(\lambda_n, \cdot, s_n \Phi_1 + w_n) - g(\lambda_n, \cdot, s_n \Phi_1)|}{|s_n|^\alpha} \Phi_1$$

$$\leq C \sup_{\tau \in [0,1]} \|g_s(\lambda_n, \cdot, s_n \Phi_1 + \tau w_n)\|_{L^\infty(\partial\Omega)}$$

$$\leq C \sup_{|s| \leq 1/n} \|g_s(\lambda_n, \cdot, s)\|_{L^\infty(\partial\Omega)},$$

which tends to 0 as $n \rightarrow \infty$.

Therefore

$$\liminf_{n \rightarrow +\infty} \int_{\partial\Omega} \frac{s_n g(\lambda_n, \cdot, s_n \Phi_1 + w_n)}{|s_n|^{1+\alpha}} \Phi_1$$

$$\geq \lim_{n \rightarrow \infty} \int_{\partial\Omega} \frac{s_n g(\lambda_n, \cdot, s_n \Phi_1 + w_n) - s_n g(\lambda_n, \cdot, s_n \Phi_1)}{|s_n|^{1+\alpha}} \Phi_1$$

$$+ \liminf_{n \rightarrow +\infty} \int_{\partial\Omega} \frac{s_n g(\lambda_n, \cdot, s_n \Phi_1)}{|s_n|^{1+\alpha}} \Phi_1$$

$$= \liminf_{n \rightarrow +\infty} \int_{\partial\Omega} \frac{s_n g(\lambda_n, \cdot, s_n \Phi_1)}{|s_n|^{1+\alpha}} \Phi_1$$

$$= \liminf_{n \rightarrow +\infty} \int_{\partial\Omega} \frac{s_n g(\sigma_1, \cdot, s_n \Phi_1)}{|s_n|^{1+\alpha}} \Phi_1,$$

where we have used (2-18) and (1-5) respectively.

Now note that, multiplying and dividing by $|\Phi_1 + w_n/s_n|^\alpha$ the integrand of the left hand side above can be written as

$$\frac{s_n g(\lambda_n, \cdot, s_n \Phi_1 + w_n)}{|s_n|^{1+\alpha}} \Phi_1 = \frac{(s_n \Phi_1 + w_n) g(\lambda_n, \cdot, s_n \Phi_1 + w_n)}{|s_n \Phi_1 + w_n|^{1+\alpha}} \left| \Phi_1 + \frac{w_n}{s_n} \right|^\alpha.$$

Then, (H2) and the fact that $\Phi_1 + w_n/s_n \rightarrow \Phi_1$ in $L^\infty(\partial\Omega)$ concludes the proof. \square

Now we prove the first main result in this paper. Roughly speaking, it states that if there are a sequence of subcritical solutions and another of supercritical solutions, since the solution set is connected, there are infinitely many turning points and infinitely many resonant solutions. We prove the result for the positive branch. The same conclusions can be attained for the connected branch of negative solutions bifurcating from zero.

Proof of Theorem 1.3. From Proposition 2.2(ii), consider any two sequences of solutions of (1-1), such that $(\lambda_n, u_n) \rightarrow (\sigma_1, 0)$ and $(\lambda'_n, u'_n) \rightarrow (\sigma_1, 0)$ in \mathcal{C}^+ with

$$P(u_n) = \frac{\int_{\partial\Omega} u_n \Phi_1}{\int_{\partial\Omega} \Phi_1^2} = s_n \quad \text{and} \quad P(u'_n) = \frac{\int_{\partial\Omega} u'_n \Phi_1}{\int_{\partial\Omega} \Phi_1^2} = s'_n.$$

Writing $u_n = s_n \Phi_1 + w_n$, with $w_n \in \text{span}[\Phi_1]^\perp$, from Proposition 2.2(i), we have $\|w_n\|_{L^\infty(\partial\Omega)} = O(|s_n|^\alpha)$. From Lemmata 2.3, and 2.4, hypotheses (1-5) and (1-7) we get that

$$\begin{aligned} \liminf_{n \rightarrow \infty} \frac{\sigma_1 - \lambda_n}{\|u_n\|_{L^\infty(\partial\Omega)}^{\alpha-1}} &\geq \frac{1}{\int_{\partial\Omega} \Phi_1^2} \liminf_{n \rightarrow \infty} \int_{\partial\Omega} \frac{(s_n \Phi_1 + w_n)g(\lambda_n, \cdot, s_n \Phi_1 + w_n)}{|s_n \Phi_1 + w_n|^{1+\alpha}} \Phi_1^{1+\alpha} \\ &\geq \frac{1}{\int_{\partial\Omega} \Phi_1^2} \liminf_{n \rightarrow +\infty} \int_{\partial\Omega} \frac{s_n g(\lambda_n, \cdot, s_n \Phi_1)}{|s_n|^{1+\alpha}} \Phi_1 \\ &= \frac{1}{\int_{\partial\Omega} \Phi_1^2} \liminf_{n \rightarrow +\infty} \int_{\partial\Omega} \frac{s_n g(\sigma_1, \cdot, s_n \Phi_1)}{|s_n|^{1+\alpha}} \Phi_1 > 0, \end{aligned}$$

and therefore $\lambda_n < \sigma_1$.

Analogously, for (λ'_n, u'_n) we get $\lambda'_n > \sigma_1$. Hence (i) is proved.

To prove (ii), assume, by choosing subsequences if necessary, that $s_n > s'_n > s_{n+1}$ for all $n \geq 0$ and that $0 < s_n, s'_n \leq S_0$ where S_0 is the one from Proposition 2.2(ii). In particular, by parts (i) and (ii) of Proposition 2.2, if $(\lambda, u) \in \mathcal{C}^+$ and $P(u) = s < S_0$ then $\|u\|_{L^\infty(\partial\Omega)} \leq (1 + C_1 \|G_1\|_{L^r(\partial\Omega)} |S_0|^{\alpha-1})s$. Taking S_0 small enough we may assume that $\|u\|_{L^\infty(\partial\Omega)} \leq 2s$.

Let

$$(2-19) \quad K_n = \{(\lambda, u) \in \mathcal{C}^+ : P(u) = s \text{ and } s_n \geq s \geq s_{n+1}\}.$$

Let us see that, for each $n \in \mathbb{N}$, K_n is a compact subset of $\mathbb{R} \times C(\bar{\Omega})$. Let $\{(\mu_k, v_k)\}$ be a sequence in K_n . Without loss of generality we may assume that $\{\mu_k\}$ converges to μ^* . Since $v_k = t_k \Phi_1 + w_k$ with $w_k = O(|t_k|^\alpha)$ and $s_n \geq t_k =: P(v_k) \geq s_{n+1}$, for all k , we have $\|v_k\|_{C(\partial\Omega)} \leq t_k + \|w_k\|_{C(\partial\Omega)} \leq C$ with C independent of k . This together with Proposition 2.3 of [Arrieta et al. 2007] yields

$$(2-20) \quad \|v_k\|_{C(\bar{\Omega})} \leq C_1(1 + \|v_k\|_{C(\partial\Omega)}) \leq C,$$

where, again, C is independent of k . Since the embedding $C^\gamma(\bar{\Omega}) \rightarrow C^{\gamma'}(\bar{\Omega})$ is compact for $0 < \gamma' < \gamma$ we may further assume that the sequence $\{v_k\}$ converges to some $u^* \in C^{\gamma'}(\bar{\Omega})$. This, hypothesis (H1) and the dominated convergence theorem imply that $\{g(\mu_k, \cdot, v_k)\}$ converges to $g(\mu^*, \cdot, u^*)$ in $L^r(\partial\Omega)$. Therefore, since

$$(2-21) \quad \begin{cases} -\Delta v_k + v_k = 0 & \text{in } \Omega \\ \frac{\partial v_k}{\partial n} = \mu_k v_k + g(\mu_k, x, v_k) & \text{on } \partial\Omega, \end{cases}$$

passing to the limit in the weak sense we have

$$(2-22) \quad \begin{cases} -\Delta u^* + u^* = 0 & \text{in } \Omega, \\ \frac{\partial u^*}{\partial n} = \mu^* u^* + g(\lambda^*, x, u^*) & \text{on } \partial\Omega. \end{cases}$$

By the continuity of the projection operator we also have $s_n \geq s^* = P(u^*) = \lim_{k \rightarrow \infty} P(v_k) \geq s_{n+1}$. Hence $(\mu^*, u^*) \in K_n$, which proves that K_n is compact.

Since $s_n > s'_n > s_{n+1}$ there exists $(\lambda, u) \in K_n$ with $\lambda > \sigma_1$. Hence, if we define

$$(2-23) \quad \lambda_n^* = \sup\{\lambda : (\lambda, u) \in K_n\},$$

then $\lambda_n^* \geq \lambda'_n > \sigma_1$ see part (i). From the compactness of K_n there exists u_n^* such that $(\lambda_n^*, u_n^*) \in K_n$. From the definition of λ_n^* if (λ, u) is a solution of (1-1) with $s_n > P(u_n) > s_{n+1}$, then $\lambda \leq \lambda_n^*$ which proves that (λ_n^*, u_n^*) is a (supercritical) turning point.

With a completely symmetric argument, using the sets

$$K'_n = \{(\lambda, u) \in \mathcal{C}^+ : P(u) = s \text{ and } s'_n \geq s \geq s'_{n+1}\}$$

and defining $\lambda'_n = \inf\{\lambda : (\lambda, u) \in K'_n\}$, we show the existence of u_n^* such that $(\lambda'_n, u_n^*) \in K'_n$ is a (subcritical) turning point.

In order to prove the existence of resonant solutions, we now show that there exists $n_0 \in \mathbb{N}$ such that for each $n \geq n_0$ both sets K_n and K'_n contain resonant solutions, that is, solutions of the form (σ_1, u) .

We use a *reductio ad absurdum* argument for the sets K_n . If this is not the case, then there will exist a sequence of integers numbers $n_j \rightarrow +\infty$ such that K_{n_j} does not contain any resonant solution. This implies that the compact sets $K_{n_j}^+ = \{(\lambda, u) \in K_{n_j} : \lambda \geq \sigma_1\}$ can be written as

$$K_{n_j}^+ = \mathcal{C}^+ \cap \{(\lambda, u) \in \mathbb{R} \times C(\partial\Omega) : \lambda > \sigma_1, s_{n_j} > P(u) > s_{n_j+1}\};$$

therefore $K_{n_j}^+$ contains at least a connected component of \mathcal{C}^+ . Moreover it is nonempty since we know that there exists at least one solution (λ, u) with $P(u) = s'_{n_j} \in (s_{n_j}, s_{n_j+1})$ and therefore $\lambda > \sigma_1$. The fact that we can construct a sequence of connected components of \mathcal{C}^+ contradicts the fact that \mathcal{C}^+ is a connected near $(\sigma_1, 0) \in \mathbb{R} \times C(\bar{\Omega})$.

A completely symmetric argument can be applied to the sets K'_n . □

3. Two examples

3.1. Resonant solutions for the oscillatory nonlinearity (1-2). In [Arrieta et al. 2007, Theorem 8.1] it is proved that if $\alpha > 1$, for any $\beta \in \mathbb{R}$, and $C \in \mathbb{R}$ there is an unbounded branch of positive solutions. Assume from now that $\beta < 0$.

Taking $|C| \leq 1$ it is not difficult to see that

$$u_k(x) := [\sin(-C) + k\pi]^{1/\beta} \Phi_1(x), \quad k \geq 0,$$

defines a sequence of resonant solutions to (1-1) such that $u_k(x) \rightarrow 0$ as $k \rightarrow \infty$.

3.2. A one-dimensional example. Now we consider the one-dimensional version of (1-1), where most computations can be made explicit.

Let $\{\sigma_i\}$ denote the sequence of *Steklov* eigenvalues of the problem (1-3). For $N > 1$ the Steklov eigenvalues form an increasing sequence of real numbers, $\{\sigma_i\}_{i=1}^{\infty}$ while for $N = 1$ there are only two Steklov eigenvalues as we made explicit below.

Observe that Equation (1-1) in the one-dimensional domain $\Omega = (0, 1)$ reads

$$\begin{cases} -u_{xx} + u = 0 & \text{in } (0, 1), \\ -u_x(0) = \lambda u + g(\lambda, 0, u(0)), \\ u_x(1) = \lambda u + g(\lambda, 1, u(1)). \end{cases}$$

The general solution of the differential equation is $u(x) = ae^x + be^{-x}$ and therefore the nonlinear boundary conditions provide two nonlinear equations in terms of two constants a and b . The function $u = ae^x + be^{-x}$ is a solution if (λ, a, b) satisfy

$$\begin{pmatrix} -(1+\lambda) & (1-\lambda) \\ (1-\lambda)e & -(1+\lambda)e^{-1} \end{pmatrix} \begin{pmatrix} a \\ b \end{pmatrix} = \begin{pmatrix} g(\lambda, 0, a+b) \\ g(\lambda, 1, ae+be^{-1}) \end{pmatrix}.$$

In this case we only have two Steklov eigenvalues,

$$\sigma_1 = \frac{e-1}{e+1} < \sigma_2 = \frac{1}{\sigma_1} = \frac{e+1}{e-1}.$$

Restricting the analysis to symmetric solutions $u_s(x) = s(e^x + e^{1-x})$, with $s \in \mathbb{R}$, and choosing $g(\lambda, x, s) = g(s)$, it is easy to prove that $u_s(x)$ is a solution if and only if λ satisfies

$$(3-1) \quad \lambda(s) = \sigma_1 - \frac{g(s(e+1))}{s(e+1)}, \quad s > 0.$$

Therefore, whenever $g(u) = o(u)$ at zero, there is a branch of solutions $(\lambda(s), u_s)$ converging to $(\sigma_1, 0)$ as $s \rightarrow 0$.

Now fix $g(s) = s^\alpha \sin(s^\beta)$ for an arbitrary $\alpha > 1$, $\beta < 0$. From the definition in (2-1) we can write

$$\underline{G}_{0^+} := \int_{\partial\Omega} \liminf_{s \rightarrow 0^+} \frac{sg(s)}{|s|^{1+\alpha}} \Phi^{1+\alpha} = \int_{\partial\Omega} \liminf_{s \rightarrow 0^+} \sin(s^\beta) \Phi^{1+\alpha} = - \int_{\partial\Omega} \Phi^{1+\alpha} < 0,$$

$$\overline{G}_{0^+} := \int_{\partial\Omega} \limsup_{s \rightarrow 0^+} \frac{sg(s)}{|s|^{1+\alpha}} \Phi^{1+\alpha} = \int_{\partial\Omega} \limsup_{s \rightarrow 0^+} \sin(s^\beta) \Phi^{1+\alpha} = \int_{\partial\Omega} \Phi^{1+\alpha} > 0,$$

and then $\underline{G}_{0^+} < 0 < \overline{G}_{0^+}$.

Moreover, by looking in (3-1) at the values of $s \in \mathbb{R}$ such that $\lambda(s) = \sigma_1$ it is easy to check that (σ_1, u_k) is a solution for any $k \in \mathbb{Z}$, where

$$u_k(x) := \frac{(k\pi)^{1/\beta}}{e+1} (e^x + e^{1-x});$$

that is, there is a sequence of solutions of the resonant problem converging to zero, as shown in Figure 3.

Moreover, computing in (3-1) the local maxima and minima of $\lambda(s)$ it is not difficult to check that (λ_k^*, u_k^*) is a sequence of turning points converging to zero, where

$$\lambda_k^* := \sigma_1 - t_k^{(\alpha-1)/\beta} \sin(t_k), \quad u_k^*(x) := t_k^{1/\beta} (e^x + e^{1-x})$$

and where t_k is such that

$$\tan(t_k) = -\frac{\beta}{\alpha-1} t_k, \quad t_k \in [-\pi/2 + k\pi, \pi/2 + k\pi]$$

with $t_k \rightarrow \infty$ and $t_k^{1/\beta} \rightarrow 0$ as $k \rightarrow \infty$ thanks to $\beta < 0$.

Let us observe that the bifurcation from zero phenomena occurs whenever $\alpha > 1$ for any β and that whenever $\alpha + \beta < 1$ the number of oscillations grows up faster than the number of oscillations of multiples of the eigenfunction and cannot be controlled; compare the two parts of Figure 3.

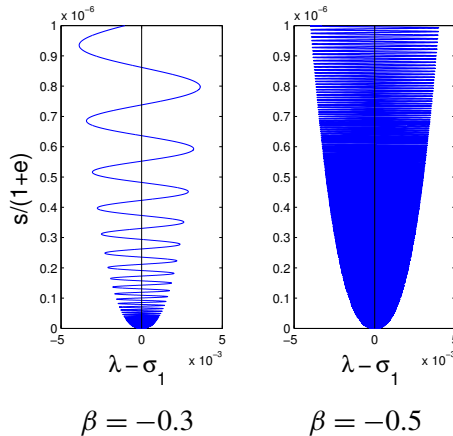


Figure 3. Bifurcation diagram in the case $\alpha = 1.4$, for two values of β .

References

- [Arrieta et al. 2007] J. M. Arrieta, R. Pardo, and A. Rodríguez-Bernal, “Bifurcation and stability of equilibria with asymptotically linear boundary conditions at infinity”, *Proc. Roy. Soc. Edinburgh Sect. A* **137**:2 (2007), 225–252. MR 2009d:35194 Zbl 1202.35119
- [Arrieta et al. 2009] J. M. Arrieta, R. Pardo, and A. Rodríguez-Bernal, “Equilibria and global dynamics of a problem with bifurcation from infinity”, *J. Differential Equations* **246**:5 (2009), 2055–2080. MR 2010c:35016 Zbl 1195.35040
- [Arrieta et al. 2010] J. M. Arrieta, R. Pardo, and A. Rodríguez-Bernal, “Infinite resonant solutions and turning points in a problem with unbounded bifurcation”, *Internat. J. Bifur. Chaos Appl. Sci. Engrg.* **20**:9 (2010), 2885–2896. MR 2011m:35107 Zbl 1202.35085
- [Castro and Pardo 2011] A. Castro and R. Pardo, “Infinitely many stability switches in a problem with sublinear oscillatory boundary conditions”, preprint MA-UCM 2011-2, Univ. Complutense de Madrid, 2011, available at <http://www.ucm.es/centros/cont/descargas/documento23873.pdf>.
- [Crandall and Rabinowitz 1971] M. G. Crandall and P. H. Rabinowitz, “Bifurcation from simple eigenvalues”, *J. Functional Analysis* **8** (1971), 321–340. MR 44 #5836 Zbl 0219.46015
- [García-Melián et al. 2009] J. García-Melián, J. D. Rossi, and J. C. Sabina de Lis, “An elliptic system with bifurcation parameters on the boundary conditions”, *J. Differential Equations* **247**:3 (2009), 779–810. MR 2010f:35073 Zbl 1171.35037
- [Korman 2008] P. Korman, “An oscillatory bifurcation from infinity, and from zero”, *NoDEA Non-linear Differential Equations Appl.* **15**:3 (2008), 335–345. MR 2009h:35137 Zbl 1166.34010
- [Rabinowitz 1971] P. H. Rabinowitz, “Some global results for nonlinear eigenvalue problems”, *J. Functional Analysis* **7** (1971), 487–513. MR 46 #745 Zbl 0212.16504
- [Rabinowitz 1973] P. H. Rabinowitz, “On bifurcation from infinity”, *J. Differential Equations* **14** (1973), 462–475. MR 48 #7047 Zbl 0272.35017

Received June 3, 2011. Revised November 22, 2011.

ALFONSO CASTRO
DEPARTMENT OF MATHEMATICS
HARVEY MUDD COLLEGE
301 PLATT BOULEVARD
CLAREMONT, CA 91711
UNITED STATES
castro@math.hmc.edu

ROSA PARDO
DEPARTAMENTO DE MATEMÁTICA APLICADA
UNIVERSIDAD COMPLUTENSE DE MADRID
AVENIDA COMPLUTENSE S/N
28040 MADRID
SPAIN
rpardo@mat.ucm.es

RELATIVE MEASURE HOMOLOGY AND CONTINUOUS BOUNDED COHOMOLOGY OF TOPOLOGICAL PAIRS

ROBERTO FRIGERIO AND CRISTINA PAGLIANTINI

Measure homology was introduced by Thurston in his notes about the geometry and topology of 3-manifolds, where it was exploited in the computation of the simplicial volume of hyperbolic manifolds. Zastrow and Hansen independently proved that there exists a canonical isomorphism between measure homology and singular homology (on the category of CW-complexes), and it was then shown by Löh that, in the absolute case, such isomorphism is in fact an isometry with respect to the L^1 -seminorm on singular homology and the total variation seminorm on measure homology. Löh's result plays a fundamental rôle in the use of measure homology as a tool for computing the simplicial volume of Riemannian manifolds.

This paper deals with an extension of Löh's result to the relative case. We prove that relative singular homology and relative measure homology are isometrically isomorphic for a wide class of topological pairs. Our results can be applied for instance in computing the simplicial volume of Riemannian manifolds with boundary.

Our arguments are based on new results about continuous (bounded) cohomology of topological pairs, which are probably of independent interest.

1. Introduction

Measure homology was introduced in [Thurston 1979], where it was exploited in the proof that the simplicial volume of a closed hyperbolic n -manifold is equal to its Riemannian volume divided by a constant only depending on n (this result is attributed in [Thurston 1979] to Gromov). In order to rely on measure homology, it is necessary to know that this theory “coincides” with the usual real singular homology, at least for a large class of spaces. The proof that measure homology and real singular homology of CW-pairs are isomorphic has appeared in [Hansen 1998; Zastrow 1998]. However, in order to exploit measure homology as a tool for computing the simplicial volume, one has to show that these homology theories are not only isomorphic, but also *isometric* (with respect to the seminorms introduced below). In the absolute case, this result is achieved in [Löh 2006]. Our paper is

MSC2010: primary 55N10, 55N35; secondary 20J06, 55U15, 57N65.

Keywords: simplicial volume, singular homology, bounded cohomology of groups, CAT(0) spaces.

devoted to extending Löh's result to the context of relative homology of topological pairs. As mentioned in [Fujiwara and Manning 2011, Appendix A] and [Löh 2007, Remark 4.22], such an extension seems to raise difficulties that suggest that Löh's argument should not admit a straightforward translation into the relative context. For a detailed account about the notion of measure homology and its applications see, e.g., the introductions of [Zastrow 1998; Berlanga 2008].

In order to achieve our main results, we develop some aspects of the theory of continuous bounded cohomology of topological pairs. More precisely, we compare such theory with the usual bounded cohomology of pairs of groups and spaces. Park [2003] provided the algebraic foundations to the theory of relative bounded cohomology, extending Ivanov's [1985] homological algebra approach to the relative case. However, Park endows the bounded cohomology of a pair of spaces with a seminorm which is *a priori* different from the seminorm considered in this paper. In fact, the most common definition of simplicial volume is based on a specific L^1 -seminorm on singular homology, whose dual is just the L^∞ -seminorm on bounded cohomology defined in [Gromov 1982, Section 4.1]. This seminorm does not coincide *a priori* with Park's seminorm, so our results cannot be deduced from Park's arguments. More precisely, it is shown in [Park 2003, Theorem 4.6] that Gromov's and Park's norms are bi-Lipschitz equivalent (see Theorem 6.1 below). In [Park 2003, page 206] it is stated that it remains unknown if this equivalence is actually an isometry. In Section 6 we answer this question in the negative, providing examples showing that Park's and Gromov's seminorms indeed do not coincide in general.

1A. Relative singular homology of pairs. Let X be a topological space and $W \subseteq X$ a (possibly empty) subspace of X . For $n \in \mathbb{N}$ we denote by $C_n(X)$ the module of singular n -chains with real coefficients, i.e., the \mathbb{R} -module freely generated by the set $S_n(X)$ of singular n -simplices with values in X . The natural inclusion of W in X induces an inclusion of $C_n(W)$ into $C_n(X)$, and we denote by $C_n(X, W)$ the quotient space $C_n(X)/C_n(W)$. The usual differential of the complex $C_*(X)$ defines a differential $d_*: C_*(X, W) \rightarrow C_{*-1}(X, W)$. The homology of the resulting complex is the usual relative singular homology of the topological pair (X, W) , and will be denoted by $H_*(X, W)$.

The real vector space $C_n(X, W)$ can be endowed with a natural L^1 -norm, as follows. If $\alpha \in C_n(X, W)$, then

$$\|\alpha\|_1 = \inf \left\{ \sum_{\sigma \in S_n(X)} |a_\sigma|, \text{ where } \alpha = \left[\sum_{\sigma \in S_n(X)} a_\sigma \sigma \right] \text{ in } C_n(X)/C_n(W) \right\}.$$

Such a norm descends to a seminorm on $H_n(X, W)$, which is defined as follows: if $[\alpha] \in H_n(X, W)$, then

$$\|[\alpha]\|_1 = \inf\{\|\beta\|_1 \mid \beta \in C_n(X, W), d_n\beta = 0, [\beta] = [\alpha]\}$$

(this seminorm can be null on nonzero elements of $H_n(X, W)$). Of course, we recover the absolute homology modules of X just by setting $H_n(X) = H_n(X, \emptyset)$.

1B. Relative measure homology of pairs. We now recall the definition of relative measure homology of the pair (X, W) . We endow $S_n(X)$ with the compact-open topology and denote by $\mathbb{B}_n(X)$ the σ -algebra of Borel subsets of $S_n(X)$. If μ is a signed measure on $\mathbb{B}_n(X)$ (in this case we say for short that μ is a Borel measure on $S_n(X)$), the *total variation of μ* is defined by the formula

$$\|\mu\|_m = \sup_{A \in \mathbb{B}_n(X)} \mu(A) - \inf_{B \in \mathbb{B}_n(X)} \mu(B) \in [0, +\infty]$$

(the subscript m stands for *measure*). For every $n \geq 0$, the measure chain module $\mathcal{C}_n(X)$ is the real vector space of the Borel measures on $S_n(X)$ having finite total variation and admitting a compact determination set. The graded module $\mathcal{C}_*(X)$ can be given the structure of a complex via the boundary operator

$$\begin{aligned} \partial_n : \mathcal{C}_n(X) &\rightarrow \mathcal{C}_{n-1}(X), \\ \mu &\mapsto \sum_{j=0}^n (-1)^j \mu^j, \end{aligned}$$

where μ^j is the push-forward of μ under the map that takes a simplex $\sigma \in S_n(X)$ into the composition of σ with the usual inclusion of the standard $(n-1)$ -simplex onto the j -th face of σ .

Let now W be a (possibly empty) subspace of X . It is proved in [Zastrow 1998, Proposition 1.10] that the σ -algebra $\mathbb{B}_n(W)$ of Borel subsets of $S_n(W)$ coincides with the set $\{A \cap S_n(W) \mid A \in \mathbb{B}_n(X)\}$. For every $\mu \in \mathcal{C}_n(W)$, the assignment

$$\nu(A) = \mu(A \cap S_n(W)), \quad A \in \mathbb{B}_n(X),$$

defines a Borel measure on $S_n(X)$, which is called the *extension* of μ . If μ has compact determination set and finite total variation then the same is true for ν , so that we have a natural inclusion $\mathcal{C}_n(W) \hookrightarrow \mathcal{C}_n(X)$ (see [Zastrow 1998, Proposition 1.10 and Lemma 1.11] for full details). The image of $\mathcal{C}_n(W)$ in $\mathcal{C}_n(X)$ will be simply denoted by $\mathcal{C}_n(W)$, and coincides with the set of the elements of $\mathcal{C}_n(X)$ which admit a compact determination set contained in $S_n(W)$. We denote by $\mathcal{C}_n(X, W)$ the quotient $\mathcal{C}_n(X)/\mathcal{C}_n(W)$.

It is readily seen that $\partial_n(\mathcal{C}_n(W)) \subseteq \mathcal{C}_{n-1}(W)$, so ∂_n induces a boundary operator $\mathcal{C}_n(X, W) \rightarrow \mathcal{C}_{n-1}(X, W)$, which will still be denoted by ∂_n . The homology of the complex $(\mathcal{C}_*(X, W), \partial_*)$ is the *relative measure homology of the pair (X, W)* , and it is denoted by $\mathcal{H}_*(X, W)$.

Just as in the case of singular homology, we may endow $\mathcal{H}_n(X, W)$ with a seminorm as follows. For every $\alpha \in \mathcal{C}_n(X, W)$ we set

$$\|\alpha\|_m = \inf \{ \|\mu\|_m, \text{ where } \mu \in \mathcal{C}_n(X), [\mu] = \alpha \text{ in } \mathcal{C}_n(X, W) = \mathcal{C}_n(X)/\mathcal{C}_n(W) \}.$$

Then, for every $[\alpha] \in \mathcal{H}_n(X, W)$ we set

$$\|[\alpha]\|_{\text{mh}} = \inf\{\|\beta\|_{\text{m}} \mid \beta \in \mathcal{C}_n(X, W), \partial_n \beta = 0, [\beta] = [\alpha]\}$$

(the subscript mh stands for *measure homology*). The absolute measure homology module $\mathcal{H}_n(X)$ can be defined just by setting $\mathcal{H}_n(X) = \mathcal{H}_n(X, \emptyset)$.

1C. Relative singular homology versus relative measure homology. For every $\sigma \in \mathcal{S}_n(X)$ let us denote by δ_σ the atomic measure supported by the singleton $\{\sigma\} \subseteq \mathcal{S}_n(X)$. The chain map

$$\begin{aligned} \iota_* : C_*(X, W) &\rightarrow \mathcal{C}_*(X, W), \\ \sum_{i=0}^k a_i \sigma_i &\mapsto \sum_{i=0}^k a_i \delta_{\sigma_i} \end{aligned}$$

induces a map

$$H_n(\iota_*) : H_n(X, W) \rightarrow \mathcal{H}_n(X, W), \quad n \in \mathbb{N},$$

which is obviously norm-nonincreasing for every $n \in \mathbb{N}$.

Theorem 1.1 [Zastrow 1998; Hansen 1998]. *Let (X, W) be a CW-pair. For every $n \in \mathbb{N}$, the map*

$$H_n(\iota_*) : H_n(X, W) \rightarrow \mathcal{H}_n(X, W)$$

is an isomorphism.

Zastrow's and Hansen's proofs of Theorem 1.1 are based on the fact that relative measure homology satisfies the Eilenberg–Steenrod axioms for homology (on suitable categories of topological pairs). Therefore, their approach avoids the explicit construction of the inverse maps $H_n(\iota_*)^{-1}$, $n \in \mathbb{N}$, and does not give much information about the behavior of such inverse maps with respect to the seminorms introduced above. In the case when $W = \emptyset$, the fact that $H_n(\iota_*)$ is indeed an isometry was proved by Löh:

Theorem 1.2 [Löh 2006]. *If X is any connected CW-complex, then for every $n \in \mathbb{N}$ the map*

$$H_n(\iota_*) : H_n(X) \rightarrow \mathcal{H}_n(X)$$

is an isometric isomorphism.

Löh's proof of Theorem 1.2 exploits deep results about the *bounded cohomology* of groups and topological spaces. In Section 3 and Section 4 we develop a suitable relative version of such results, which we use on page 125 to prove this:

Theorem 1.3. *Let (X, W) be a CW-pair, and let us suppose that the following conditions hold:*

- (1) X (whence W) is countable, and both X and W are connected;

(2) the map $\pi_j(W) \rightarrow \pi_j(X)$ induced by the inclusion $W \hookrightarrow X$ is injective for $j = 1$, and it is an isomorphism for $j \geq 2$.

Then, for every $n \in \mathbb{N}$ the isomorphism

$$H_n(\iota_*) : H_n(X, W) \rightarrow \mathcal{H}_n(X, W)$$

is isometric.

In fact, we will deduce Theorem 1.3 from Theorem 1.7 below concerning the relationships between continuous (bounded) cohomology and singular (bounded) cohomology of topological pairs.

Definition 1.4. A CW-pair (X, W) is *good* if it satisfies conditions (1) and (2) in the statement of Theorem 1.3.

We conjecture that Theorem 1.3 holds even without the hypothesis that the pair (X, W) is good, so a brief comment about the places where this assumption comes into play is in order. The fact that W is connected and π_1 -injective in X allows us to exploit results regarding the bounded cohomology of a pair (G, A) , where G is a group and A is a subgroup of G . In order to deal with the case when W is *not* assumed to be π_1 -injective, one could probably build on results regarding the bounded cohomology of a pair (G, A) , where A, G are groups and $\varphi : A \rightarrow G$ is a homomorphism of A into G . This case is treated in [Park 2003] by means of a mapping cone construction. However, the mapping cone introduced there does not admit a norm inducing Gromov's seminorm in bounded cohomology, so Park's approach seems to be of no help to our purposes. Perhaps it is easier to drop from the hypotheses of Theorem 1.3 the requirement that W be connected (provided that we still assume that every component of W is π_1 -injective in X). Several arguments in our proofs make use of cone constructions which are based on the choice of a basepoint in the universal coverings \tilde{X}, \tilde{W} of X, W . When W is connected (and π_1 -injective in X), the space \tilde{W} is realized as a connected subset of \tilde{X} , and this allows us to define compatible cone constructions on \tilde{X} and \tilde{W} . It is not clear how to replace these constructions when W is disconnected: one could probably build on the theory of homology and cohomology of a group with respect to any system of subgroups, as described for instance in [Bieri and Eckmann 1978] (see also [Mineyev and Yaman 2007]), but several difficulties arise which we have not been able to overcome. Finally, the assumption that $\pi_i(W)$ is isomorphic to $\pi_i(X)$ for every $i \geq 2$ plays a fundamental rôle in our proof of Proposition 4.7 below. One could get rid of this assumption by using a result stated without proof in [Park 2003, Lemma 4.2], but at the moment we are not able to provide a proof for Park's statement (see Remark 4.9 for a brief discussion of this issue).

1D. *Locally convex pairs.* We are able to prove that measure homology is isometric to singular homology also for a large family of pairs of metric spaces, namely for those pairs which support a *relative straightening* for simplices.

The *straightening procedure* for simplices was introduced in [Thurston 1979], and establishes an isometric isomorphism between the usual singular homology of a space and the homology of the complex of *straight* chains. Such a procedure was originally defined on hyperbolic manifolds, and has then been extended to the context of nonpositively curved Riemannian manifolds. In Section 2 we give the precise definition of *locally convex pair of metric spaces*. Then, following some ideas described in [Löh and Sauer 2009], for every locally convex pair (X, W) we define a straightening procedure which induces a chain map between relative measure chains and relative singular chains. It turns out that such a straightening induces a well-defined norm-nonincreasing map $\mathcal{H}_n(X, W) \rightarrow H_n(X, W)$. This map provides the desired norm-nonincreasing inverse of $H_n(\iota_*)$, so that we can prove (in Section 2D) the following:

Theorem 1.5. *Let (X, W) be a locally convex pair of metric spaces. Then the map*

$$H_n(\iota_*) : H_n(X, W) \rightarrow \mathcal{H}_n(X, W)$$

is an isometric isomorphism for every $n \in \mathbb{N}$.

The class of locally convex pairs is indeed quite large, including for example all the pairs $(M, \partial M)$, where M is a nonpositively curved complete Riemannian manifold with geodesic boundary ∂M .

Remark 1.6. Suppose that (X, W) is a locally convex pair, and let K be a connected component of W . An easy application of a metric version of Cartan–Hadamard theorem (see [Bridson and Haefliger 1999, II.4.1]) shows that $\pi_1(K)$ injects into $\pi_1(X)$, and $\pi_i(K) = \pi_i(X) = 0$ for every $i \geq 2$. In particular, if (X, W) is also a countable CW-pair and W is connected, then (X, W) is good, and the conclusion of Theorem 1.5 also descends from Theorem 1.3. Note however that the request that W be connected could be quite restrictive in several applications of our results. For example, it is well-known that the natural compactification of a complete finite-volume hyperbolic manifold with geodesic boundary and/or cusps is a manifold with boundary N admitting a locally CAT(0) (whence locally convex) metric that turns the pair $(N, \partial N)$ into a locally convex pair (see [Bridson and Haefliger 1999, pages 362–366], for example). We have discussed in [Frigerio and Pagliantini 2010] some properties of the simplicial volume of such manifolds, and in that context several interesting examples have in fact disconnected boundary. In [Pagliantini 2012] it is shown how to apply Theorem 1.5 for getting shorter proofs of the main results of [Frigerio and Pagliantini 2010].

1E. (Continuous) relative bounded cohomology. As mentioned above, our proof of Theorem 1.3 involves the study of the relative bounded cohomology of topological pairs. Introduced in [Gromov 1982], the relative bounded cohomology of pairs (of groups or spaces) seems to be less clearly understood than absolute bounded cohomology. Here below we define the *continuous* (bounded) cohomology of topological pairs, and we put on (continuous) bounded cohomology Gromov's L^∞ -seminorm which is "dual" (in a sense to be specified below) to the seminorm on (measure) homology described above. Then, in Section 4 we compare the continuous bounded cohomology of a good CW-pair to its usual singular bounded cohomology (see Theorem 1.7 below). In Section 5 we show how this result implies Theorem 1.3.

Let us now state more precisely our results. For every $n \in \mathbb{N}$ we denote by $C^n(X)$ and $C^n(X, W)$ the algebraic duals of $C_n(X)$ and $C_n(X, W)$ (that is, the respective modules of singular n -cochains with real coefficients). We will often identify $C^n(X, W)$ with a submodule of $C^n(X)$ via the canonical isomorphism

$$C^n(X, W) \cong \{f \in C^n(X) \mid f|_{C_n(W)} = 0\}.$$

If $\delta^* : C^*(X, W) \rightarrow C^{*+1}(X, W)$ is the usual differential, the homology of the complex $(C^*(X, W), \delta^*)$ is the relative singular cohomology of the pair (X, W) , and it is denoted by $H^*(X, W)$.

We regard $S_n(X)$ as a subset of $C_n(X)$, so that for every cochain $\varphi \in C^n(X, W) \subseteq C^n(X)$ it makes sense to consider the restriction $\varphi|_{S_n(X)}$. In particular, we say that φ is *continuous* if $\varphi|_{S_n(X)}$ is (recall that $S_n(X)$ is endowed with the compact-open topology). If we set

$$C_c^*(X, W) = \{\varphi \in C^*(X, W) \mid \varphi \text{ is continuous}\},$$

then it is readily seen that $\delta^n(C_c^n(X, W)) \subseteq C_c^{n+1}(X, W)$, so $C_c^*(X, W)$ is a subcomplex of $C^*(X, W)$, whose homology is denoted by $H_c^*(X, W)$.

We now come to the definition of (continuous) bounded cohomology. We endow $C^n(X, W)$ with the L^∞ -norm defined by

$$\|f\|_\infty = \sup_{\sigma \in S_n(X)} |f(\sigma)| \in [0, \infty], \quad f \in C^n(X, W),$$

and introduce the following submodules of $C^*(X, W)$:

$$C_b^*(X, W) = \{f \in C^*(X, W) \mid \|f\|_\infty < \infty\},$$

$$C_{cb}^*(X, W) = C_b^*(X, W) \cap C_c^*(X, W).$$

The coboundary map δ^n is bounded, so $C_b^*(X, W)$ (resp. $C_{cb}^*(X, W)$) is a subcomplex of $C^*(X, W)$ (resp. of $C_c^*(X, W)$). Its homology is denoted by $H_b^*(X, W)$ (resp. $H_{cb}^*(X, W)$), and it is called the *bounded cohomology* (resp. *continuous*

bounded cohomology) of (X, W) . The L^∞ -norm on $C^*(X, W)$ descends (after suitable restrictions) to a seminorm on each of the modules $H^*(X, W)$, $H_c^*(X, W)$, $H_b^*(X, W)$, $H_{cb}^*(X, W)$. These seminorms will still be denoted by $\|\cdot\|_\infty$. The inclusion maps

$$\rho_b^* : C_{cb}^*(X, W) \hookrightarrow C_b^*(X, W), \quad \rho_c^* : C_c^*(X, W) \hookrightarrow C^*(X, W)$$

induce maps

$$H^*(\rho_b^*) : H_{cb}^*(X, W) \rightarrow H_b^*(X, W), \quad H^*(\rho_c^*) : H_c^*(X, W) \rightarrow H^*(X, W),$$

that are a priori neither injective nor surjective.

We are now ready to state our main result about (continuous) bounded cohomology of pairs, which is proved in Section 4E:

Theorem 1.7. *Let (X, W) be a good CW-pair. Then the map*

$$H^n(\rho_b^*) : H_{cb}^n(X, W) \rightarrow H_b^n(X, W)$$

admits a right inverse which is an isometric embedding (in particular, $H^n(\rho_b^)$ is surjective) for every $n \in \mathbb{N}$.*

In the absolute case, when $W = \emptyset$, Theorem 1.7 is proved in [Frigerio 2011, Theorem 1.2]. In order to prove Theorem 1.7 we suitably develop the theory of relative bounded cohomology of pairs of groups. In particular, our Theorem 4.1 implies the following result, which is maybe of independent interest (see Section 3 for the definition of $H_b^*(G, A)$, where G is a group and A is a subgroup of G):

Theorem 1.8. *Let (X, W) be a countable CW-pair. Also suppose that X, W are connected, and that the map $\pi_1(W) \rightarrow \pi_1(X)$ induced by the inclusion $W \hookrightarrow X$ is injective. Then for every $n \in \mathbb{N}$ there exists a norm-nonincreasing isomorphism*

$$H_b^n(\pi_1(X), \pi_1(W)) \rightarrow H_b^n(X, W).$$

If in addition the pair (X, W) is good, then this isomorphism is isometric.

In Section 4F we show how Theorem 1.7 and [Frigerio 2011, Theorem 1.1] can be exploited to prove the following:

Theorem 1.9. *Let (X, W) be a locally finite good CW-pair. Then the map*

$$H^n(\rho_c^*) : H_c^n(X, W) \rightarrow H^n(X, W)$$

is an isometric isomorphism for every $n \in \mathbb{N}$.

2. The case of locally convex pairs

The following definitions can be found for instance in [Bridson and Haefliger 1999]. Let (X, d) be a metric space (when d is fixed, we denote (X, d) simply by X). A *geodesic segment* in X is an isometric embedding of a bounded closed interval into X . The metric d (or the metric space $X = (X, d)$) is *geodesic* if every two points in X are joined by a geodesic segment (in particular, X is path-connected and locally path connected). Moreover, d (or $X = (X, d)$) is *globally convex* if it is geodesic and if any two geodesic segments $c_1 : [0, a] \rightarrow X$, $c_2 : [0, a] \rightarrow X$ such that $c_1(0) = c_2(0)$ satisfy the condition $d(c_1(ta), c_2(ta)) \leq td(c_1(a), c_2(a))$ for every $t \in [0, 1]$ (and in this case, X is contractible, see Lemma 2.1 below). We say that d (or $X = (X, d)$) is *locally convex* if every point in X has a neighborhood in which the restriction of d is convex (in particular, it is geodesic). A subspace $Y \subseteq X$ is *convex* if every geodesic segment (in X) joining any two points of Y is entirely contained in Y (in particular, if X is geodesic, then Y is path-connected).

Suppose that X is geodesic, complete and locally convex. Then it is locally contractible, hence it admits a universal covering $p : \tilde{X} \rightarrow X$. We endow \tilde{X} with the length metric induced by p , that is, the unique length metric \tilde{d} such that $p : (\tilde{X}, \tilde{d}) \rightarrow (X, d)$ is a local isometry (see [Bridson and Haefliger 1999, Proposition I.3.25]). Since (X, d) is complete and geodesic, the same is true for (\tilde{X}, \tilde{d}) . Moreover, the Cartan–Hadamard theorem for metric spaces (see [loc. cit., II.4.1]) implies that the space (\tilde{X}, \tilde{d}) is globally convex.

Let W be any subset of X . We say that (X, W) is a *locally convex pair of metric spaces* (or simply a *locally convex pair*) if the following conditions hold:

- (1) X is geodesic, complete and locally convex;
- (2) W is closed in X and locally path-connected;
- (3) every path-connected component of $p^{-1}(W) \subseteq \tilde{X}$ is convex in \tilde{X} .

Throughout the whole section we denote by (X, W) a locally convex pair of metric spaces, we fix a universal covering $p : \tilde{X} \rightarrow X$ (where \tilde{X} is endowed with the induced metric), and we denote by \tilde{W} the subset $p^{-1}(W) \subseteq \tilde{X}$ (on the contrary, in Section 4 we will denote by \tilde{W} a fixed connected component of $p^{-1}(W)$).

2A. Straight simplices. In order to properly define straight simplices we first need the following result, which is an immediate consequence of the Cartan–Hadamard theorem for metric spaces:

Lemma 2.1 [Bridson and Haefliger 1999, II.4.1]. *For every pair of points $p, q \in \tilde{X}$ there exists a unique geodesic segment in \tilde{X} joining p to q . Moreover, if $\alpha_{p,q} : [0, 1] \rightarrow \tilde{X}$ is a constant-speed parametrization of such a segment, then $\alpha_{p,q}$ continuously depends (with respect to the compact-open topology) on p and q . In particular, \tilde{X} is contractible.*

For $i \in \mathbb{N}$ we denote by e_i the point $(0, 0, \dots, 1, \dots, 0, 0, \dots) \in \mathbb{R}^{\mathbb{N}}$ where the unique nonzero coefficient is at the i -th entry (entries are indexed by \mathbb{N} , so $(1, 0, \dots) = e_0$). We denote by Δ_p the standard p -simplex, that is, the convex hull of e_0, \dots, e_p , and we observe that with these notations we have $\Delta_p \subseteq \Delta_{p+1}$.

Let $k \in \mathbb{N}$, and let x_0, \dots, x_k be points in \tilde{X} . We recall here the well-known definition of *straight simplex* $[x_0, \dots, x_k] \in S_k(\tilde{X})$ with vertices x_0, \dots, x_k : if $k = 0$, then $[x_0]$ is the 0-simplex with image x_0 ; if straight simplices have been defined for every $h \leq k$, then $[x_0, \dots, x_{k+1}] : \Delta_{k+1} \rightarrow \tilde{X}$ is determined by the following condition: for every $z \in \Delta_k \subseteq \Delta_{k+1}$, the restriction of $[x_0, \dots, x_{k+1}]$ to the segment with endpoints z, e_{k+1} is a constant speed parametrization of the geodesic joining $[x_0, \dots, x_k](z)$ to x_{k+1} (the fact that $[x_0, \dots, x_{k+1}]$ is well-defined and continuous is an immediate consequence of Lemma 2.1).

2B. Nets. Let $\Gamma \cong \pi_1(X)$ be the group of covering automorphisms of $p : \tilde{X} \rightarrow X$, and observe that, since p is a local isometry, every element of Γ is an isometry of \tilde{X} .

Definition 2.2. A *net* in \tilde{X} is given by a subset $\tilde{\Lambda} \subseteq \tilde{X}$ and a locally finite collection of Borel sets $\{\tilde{B}_x\}_{x \in \tilde{\Lambda}}$ such that the following conditions hold:

- (1) $\tilde{X} = \bigcup_{x \in \tilde{\Lambda}} \tilde{B}_x$ and $\tilde{B}_x \cap \tilde{B}_y = \emptyset$ for every $x, y \in \tilde{\Lambda}$ with $x \neq y$.
- (2) $\gamma(\tilde{\Lambda}) = \tilde{\Lambda}$ for every $\gamma \in \Gamma$ and $\gamma(\tilde{B}_x) = \tilde{B}_{\gamma(x)}$ for every $x \in \tilde{\Lambda}, \gamma \in \Gamma$.
- (3) If \tilde{K} is a path-connected component of \tilde{W} , then $\tilde{K} \subseteq \bigcup_{x \in \tilde{\Lambda} \cap \tilde{K}} \tilde{B}_x$.

Lemma 2.3. *There exists a net.*

Proof. For every $q \in X$ let us denote by U_q an evenly covered open neighborhood of q in X (with respect to the universal covering $\tilde{X} \rightarrow X$). Since W is closed and locally path-connected, we may also suppose that $W \cap U_q$ is path-connected. Being metrizable, X is paracompact, so the open covering $\{U_q\}_{q \in X}$ admits a locally finite open refinement $\{V_i\}_{i \in I}$. Now fix a total ordering \leq on I in such a way that $i \leq j$ whenever $V_i \cap W \neq \emptyset$ and $V_j \cap W = \emptyset$, and let us set

$$B_i = V_i \setminus \left(\bigcup_{j < i} V_j \right).$$

By construction, the family $\{B_i\}_{i \in I}$ is locally finite in X . Moreover, every B_i is the intersection of an open set and a closed set, so it is a Borel subset of X . Therefore, up to replacing I with the subset $\{i \in I \mid B_i \neq \emptyset\}$, the family $\{B_i\}_{i \in I}$ provides a locally finite cover of X by nonempty Borel sets. For every $i \in I$ let us choose $x_i \in B_i$ in such a way that $x_i \in W$ whenever $B_i \cap W \neq \emptyset$, and let us set $\Lambda = \bigcup_{i \in I} \{x_i\}$. We also set $B_{x_i} = B_i$ for every $i \in I$.

We now define $\tilde{\Lambda} = p^{-1}(\Lambda)$. For every $i \in I$ we choose an element $\tilde{x}_i \in p^{-1}(x_i)$, and we take $q_i \in X$ in such a way that $B_{x_i} \subseteq U_{q_i}$. Being evenly covered, U_{q_i} lifts to

the disjoint union $p^{-1}(U_{q_i}) = \bigcup_{\gamma \in \Gamma} \gamma(\tilde{U}_{q_i})$, where \tilde{U}_{q_i} is the connected component of $p^{-1}(U_{q_i})$ containing \tilde{x}_i .

We are now ready to define \tilde{B}_x , where x is any element of $\tilde{\Lambda}$. In fact, every $x \in \tilde{\Lambda}$ uniquely determines an index $i \in I$ and an element $\gamma \in \Gamma$ such that $x = \gamma(\tilde{x}_i)$, and we can set $\tilde{B}_x = \gamma(\tilde{U}_{q_i} \cap p^{-1}(B_{x_i}))$. Of course \tilde{B}_x is a Borel subset of \tilde{X} .

It is now easy to check that the pair $(\tilde{\Lambda}, \{\tilde{B}_x\}_{x \in \tilde{\Lambda}})$ provides a net: the local finiteness of the family $\{\tilde{B}_x, x \in \tilde{\Lambda}\}$ readily descends from the fact p is a covering and $\{B_x, x \in \Lambda\}$ is locally finite in X , and conditions (1) and (2) of Definition 2.2 are an obvious consequence of our choices. We now show that condition (3) also holds. We fix $x \in \tilde{\Lambda}$ such that $\tilde{W} \cap \tilde{B}_x \neq \emptyset$. By construction we have $x \in \tilde{W}$, and there exist $\gamma \in \Gamma$ and $i \in I$ such that $\tilde{B}_x \subseteq \gamma(\tilde{U}_{q_i})$. Our assumption that $U_q \cap W$ is path-connected implies that $\gamma(\tilde{U}_{q_i}) \cap \tilde{W}$ is also path-connected, so the set $\tilde{B}_x \cap \tilde{W}$ is entirely contained in the path-connected component of \tilde{W} containing x , whence the conclusion. \square

2C. Straightening. We are now ready to define our straightening operator. Let $(\tilde{\Lambda}, \{\tilde{B}_x\}_{x \in \tilde{\Lambda}})$ be a net. We denote by $S_n^{\tilde{\Lambda}}(\tilde{X}) \subseteq S_n(\tilde{X})$ the set of straight n -simplices in \tilde{X} with vertices in $\tilde{\Lambda}$. Then we let $\tilde{\text{str}}_n : C_n(\tilde{X}) \rightarrow C_n(\tilde{X})$ be the unique linear map such that for $\tilde{\sigma} \in S_n(\tilde{X})$

$$\tilde{\text{str}}_n(\tilde{\sigma}) = [x_0, \dots, x_n] \in S_n^{\tilde{\Lambda}}(\tilde{X}),$$

where $x_i \in \tilde{\Lambda}$ is such that $\tilde{\sigma}(e_i) \in \tilde{B}_{x_i}$ for $i = 0, \dots, n$.

Proposition 2.4. *The map $\tilde{\text{str}}_* : C_*(\tilde{X}) \rightarrow C_*(\tilde{X})$ satisfies the following properties:*

- (1) $d_{n+1} \circ \tilde{\text{str}}_{n+1} = \tilde{\text{str}}_n \circ d_{n+1}$ for every $n \in \mathbb{N}$.
- (2) $\tilde{\text{str}}_n(\gamma \circ \tilde{\sigma}) = \gamma \circ \tilde{\text{str}}_n(\tilde{\sigma})$ for every $n \in \mathbb{N}$, $\gamma \in \Gamma$, $\tilde{\sigma} \in S_n(\tilde{X})$.
- (3) $\tilde{\text{str}}_*(C_*(\tilde{W})) \subseteq C_*(\tilde{W})$.
- (4) *The induced chain map $C_*(\tilde{X}, \tilde{W}) \rightarrow C_*(\tilde{X}, \tilde{W})$, which we will still denote by $\tilde{\text{str}}_*$, is Γ -equivariantly homotopic to the identity.*

Proof. If $x_0, \dots, x_n \in \tilde{X}$, then it is easily seen that for every $i \leq n$ the i -th face of $[x_0, \dots, x_n]$ is given by $[x_0, \dots, \hat{x}_i, \dots, x_n]$; moreover since isometries preserve geodesics we have $\gamma \circ [x_0, \dots, x_n] = [\gamma(x_0), \dots, \gamma(x_n)]$ for every $\gamma \in \text{Isom}(\tilde{X})$. Together with property (2) in the definition of net, these facts readily imply points (1) and (2) of the proposition.

If $\tilde{\sigma} \in S_n(\tilde{W})$, then all the vertices of $\tilde{\sigma}$ lie in the same connected component \tilde{K} of \tilde{W} . By property (3) in the definition of net, the vertices of $\tilde{\text{str}}_n(\tilde{\sigma})$ still lie in \tilde{K} . Since (X, W) is a locally convex pair, the subset \tilde{K} is convex in \tilde{X} , so $\tilde{\text{str}}_n(\tilde{\sigma})$ belongs to $S_n(\tilde{W})$, whence (3).

Finally, for $\tilde{\sigma} \in S_n(\tilde{X})$, let $F_{\tilde{\sigma}} : \Delta_n \times [0, 1] \rightarrow \tilde{X}$ be defined by $F_{\tilde{\sigma}}(x, t) = \beta_x(t)$, where $\beta_x : [0, 1] \rightarrow \tilde{X}$ is the constant-speed parametrization of the geodesic

segment joining $\tilde{\sigma}(x)$ with $\tilde{\text{str}}(\tilde{\sigma})(x)$. We set $T_n(\tilde{\sigma}) = (F_{\tilde{\sigma}})_*(c)$, where c is the standard chain triangulating the prism $\Delta_n \times [0, 1]$ by $(n+1)$ -simplices. The fact that $d_{n+1}T_n + T_{n-1}d_n = \text{Id} - \tilde{\text{str}}_n$ is now easily checked, while the Γ -equivariance of T_* is a consequence of property (2) of nets together with the fact that geodesics are preserved by isometries. As above, the fact that $T_n(C_n(\tilde{W})) \subseteq C_{n+1}(\tilde{W})$ is a consequence of the convexity of the components of \tilde{W} . \square

Let $\Lambda = p(\tilde{\Lambda})$, and let $S_*^\Lambda(X)$ be the subset of $S_*(X)$ given by those singular simplices which are obtained by composing a simplex in $S_*^\Lambda(\tilde{X})$ with the covering projection p . As a consequence of Proposition 2.4 we get the following:

Proposition 2.5. *The map $\tilde{\text{str}}_*$ induces a chain map $\text{str}_* : C_*(X, W) \rightarrow C_*(X, W)$ which is homotopic to the identity.*

Remark 2.6. The maps $\tilde{\text{str}}_*$, str_* obviously depend on the net chosen for their construction. Such a dependence is however somewhat inessential in our arguments below. Henceforth we understand that a net $(\tilde{\Lambda}, \{\tilde{B}_x\}_{x \in \tilde{\Lambda}})$ is fixed, and we denote by $\tilde{\text{str}}_*$, str_* the corresponding straightening operators.

We are now ready to construct a chain map $\theta_* : \mathcal{C}_*(X, W) \rightarrow C_*(X, W)$ whose induced map in homology will provide the desired norm-nonincreasing inverse of $H_*(t_*)$.

Fix a simplex $\sigma \in S_n^\Lambda(X)$. It is readily seen that the set $\text{str}_n^{-1}(\sigma)$ is a Borel subset of $S_n(X)$. Therefore, for every measure $\mu \in \mathcal{C}_n(X)$ it makes sense to set

$$c_\sigma(\mu) = \mu(\text{str}_n^{-1}(\sigma)) \in \mathbb{R}.$$

Lemma 2.7. *For every measure $\mu \in \mathcal{C}_n(X)$, the set*

$$\{\sigma \in S_n^\Lambda(X) \mid c_\sigma(\mu) \neq 0\}$$

is finite.

Proof. Since μ admits a compact determination set, it is sufficient to show that the family $\{\text{str}_n^{-1}(\sigma), \sigma \in S_n^\Lambda(X)\}$ is locally finite in $S_n(X)$. So, let us take $\sigma_0 \in S_n(X)$, and let $\tilde{\sigma}_0 \in S_n(\tilde{X})$ be a lift of σ_0 to \tilde{X} . For every $j = 0, \dots, n$, let Z_j be an open neighborhood of $\tilde{\sigma}_0(e_j)$ that intersects only a finite number of \tilde{B}_{x_i} 's, and let $\tilde{\Omega} \subseteq S_n(\tilde{X})$ be the set of n -simplices whose i -th vertex belongs to Z_j for every $i = 0, \dots, n$. Then $\tilde{\Omega}$ is an open neighborhood of $\tilde{\sigma}_0$ in $S_n(\tilde{X})$.

Let $p_n : S_n(\tilde{X}) \rightarrow S_n(X)$ be the map taking every $\tilde{\sigma} \in S_n(\tilde{X})$ into $p \circ \tilde{\sigma}$. It is proved in [Frigerio 2011, Lemma A.4] (see also [Löh 2006]) that p_n is a covering, whence an open map, so $\Omega = p_n(\tilde{\Omega})$ is an open neighborhood of σ_0 in $S_n(X)$. Moreover, by construction the set $\text{str}_n(\Omega) = \text{str}_n(p_n(\tilde{\Omega})) = p_n(\tilde{\text{str}}_n(\tilde{\Omega}))$ is finite, whence the conclusion. \square

By Lemma 2.7 we can define the map

$$\theta_n : \mathcal{C}_n(X) \rightarrow C_n(X), \quad \theta_n(\mu) = \sum_{\sigma \in S_n^\Delta(X)} c_\sigma(\mu)\sigma.$$

Lemma 2.8. (1) $\theta_n \circ \partial_{n+1} = d_{n+1} \circ \theta_{n+1}$ for every $n \in \mathbb{N}$.

(2) $\theta_n(\mathcal{C}_n(W)) \subseteq C_n(W)$ for every $n \in \mathbb{N}$.

(3) $\|\theta_n(\mu)\|_1 \leq \|\mu\|_m$ for every $\mu \in \mathcal{C}_n(X)$, $n \in \mathbb{N}$.

Proof. Point (1) is a direct consequence of the fact that str_* is a chain map.

Since $\text{str}_n(C_n(W)) \subseteq C_n(W)$, if $\sigma \in S_n^\Delta(X) \setminus S_n(W)$, then $\text{str}_n^{-1}(\sigma) \cap S_n(W) = \emptyset$. Therefore, if $\mu \in \mathcal{C}_n(W) \subseteq \mathcal{C}_n(X)$, then $c_\sigma(\mu) = \mu(\text{str}_n^{-1}(\sigma)) = 0$, whence point (2).

Point (3) is a consequence of the fact that, if $\{Z_j\}_{j \in J}$ is a finite collection of pairwise disjoint Borel subsets of $S_n(X)$, then $\sum_{j \in J} |\mu(Z_j)| \leq \|\mu\|_m$. \square

2D. Concluding the proof of Theorem 1.5. As a consequence of Lemma 2.8, the map $\theta_* : \mathcal{C}_*(X) \rightarrow C_*(X)$ induces norm-nonincreasing maps

$$\bar{\theta}_* : \mathcal{C}_*(X, W) \rightarrow C_*(X, W), \quad H_*(\bar{\theta}_*) : \mathcal{H}_*(X, W) \rightarrow H_*(X, W).$$

Since we have already seen that $H_*(\iota_*) : H_*(X, W) \rightarrow \mathcal{H}_*(X, W)$ is a norm-nonincreasing isomorphism, in order to prove that $H_*(\iota_*)$ is an isometry it is sufficient to show that $H_n(\bar{\theta}_*) \circ H_n(\iota_*)$ is the identity of $H_n(X, W)$ for every $n \in \mathbb{N}$. However, we have from the very definitions that $\bar{\theta}_n \circ \iota_n = \text{str}_n$ for every $n \in \mathbb{N}$, so the conclusion follows from Proposition 2.5.

3. Relative bounded cohomology of groups

Let us recall some basic definitions and results about the bounded cohomology of groups. For full details we refer the reader to [Gromov 1982; Ivanov 1985; Monod 2001]. Henceforth, we denote by G a fixed group, which has to be thought as endowed with the discrete topology.

Definition 3.1 [Ivanov 1985; Monod 2001]. A *Banach G -module* is a Banach space V with a (left) action of G such that $\|g \cdot v\| \leq \|v\|$ for every $g \in G$ and every $v \in V$. A G -morphism of Banach G -modules is a bounded G -equivariant linear operator.

From now on we refer to a Banach G -module simply as a G -module.

3A. Relative injectivity. A bounded linear map $\iota : A \rightarrow B$ of Banach spaces is *strongly injective* if there is a bounded linear map $\sigma : B \rightarrow A$ with $\|\sigma\| \leq 1$ and $\sigma \circ \iota = \text{Id}_A$ (in particular, ι is injective). We emphasize that, even when A and B are G -modules, the map σ is *not* required to be G -equivariant.

Definition 3.2. A G -module E is *relatively injective* if for every strongly injective G -morphism $\iota : A \rightarrow B$ of Banach G -modules and every G -morphism $\alpha : A \rightarrow E$ there is a G -morphism $\beta : B \rightarrow E$ satisfying $\beta \circ \iota = \alpha$ and $\|\beta\| \leq \|\alpha\|$.

$$\begin{array}{ccccc}
 0 & \longrightarrow & A & \begin{array}{c} \xleftarrow{\sigma} \\ \xrightarrow{\iota} \end{array} & B \\
 & & \alpha \downarrow & \swarrow \beta & \\
 & & E & &
 \end{array}$$

3B. Resolutions. A G -complex (or simply a *complex*) is a sequence of G -modules E^i and G -maps $\delta^i : E^i \rightarrow E^{i+1}$ such that $\delta^{i+1} \circ \delta^i = 0$ for every i , where i runs over $\mathbb{N} \cup \{-1\}$:

$$0 \rightarrow E^{-1} \xrightarrow{\delta^{-1}} E^0 \xrightarrow{\delta^0} E^1 \xrightarrow{\delta^1} \dots \xrightarrow{\delta^n} E^{n+1} \xrightarrow{\delta^{n+1}} \dots$$

Such a sequence will often be denoted by (E^*, δ^*) .

A G -chain map (or simply a *chain map*) between G -complexes (E^*, δ_E^*) and (F^*, δ_F^*) is a sequence of G -maps $\{\alpha^i : E^i \rightarrow F^i \mid i \geq -1\}$ such that $\delta_F^i \circ \alpha^i = \alpha^{i+1} \circ \delta_E^i$ for every $i \geq -1$. If α^* , β^* are chain maps between (E^*, δ_E^*) and (F^*, δ_F^*) which coincide in degree -1 , a G -homotopy between α^* and β^* is a sequence of G -maps $\{T^i : E^i \rightarrow F^{i-1} \mid i \geq 0\}$ such that $\delta_F^{i-1} \circ T^i + T^{i+1} \circ \delta_E^i = \alpha^i - \beta^i$ for every $i \geq 0$, and $T^0 \circ \delta_E^{-1} = 0$. We recall that, according to our definition of G -maps, both chain maps between G -complexes and G -homotopies between such chain maps have to be bounded in every degree.

A complex is *exact* if δ^{-1} is injective and $\ker \delta^{i+1} = \text{Im } \delta^i$ for every $i \geq -1$. A G -resolution (or simply a *resolution*) of a G -module E is an exact G -complex (E^*, δ^*) with $E^{-1} = E$. A resolution (E^*, δ^*) is *relatively injective* if E^n is relatively injective for every $n \geq 0$.

A *contracting homotopy* for a resolution (E^*, δ^*) is a sequence of linear maps $k^i : E^i \rightarrow E^{i-1}$ such that $\|k^i\| \leq 1$ for every $i \in \mathbb{N}$, $\delta^{i-1} \circ k^i + k^{i+1} \circ \delta^i = \text{Id}_{E^i}$ if $i \geq 0$, and $k^0 \circ \delta^{-1} = \text{Id}_E$.

$$0 \longrightarrow E^{-1} \begin{array}{c} \xleftarrow{k^0} \\ \xrightarrow{\delta^{-1}} \end{array} E^0 \begin{array}{c} \xleftarrow{k^1} \\ \xrightarrow{\delta^0} \end{array} E^1 \begin{array}{c} \xleftarrow{k^2} \\ \xrightarrow{\delta^1} \end{array} \dots \begin{array}{c} \xleftarrow{k^n} \\ \xrightarrow{\delta^{n-1}} \end{array} E^n \begin{array}{c} \xleftarrow{k^{n+1}} \\ \xrightarrow{\delta^n} \end{array} \dots$$

Note however that it is not required that k^i be G -equivariant. A resolution is *strong* if it admits a contracting homotopy.

The following result can be proved by means of standard homological algebra arguments (see [Ivanov 1985] and [Monod 2001, Lemmas 7.2.4 and 7.2.6]).

Proposition 3.3. *Let $\alpha : E \rightarrow F$ be a G -map between G -modules, let (E^*, δ_E^*) be a strong resolution of E , and suppose (F^*, δ_F^*) is a G -complex such that $F^{-1} = F$*

and F^i is relatively injective for every $i \geq 0$. Then α extends to a chain map α^* , and any two extensions of α to chain maps are G -homotopic.

3C. Absolute bounded cohomology of groups. If E is a G -module, we denote by $E^G \subseteq E$ the submodule of G -invariant elements in E .

Let (E^*, δ^*) be a relatively injective strong resolution of the trivial G -module \mathbb{R} (such a resolution exists, see Section 3D). Since coboundary maps are G -maps, they restrict to the G -invariant submodules of the E^i 's. Thus $((E^*)^G, \delta^*|)$ is a subcomplex of (E^*, δ^*) . A standard application of Proposition 3.3 now shows that the isomorphism type of the homology of $((E^*)^G, \delta^*|)$ does not depend on the chosen resolution (while the seminorm induced on such homology module by the norms on the E^i 's could depend on it). What is more, there exists a canonical isomorphism between the homology of any two such resolutions, which is induced by any extension of the identity of \mathbb{R} . For every $n \geq 0$, we now define the n -dimensional *bounded cohomology* module $H_b^n(G)$ of G as follows: if $n \geq 1$, then $H_b^n(G)$ is the n -th homology module of the complex $((E^*)^G, \delta^*|)$, while if $n = 0$ then $H_b^n(G) = \ker \delta^0 \cong \mathbb{R}$.

3D. The standard resolution. For every $n \in \mathbb{N}$, let $B^n(G)$ be the space of bounded real maps on G^{n+1} . We endow $B^n(G)$ with the supremum norm and with the diagonal action of G defined by $(g \cdot f)(g_0, \dots, g_n) = f(g^{-1}g_0, \dots, g^{-1}g_n)$, thus defining on $B^n(G)$ a structure of G -module. For $n \geq 0$ we define $\delta^n : B^n(G) \rightarrow B^{n+1}(G)$ by setting:

$$\delta^n(f)(g_0, g_1, \dots, g_{n+1}) = \sum_{i=0}^{n+1} (-1)^i f(g_0, \dots, \widehat{g}_i, \dots, g_{n+1}).$$

Moreover, we let $B^{-1}(G) = \mathbb{R}$ be the trivial G -module, and we define $\delta^{-1} : \mathbb{R} \rightarrow B^0(G)$ by setting $\delta^{-1}(t)(g) = t$ for every $g \in G$. The complex $(B^*(G), \delta^*)$ admits the following contracting homotopy:

$$(1) \quad v^n : B^n(G) \rightarrow B^{n-1}(G), \quad v^n(f)(g_0, \dots, g_{n-1}) = f(e, g_0, \dots, g_{n-1})$$

(for $n = 0$ we understand that $v^0(f) = f(e) \in \mathbb{R} = B^{-1}(G)$ for every $f \in B^0(G)$). Therefore, the complex $(B^*(G), \delta^*)$ provides a strong resolution of the trivial G -module \mathbb{R} , and we will see in Proposition 3.5 below that such a resolution is also relatively injective. In fact, the complex $(B^*(G), \delta^*)$ is usually known as the *standard resolution of the trivial G -module* \mathbb{R} .

Remark 3.4. We briefly compare our notion of standard resolution with Ivanov's and Monod's ones. In [Ivanov 1985], for every $n \in \mathbb{N}$ the set $B^n(G)$ is denoted by $B(G^{n+1})$, and is turned into a Banach G -module by the action $g \cdot f(g_0, \dots, g_n) =$

$f(g_0, \dots, g_n \cdot g)$. Moreover, the sequence of modules $B(G^n)$, $n \in \mathbb{N}$, is equipped with a structure of G -complex

$$0 \rightarrow \mathbb{R} \xrightarrow{d_{-1}} B(G) \xrightarrow{d_0} B(G^2) \xrightarrow{d_1} \dots \xrightarrow{d_n} B(G^{n+2}) \xrightarrow{d_{n+1}} \dots,$$

where $d_{-1}(t)(g) = t$ and

$$\begin{aligned} d_n(f)(g_0, \dots, g_{n+1}) \\ = (-1)^{n+1} f(g_1, \dots, g_{n+1}) + \sum_{i=0}^n (-1)^{n-i} f(g_0, \dots, g_i g_{i+1}, \dots, g_{n+1}) \end{aligned}$$

for every $n \geq 0$ (here we are using Ivanov's notation also for the differential). Now, it is readily seen that Ivanov's resolution is isomorphic to our standard resolution via the isometric G -chain isomorphism $\varphi^* : B^*(G) \rightarrow B(G^{*+1})$ defined by

$$\varphi^n(f)(g_0, \dots, g_n) = f(g_n^{-1}, g_n^{-1} g_{n-1}^{-1}, \dots, g_n^{-1} g_{n-1}^{-1} \cdots g_1^{-1} g_0^{-1});$$

its inverse is given by

$$(\varphi^n)^{-1}(f)(g_0, \dots, g_n) = f(g_n^{-1} g_{n-1}, g_{n-1}^{-1} g_{n-2}, \dots, g_1^{-1} g_0, g_0^{-1}).$$

We observe that our contracting homotopy (1) is conjugated by φ^* into Ivanov's contracting homotopy [1985] for the complex $(B(G^*), d_*)$.

Our notation is much closer to Monod's one. In fact, in [Monod 2001] the more general case of a topological group G is addressed, and the n -th module of the standard G -resolution of \mathbb{R} is inductively defined by setting

$$C_b^0(G, \mathbb{R}) = C_b(G, \mathbb{R}), \quad C_b^n(G, \mathbb{R}) = C_b(G, C_b^{n-1}(G, \mathbb{R})),$$

where $C_b(G, E)$ denotes the space of *continuous* bounded maps from G to the Banach space E . However, as observed in [Monod 2001, Remarks 6.1.2 and 6.1.3], the case when G is an abstract group may be recovered from the general case just by equipping G with the discrete topology. In that case, our notion of standard resolution coincides with Monod's. (See also [Monod 2001, Remark 7.4.9].)

Proposition 3.5 [Ivanov 1985; Monod 2001]. *The standard resolution of \mathbb{R} as a G -module is relatively injective and strong.*

Proof. We have already shown that the standard resolution is strong. The fact that it is also relatively injective is proved in [Monod 2001, Proposition 4.4.1] (see also Remark 7.4.9 of the same reference). Alternatively, since our standard resolution is isometrically isomorphic to Ivanov's one (see Remark 3.4), the relative injectivity of the standard resolution may be deduced from [Ivanov 1985, Lemma 3.2.2]. \square

The seminorm induced on $H_b^*(G)$ by the standard resolution is called the *canonical seminorm*. It is shown in [Ivanov 1985] that the canonical seminorm coincides

with the infimum of all the seminorms induced on $H_b^*(G)$ by any relatively injective strong resolution of the trivial G -module \mathbb{R} (see also Proposition 3.10 below).

3E. Relative bounded cohomology of groups. Let A be a subgroup of G . Henceforth, whenever E is a G -module we understand that E is endowed also with the natural structure of A -module induced by the inclusion of A in G .

Definition 3.6 [Park 2003, Definitions 3.1 and 3.5]. Let (U^*, δ_U^*) be a relatively injective strong G -resolution of the trivial G -module \mathbb{R} and (V^*, δ_V^*) be a relatively injective strong A -resolution of the trivial A -module \mathbb{R} . By Proposition 3.3, the identity of \mathbb{R} may be extended to an A -chain map $\lambda^* : U^* \rightarrow V^*$. The pair of resolutions (U^*, δ_U^*) , (V^*, δ_V^*) , together with the chain map λ^* , provides a *pair of resolutions* for $(G, A; \mathbb{R})$. We say that such a pair is

- (1) *allowable*, if the chain map λ^* commutes with the contracting homotopies of (U^*, δ_U^*) and (V^*, δ_V^*) ;
- (2) *proper*, if the map λ^n restricts to a surjective map $\widehat{\lambda}^n : (U^n)^G \rightarrow (V^n)^A$ for every $n \in \mathbb{N}$.

We denote by $\ker(U^n \rightarrow V^n)$ the kernel of λ^n . It is readily seen that the module $\ker(U^n \rightarrow V^n)^G \subseteq (U^n)^G$ coincides with the kernel of $\widehat{\lambda}^n$.

If the pair of resolutions (U^*, δ_U^*) , (V^*, δ_V^*) is proper, there exists an exact sequence

$$0 \longrightarrow \ker(U^n \rightarrow V^n)^G \longrightarrow (U^n)^G \xrightarrow{\widehat{\lambda}^n} (V^n)^A \longrightarrow 0,$$

which induces the long exact sequence

$$\cdots \longrightarrow H_b^{n-1}(A) \longrightarrow H^n(\ker(U^* \rightarrow V^*)^G) \longrightarrow H_b^n(G) \longrightarrow H_b^n(A) \longrightarrow \cdots$$

As observed in [Park 2003], if the pair (U^*, δ_U^*) , (V^*, δ_V^*) is also allowable, then the isomorphism type of $H^n(\ker(U^* \rightarrow V^*)^G)$ does not depend on the chosen proper allowable pair of resolutions (see also Proposition 3.10 below). Such a module is called the *n -th bounded cohomology group of the pair (G, A)* , and it is denoted by $H_b^n(G, A)$.

3F. The standard pair of resolutions. The following result is proved in [Park 2003, Propositions 3.1 and 3.18], and shows that, just as in the absolute case, there exists a canonical proper allowable pair of resolutions for $(G, A; \mathbb{R})$. Strictly speaking, Park's notion of standard pair of resolutions is different from ours, since it is based on Ivanov's definition of standard resolution. However, the isomorphism described in Remark 3.4 translates Park's results into the following:

Proposition 3.7. *The standard resolutions $B^*(G)$ and $B^*(A)$ of the trivial G - and A -module \mathbb{R} , together with the obvious restriction map $B^*(G) \rightarrow B^*(A)$, provide a proper allowable pair of resolutions for $(G, A; \mathbb{R})$.*

The seminorm induced on $H_b^*(G, A; \mathbb{R})$ by this resolution is called the *canonical seminorm*. In order to save some words, from now on we fix the following notation:

$$B^n(G, A) = \ker(B^n(G) \rightarrow B^n(A)).$$

3G. Morphisms of pairs of resolutions. Let (U^*, δ_U^*) , (V^*, δ_V^*) and (E^*, δ_E^*) , (F^*, δ_F^*) be pairs of resolutions for $(G, A; \mathbb{R})$. A *morphism* between such pairs is a pair of chain maps (α_G^*, α_A^*) such that:

- (1) $\alpha_G^* : U^* \rightarrow E^*$ (resp. $\alpha_A^* : V^* \rightarrow F^*$) is a G -chain map (resp. an A -chain map) extending the identity of $\mathbb{R} = U^{-1} = E^{-1}$ (resp. the identity of $\mathbb{R} = V^{-1} = F^{-1}$);
- (2) for every $n \in \mathbb{N}$, the following diagram commutes

$$\begin{array}{ccc} U^n & \longrightarrow & V^n \\ \downarrow \alpha_G^n & & \downarrow \alpha_A^n \\ E^n & \longrightarrow & F^n, \end{array}$$

where the horizontal rows represent the A -morphisms involved in the definition of a pair of resolutions.

By condition (2), if (α_G^*, α_A^*) is a morphism of pairs of resolutions, then α_G^* restricts to a chain map

$$\alpha_{G,A}^* : \ker(U^* \rightarrow V^*) \rightarrow \ker(E^* \rightarrow F^*),$$

which induces in turn a map

$$H^*(\alpha_{G,A}^*) : H^*(\ker(U^* \rightarrow V^*)^G) \rightarrow H^*(\ker(E^* \rightarrow F^*)^G).$$

Proposition 3.8. *If the pairs of resolutions*

$$(U^*, \delta_U^*), (V^*, \delta_V^*) \quad \text{and} \quad (E^*, \delta_E^*), (F^*, \delta_F^*)$$

are proper, the map $H^(\alpha_{G,A}^*)$ is an isomorphism.*

Proof. Our hypothesis ensures that we have the commutative diagram

$$\begin{array}{ccccccc} \dots & H^{n-1}((V^*)^A) & \longrightarrow & H^n(\ker(U^* \rightarrow V^*)^G) & \longrightarrow & H^n((U^*)^G) & \longrightarrow & H^n((V^*)^A) & \dots \\ & \downarrow H^{n-1}(\alpha_A^*) & & \downarrow H^n(\alpha_{G,A}^*) & & \downarrow H^n(\alpha_G^*) & & \downarrow H^n(\alpha_A^*) & \\ \dots & H^{n-1}((F^*)^A) & \longrightarrow & H^n(\ker(E^* \rightarrow F^*)^G) & \longrightarrow & H^n((E^*)^G) & \longrightarrow & H^n((F^*)^A) & \dots \end{array}$$

The discussion carried out in Section 3C implies that the vertical arrows corresponding to $H^*(\alpha_G^*)$ and $H^*(\alpha_A^*)$ are isomorphisms, so the conclusion follows from the Five Lemma. \square

Remark 3.9. At the moment we are not able to prove either that every two proper allowable pairs of resolutions for $(G, A; \mathbb{R})$ are related by a morphism of pairs of resolutions, or that any two such morphisms induce the same map in cohomology. In fact, whenever two proper allowable pairs of resolutions are given, using Proposition 3.3 one can easily construct the needed chain maps α_G^* and α_A^* . However, some troubles arise in proving that such chain maps can be chosen so to fulfill condition (2) in the above definition of morphism of pairs of resolutions. Despite these difficulties, the results proved in Propositions 3.8 and 3.10 are sufficient to our purposes.

Also observe that in the statement of Proposition 3.8 we do not require the involved pairs of resolutions to be allowable. However, allowability plays a fundamental rôle in constructing a morphism of pairs of resolutions between any generic proper allowable pair of resolutions and the standard pair of resolutions (see Proposition 3.10 below), and in getting explicit bounds on the norm of such a morphism.

The following result shows that, just as in the absolute case, the bounded cohomology of (G, A) is computed by any proper allowable pair of resolutions for $(G, A; \mathbb{R})$. Moreover, the canonical seminorm coincides with the infimum of all the seminorms induced on $H_b^*(G, A)$ by any such pair of resolutions.

Proposition 3.10. *Let (U^*, δ_U^*) , (V^*, δ_V^*) be a proper allowable pair of resolutions for $(G, A; \mathbb{R})$. Then there exists a morphism (α_G^*, α_A^*) between this pair of resolutions and the canonical pair of resolutions introduced in Section 3F. Moreover, one may choose α_G^*, α_A^* in such a way that the induced map*

$$H^*(\alpha_{G,A}^*) : H^*(\ker(U^* \rightarrow V^*)^G) \rightarrow H^*(B^*(G, A)^G) \cong H_b^*(G, A)$$

is a norm-nonincreasing isomorphism.

Proof. Let k_G^* and k_A^* be the contracting homotopies of (U^*, δ_U^*) and (V^*, δ_V^*) , respectively. Define α_G^n and α_A^n by induction as follows:

$$(2) \quad \begin{aligned} \alpha_G^n(f)(g_0, \dots, g_n) &= \alpha_G^{n-1}(g_0(k_G^n g_0^{-1}(f)))(g_1, \dots, g_n) \in \mathbb{R}, \\ \alpha_A^n(f)(g_0, \dots, g_n) &= \alpha_A^{n-1}(g_0(k_A^n g_0^{-1}(f)))(g_1, \dots, g_n) \in \mathbb{R}. \end{aligned}$$

That α_G^* is indeed a G -chain map and α_A^* is an A -chain map is showed in the proof of [Monod 2001, Theorem 7.3.1]. (Alternatively, one may easily check that the maps α_G^* and α_A^* are related to the maps given in [Ivanov 1985, Theorem 3.6]

via the isomorphism described in Remark 3.4.) Moreover, it is clear from the definitions that α_G^* and α_A^* are norm-nonincreasing in every degree.

Since the chain map $U^* \rightarrow V^*$ commutes with the contracting homotopies of (U^*, δ_U^*) and (V^*, δ_V^*) , the following diagram commutes:

$$\begin{array}{ccc} U^n & \longrightarrow & V^n \\ \downarrow \alpha_G^n & & \downarrow \alpha_A^n \\ B^n(G) & \longrightarrow & B^n(A). \end{array}$$

This implies that (α_G^*, α_A^*) is a morphism of pairs of resolutions. Now the conclusion follows from Proposition 3.8. \square

4. Relative (continuous) bounded cohomology of spaces

Throughout the whole section we denote by (X, W) a countable CW-pair. We also make the following:

Standing assumption: Both X and W are connected, and the inclusion of W in X induces an injective map on fundamental groups.

Being locally contractible, the space X admits a universal covering $p : \tilde{X} \rightarrow X$. We denote by \tilde{W} a fixed connected component of $p^{-1}(W) \subseteq \tilde{X}$. We also choose a basepoint $b_0 \in \tilde{W}$. This choice determines a canonical isomorphism between $\pi_1(X, p(b_0))$ and the group G of the covering automorphisms of \tilde{X} . We denote by $A \subseteq G$ the subgroup corresponding to $i_*(\pi_1(W, p(b_0)))$ under this isomorphism, where $i : W \rightarrow X$ is the inclusion. Observe that A coincides with the group of automorphisms of \tilde{X} that leave \tilde{W} invariant. In particular, for every $n \in \mathbb{N}$ the module $C_b^n(\tilde{X})$ (resp. $C_b^n(\tilde{W})$) admits a natural structure of G -module (resp. A -module). Moreover, the covering projection $p : \tilde{X} \rightarrow X$ defines a pull-back map $p^* : C_b^*(X, W) \rightarrow C_b^*(\tilde{X}, \tilde{W})$ which induces in turn an isometric isomorphism $C_b^*(X, W) \rightarrow C_b^*(\tilde{X}, \tilde{W})^G$. As a consequence, we get the natural identification

$$H_b^*(X, W) \cong H^*(C_b^*(\tilde{X}, \tilde{W})^G).$$

The straightening procedure described in Section 2 shows that, when (X, W) is a locally convex pair of metric spaces, in order to compute the relative singular homology of (X, W) one may replace the singular complex $C_*(X, W)$ with the subcomplex of straight chains. As a consequence, it is easily seen that in order to compute the cohomology (resp. the bounded cohomology) of (X, W) one may replace the complex $C^*(\tilde{X}, \tilde{W})^G$ (resp. $C_b^*(\tilde{X}, \tilde{W})^G$) with the subcomplex of those invariant cochains whose value on each simplex only depends on the vertices of the simplex (recall that straight simplices in \tilde{X} only depend on their vertices). Following [Gromov 1982], we say that any such cochain is *straight*.

Observe that the definition of straight cochain makes sense even when it is not possible to properly define a straightening on singular chains. Let us briefly describe some known results about straight cochains in the absolute case (when $W = \emptyset$). If \tilde{X} is contractible, a classical result ensures that both straight cochains and singular cochains compute the cohomology of G , so the cohomology of straight cochains is isomorphic to the singular cohomology of X . An important result in [Gromov 1982, Section 2.3] shows that the same is true for bounded cohomology, even without the assumption that \tilde{X} is contractible. More precisely, both bounded straight cochains and bounded singular cochains compute the bounded cohomology of G , and they both induce the canonical seminorm on $H_b^*(G)$, so the cohomology of bounded straight cochains is isometrically isomorphic to the bounded cohomology of X . Moreover by [Monod 2001, Theorem 7.4.5], the bounded cohomology of G (whence of X) is computed also by *continuous* bounded straight cochains. Monod's result plays a fundamental rôle in Löh's description of the isometric isomorphism between measure homology and singular homology in the absolute case.

In this section we show that, in the case when $W \neq \emptyset$, continuous bounded straight cochains compute the bounded cohomology of the pair (G, A) , thus extending Monod's result to the relative case (see Theorem 4.1).

Moreover, in the case when the pair (X, W) is good we prove that also $H_b^*(X, W)$ is isometrically isomorphic to $H_b^*(G, A)$, thus obtaining that the bounded cohomology of (X, W) is computed by continuous bounded straight cochains. Finally, in Section 4E we show that this result easily implies our Theorem 1.7.

4A. Bounded cochains versus continuous bounded straight cochains. We next give the precise definition of the complex of continuous bounded straight cochains. For every $n \in \mathbb{N}$ we consider the following Banach spaces:

$$\begin{aligned} C_{cbs}^n(\tilde{X}) &= \{f : \tilde{X}^{n+1} \rightarrow \mathbb{R}, f \text{ continuous and bounded}\}, \\ C_{cbs}^n(\tilde{W}) &= \{f : \tilde{W}^{n+1} \rightarrow \mathbb{R}, f \text{ continuous and bounded}\}, \end{aligned}$$

both endowed with the supremum norm. The diagonal G -action such that $g \cdot f(x_0, \dots, x_n) = f(g^{-1}x_0, \dots, g^{-1}x_n)$ for every $g \in G$ endows $C_{cbs}^n(\tilde{X})$ with a structure of G -module. The obvious coboundary maps $\delta^n : C_{cbs}^n(\tilde{X}) \rightarrow C_{cbs}^{n+1}(\tilde{X})$ given by

$$\delta^n(f)(x_0, \dots, x_{n+1}) = \sum_{i=0}^{n+1} (-1)^i f(x_0, \dots, \hat{x}_i, \dots, x_{n+1})$$

define on $C_{cbs}^*(\tilde{X})$ a structure of G -complex. In the very same way one endows $C_{cbs}^*(\tilde{W})$ with a structure of A -complex. For every $n \in \mathbb{N}$, the inclusion $\tilde{W}^{n+1} \hookrightarrow \tilde{X}^{n+1}$ induces an obvious restriction $C_{cbs}^n(\tilde{X}) \rightarrow C_{cbs}^n(\tilde{W})$, whose kernel will be

denoted by $C_{cbs}^n(\tilde{X}, \tilde{W})$. Finally, for every $n \in \mathbb{N}$ we set

$$(3) \quad H_{cbs}^n(X, W) = H^n(C_{cbs}^*(\tilde{X}, \tilde{W})^G).$$

We will prove in Propositions 4.3 and 4.7 that both $C_b^*(\tilde{X})$, $C_b^*(\tilde{W})$ and $C_{cbs}^*(\tilde{X})$, $C_{cbs}^*(\tilde{W})$ provide proper pairs of resolutions for $(G, A; \mathbb{R})$. The pair of norm-nonincreasing chain maps

$$(4) \quad \begin{aligned} \eta_G^* : C_{cbs}^*(\tilde{X}) &\rightarrow C_b^*(\tilde{X}), & \eta_G^n(f)(\sigma) &= f(\sigma(e_0), \dots, \sigma(e_n)), \\ \eta_A^* : C_{cbs}^*(\tilde{W}) &\rightarrow C_b^*(\tilde{W}), & \eta_A^n(f)(\sigma) &= f(\sigma(e_0), \dots, \sigma(e_n)) \end{aligned}$$

allows us to identify $C_{cbs}^*(\tilde{X})$ with the subcomplex of $C_b^*(\tilde{X})$ of continuous bounded straight cochains on \tilde{X} , and likewise with \tilde{W} in place of \tilde{X} . Moreover, it is readily seen that the pair (η_G^*, η_A^*) is a morphism of resolutions. Therefore, Proposition 3.8 implies that the induced map in cohomology

$$H^*(\eta_{G,A}^*) : H_{cbs}^*(X, W) = H^*(C_{cbs}^*(\tilde{X}, \tilde{W})^G) \rightarrow H^*(C_b^*(\tilde{X}, \tilde{W})^G) = H_b^*(X, W)$$

is an isomorphism. Moreover, the explicit description of $\eta_{G,A}^*$ shows that $H^*(\eta_{G,A}^*)$ is norm-nonincreasing.

Under the assumption that the pair (X, W) is good, the isomorphism $H^*(\eta_{G,A}^*)$ is in fact an isometry. This fact is proved in the following subsections, and will play a fundamental rôle in our proof of Theorem 1.7.

We now describe briefly the content of the following subsections. In Section 4B we define a morphism of resolutions (β_G^*, β_A^*) between the standard pair of resolutions and continuous bounded straight cochains via an *ad hoc* construction, and we show that this morphism induces an isometric isomorphism in cohomology. Then, under the assumption that (X, W) is good, we prove in Proposition 4.7 that bounded cochains provide a proper allowable pair of resolutions for $(G, A; \mathbb{R})$, so we may exploit Proposition 3.10 to construct a morphism of pairs of resolutions (α_G^*, α_A^*) between bounded cochains and the standard pair of resolutions for $(G, A; \mathbb{R})$. This morphism induces a norm-nonincreasing isomorphism in cohomology, so in order to prove that the isomorphism $H^*(\eta_{G,A}^*)$ is isometric we will be left to show that the composition $\beta_{G,A}^* \circ \alpha_{G,A}^*$ induces the inverse of $H^*(\eta_{G,A}^*)$ in cohomology; in other words, that the following diagram commutes:

$$\begin{array}{ccc} & H_b^*(G, A) & \\ H^*(\beta_{G,A}^*) \swarrow & & \nwarrow H^*(\alpha_{G,A}^*) \\ H_{cbs}^*(X, W) & \xrightarrow{H^*(\eta_{G,A}^*)} & H_b^*(X, W). \end{array}$$

We can summarize the results just described in the following theorem, whose proof is carried out in Subsections 4B, 4C, 4D.

Theorem 4.1. *For every $n \in \mathbb{N}$ the map*

$$H^n(\beta_{G,A}^*) : H_b^n(G, A) \rightarrow H_{cbs}^n(X, W)$$

is an isometric isomorphism, and the map

$$H^n(\eta_{G,A}^*) : H_{cbs}^n(X, W) \rightarrow H_b^n(X, W)$$

is a norm-nonincreasing isomorphism. In particular, the composition

$$H^n(\eta_{G,A}^*) \circ H^n(\beta_{G,A}^*)$$

is a norm-nonincreasing isomorphism between $H_b^n(G, A)$ and $H_b^n(X, W)$. If, in addition, (X, W) is good, then $H^n(\eta_{G,A}^)$ is an isometry, and $H_b^n(G, A)$ and $H_b^n(X, W)$ are isometrically isomorphic.*

In fact, one may notice that the proof that $H^n(\beta_{G,A}^*)$ is an isometric isomorphism still works without the assumption that X and W are countable.

4B. Mapping standard resolutions into continuous bounded straight cochains.

We begin with a generalization of [Frigerio 2011, Lemma 5.1]:

Lemma 4.2. *There exists a continuous map $\chi : \tilde{X} \rightarrow [0, 1]$ with the following properties:*

- (1) *For every $x \in \tilde{X}$ there exists a neighborhood U_x of $x \in \tilde{X}$ such that the set $\{g \in G \mid \text{supp}(\chi) \cap g(U_x) \neq \emptyset\}$ is finite.*
- (2) *For every $x \in \tilde{X}$, we have $\sum_{g \in G} \chi(g \cdot x) = 1$. (Note that the sum on the left-hand side is finite by (1).)*
- (3) *For every $w \in \tilde{W}$ and every $g \in G \setminus A$, we have $\chi(g \cdot w) = 0$, whence $\sum_{g \in A} \chi(g \cdot w) = 1$.*
- (4) *We have $\chi(b_0) = 1$, so $\chi(g \cdot b_0) = 0$ for every $g \neq 1$.*

Proof. Recall that $p : \tilde{X} \rightarrow X$ is the universal covering of X . Using that W is a subcomplex of X , one can easily construct an open covering $\mathcal{U} = \{U_i\}_{i \in I}$ of X such that every U_i is contractible (whence evenly covered with respect to $p : \tilde{X} \rightarrow X$) and $U_i \cap W$ is path-connected for every $i \in I$ (for example, if $\epsilon > 0$ is small enough and $x \in X$, the contractible ϵ -neighborhood $N_\epsilon(x)$ of x constructed in [Hatcher 2002, page 522] intersects any subcomplex of X in a contractible, whence path-connected, subset). Now choose $i_0 \in I$ such that $p(b_0) \in U_{i_0}$, and set $J = \{i \in I \mid U_i \cap W \neq \emptyset\}$ (so $i_0 \in J$).

For every U_i we choose an open subset $H_i \subseteq \tilde{X}$ in such a way that the following conditions hold:

- (a) $p|_{H_i} : H_i \rightarrow U_i$ is a homeomorphism.
- (b) $p^{-1}(U_i) = \bigcup_{g \in G} g(H_i)$ and $g(H_i) \cap g'(H_i) = \emptyset$ for every $g \neq g'$.

(c) $b_0 \in H_{i_0}$.

(d) $H_i \cap \tilde{W} \neq \emptyset$ for every $i \in J$.

We now set $U'_i = U_i \setminus \{p(b_0)\}$ for every $i \neq i_0$, $U'_{i_0} = U_{i_0}$, and $\mathcal{U}' = \{U'_i\}_{i \in I}$. Let also $H'_i = H_i \cap p^{-1}(U'_i)$. Since $U_i \cap W$ is path-connected, condition (d) easily implies that

$$H_i \cap p^{-1}(W) = H_i \cap \tilde{W} \quad \text{for every } i \in I,$$

whence

$$(5) \quad H'_i \cap p^{-1}(W) = H'_i \cap \tilde{W} \quad \text{for every } i \in I.$$

Since every CW-complex is paracompact (see [Miyazaki 1952; Bourgin 1952], for instance), we may now take a partition of unity $\{\varphi_i\}_{i \in I}$ adapted to \mathcal{U}' , and let $\psi_i : \tilde{X} \rightarrow \mathbb{R}$ be the map which coincides with $\varphi_i \circ p$ on H'_i and is null outside H'_i . We finally set

$$\chi = \sum_{i \in I} \psi_i.$$

The fact that χ satisfies properties (1) and (2) of the statement is proved in [Frigerio 2011, Lemma 5.1]. Moreover, for every $w \in \tilde{W}$ and $g \in G \setminus A$ we have $g \cdot w \in p^{-1}(W) \setminus \tilde{W}$, so Equation (5) implies that $g \cdot w$ does not belong to any H'_i . This implies point (3). Finally, since $p(b_0) \notin U'_i$ for every $i \neq i_0$, we have necessarily $\varphi_i(p(b_0)) = 0$ for every $i \neq i_0$, and $\varphi_{i_0}(p(b_0)) = 1$. By (c) this implies that $\psi_{i_0}(b_0) = 1$, whence $\chi(b_0) = 1$, as desired. \square

Proposition 4.3. *The pair $(C_{cbs}^*(\tilde{X}), \delta^*)$, $(C_{cbs}^*(\tilde{W}), \delta^*)$ provides a proper allowable pair of resolutions for $(G, A; \mathbb{R})$.*

Proof. The fact that $(C_{cbs}^*(\tilde{X}), \delta^*)$ (resp. $(C_{cbs}^*(\tilde{W}), \delta^*)$) provides a relatively injective resolution of \mathbb{R} as a trivial G -module (resp. A -module) is proved in [Monod 2001, Theorem 7.4.5]. (To apply that result our CW-complexes X and W should be locally compact, whence locally finite; but these conditions are used in Monod's proof only to ensure the existence of a suitable *Bruhat function* on \tilde{X} and on \tilde{W} ; in our case of interest the fact that G and A are discrete allows us to explicitly describe such a map; see Lemma 4.2.)

It is readily seen that these resolutions admit the contracting homotopies

$$(6) \quad \begin{aligned} t_G^n(f)(x_1, \dots, x_n) &= f(b_0, x_1, \dots, x_n), \quad f \in C_{cbs}^n(\tilde{X}), \quad (x_1, \dots, x_n) \in \tilde{X}^n, \\ t_A^n(f)(w_1, \dots, w_n) &= f(b_0, w_1, \dots, w_n), \quad f \in C_{cbs}^n(\tilde{W}), \quad (w_1, \dots, w_n) \in \tilde{W}^n. \end{aligned}$$

This clearly implies that the A -chain map $\gamma^* : C_{cbs}^*(\tilde{X}) \rightarrow C_{cbs}^*(\tilde{W})$ induced by the inclusion $\tilde{W} \hookrightarrow \tilde{X}$ commutes with the contracting homotopies.

In order to conclude we have to show that γ^* restricts to a surjective map

$$\widehat{\gamma}^* : C_{cbs}^*(\widetilde{X})^G \rightarrow C_{cbs}^*(\widetilde{W})^A.$$

Let $f : \widetilde{W}^{n+1} \rightarrow \mathbb{R}$ be an A -invariant bounded continuous map. The inclusion $\widetilde{W}^{n+1} \hookrightarrow \widetilde{X}^{n+1}$ induces a homeomorphism ψ between \widetilde{W}^{n+1}/A and a closed subset K of \widetilde{X}^{n+1}/G (recall that W is a CW-subcomplex of X , so it is closed in X). Therefore, f defines a bounded continuous map \overline{f} on K , and by Tietze's theorem we may extend \overline{f} to a bounded continuous map $\overline{g} : \widetilde{X}^{n+1}/G \rightarrow \mathbb{R}$. If g is obtained by precomposing \overline{g} with the projection $\widetilde{X}^{n+1} \rightarrow \widetilde{X}^{n+1}/G$, then $g \in C_{cbs}^n(\widetilde{X})^G$, and $\widehat{\gamma}^n(g) = f$. We have thus shown that $\widehat{\gamma}^*$ is surjective, and this concludes the proof. \square

We are now ready to describe a morphism of pairs of resolutions (β_G^*, β_A^*) between the standard pair of resolutions for $(G, A; \mathbb{R})$ and the complexes of straight cochains. Let

$$\beta_G^n : B^n(G) \rightarrow C_{cbs}^n(\widetilde{X}), \quad \beta_A^n : B^n(A) \rightarrow C_{cbs}^n(\widetilde{W})$$

be defined as follows:

$$\begin{aligned} \beta_G^n(f)(x_0, \dots, x_n) &= \sum_{(g_0, \dots, g_n) \in G^{n+1}} \chi(g_0^{-1}x_0) \cdots \chi(g_n^{-1}x_n) \cdot f(g_0, \dots, g_n), \\ \beta_A^n(f)(w_0, \dots, w_n) &= \sum_{(g_0, \dots, g_n) \in A^{n+1}} \chi(g_0^{-1}w_0) \cdots \chi(g_n^{-1}w_n) \cdot f(g_0, \dots, g_n). \end{aligned}$$

Lemma 4.4. *For every $f \in B^n(G)$, $(g_0, \dots, g_n) \in G^{n+1}$ we have*

$$\beta_G^n(f)(g_0b_0, \dots, g_nb_0) = f(g_0, \dots, g_n).$$

Proof. By Lemma 4.2(4), for every $(\gamma_0, \dots, \gamma_n) \in G^{n+1}$ we have

$$\begin{aligned} \chi(\gamma_0^{-1}g_0b_0) \cdots \chi(\gamma_n^{-1}g_nb_0) \cdot f(\gamma_0, \dots, \gamma_n) \\ = \begin{cases} f(g_0, \dots, g_n) & \text{if } \gamma_i = g_i \text{ for every } i, \\ 0 & \text{otherwise,} \end{cases} \end{aligned}$$

and this readily implies the conclusion. \square

Proposition 4.5. *The pair (β_G^*, β_A^*) provides a well-defined morphism of pairs of resolutions. For every $n \in \mathbb{N}$ the induced map*

$$H^n(\beta_{G,A}^*) : H_b^n(G, A) \rightarrow H_{cbs}^n(X, W)$$

is an isometric isomorphism.

Proof. We begin by showing that β_G^* is a G -map. So, take $f \in B^n(G)$, $g \in G$, and $(x_0, \dots, x_n) \in \tilde{X}^{n+1}$. By definition we have

$$\begin{aligned}\beta_G^n(g \cdot f)(x_0, \dots, x_n) &= \sum_{(g_0, \dots, g_n) \in G^{n+1}} \chi(g_0^{-1}x_0) \cdots \chi(g_n^{-1}x_n) \cdot f(g^{-1}g_0, \dots, g^{-1}g_n), \\ (g \cdot \beta_G^n(f))(x_0, \dots, x_n) &= \sum_{(g_0, \dots, g_n) \in G^{n+1}} \chi(g_0^{-1}g^{-1}x_0) \cdots \chi(g_n^{-1}g^{-1}x_n) \cdot f(g_0, \dots, g_n),\end{aligned}$$

and an easy change of variables implies that β_G^n is a G -map. A similar argument shows that β_A^n is an A -map. We now check that β_G^* is a chain map. By Lemma 4.2(2), for every $x_i \in \tilde{X}$ we have $\sum_{g \in G} \chi(g^{-1}x_i) = 1$, so if $(g_0, \dots, g_{n+1}) \in G^{n+2}$ and $(x_0, \dots, x_{n+1}) \in \tilde{X}^{n+2}$ are fixed, then

$$\begin{aligned}\chi(g_0^{-1}x_0) \cdots \chi(\widehat{g_i^{-1}x_i}) \cdots \chi(g_{n+1}^{-1}x_{n+1}) \\ = \sum_{g \in G} \chi(g_0^{-1}x_0) \cdots \chi(g^{-1}x_i) \cdots \chi(g_{n+1}^{-1}x_{n+1})\end{aligned}$$

and $\beta_G^n(f)(x_0, \dots, \widehat{x_i}, \dots, x_{n+1})$ is equal to

$$\sum_{(g_0, \dots, \widehat{g_i}, \dots, g_{n+1}) \in G^{n+1}} \chi(g_0^{-1}x_0) \cdots \chi(\widehat{g_i^{-1}x_i}) \cdots \chi(g_{n+1}^{-1}x_{n+1}) \cdot f(g_0, \dots, \widehat{g_i}, \dots, g_{n+1}),$$

which in turn equals

$$\sum_{(g_0, \dots, g_i, \dots, g_{n+1}) \in G^{n+2}} \chi(g_0^{-1}x_0) \cdots \chi(g_i^{-1}x_i) \cdots \chi(g_{n+1}^{-1}x_{n+1}) \cdot f(g_0, \dots, \widehat{g_i}, \dots, g_{n+1}).$$

From this equality it is easy to deduce that $\delta^n(\beta_G^n(f)) = \beta_G^{n+1}(\delta^n(f))$, and this proves that β_G^* is a chain map. Since χ has been chosen in such a way that Lemma 4.2(3) holds, the same argument may be exploited to show that β_A^* is also a chain map.

Using again Lemma 4.2(3), it is easily checked that the restriction $\beta_G^n(f)|_{\tilde{W}^{n+1}}$ coincides with the map $\beta_A^n(f|_{A^{n+1}})$ for every $f \in B^n(G)$. As a consequence, the pair (β_G^*, β_A^*) is a morphism of pairs of resolutions, and Proposition 3.8 implies that $H^*(\beta_{G,A}^*)$ is an isomorphism. Moreover, $H^n(\beta_{G,A}^*)$ is obviously norm-non-increasing for every $n \in \mathbb{N}$.

Recall now that Proposition 3.10 provides a morphism of pairs of resolutions

$$\zeta_G^* : C_{cbs}^*(\tilde{X}) \rightarrow B^*(G), \quad \zeta_A^* : C_{cbs}^*(\tilde{W}) \rightarrow B^*(A),$$

which induces a norm-nonincreasing isomorphism

$$H^*(\zeta_{G,A}^*) : H_{cbs}^*(X, W) \rightarrow H_b^*(G, A).$$

In order to conclude it is sufficient to show that for every $n \in \mathbb{N}$ the composition $\zeta_G^n \circ \beta_G^n$ is the identity of $B^n(G)$.

The proof of Proposition 3.10 implies that the map ζ_G^n can be described by the following inductive formula:

$$\zeta_G^n(f)(g_0, \dots, g_n) = \zeta_G^{n-1}(g_0(t_G^n(g_0^{-1}(f))))(g_1, \dots, g_n),$$

where t_G^* is the contracting homotopy for the resolution $C_{cbs}^*(\tilde{X})$ described in Equation (6). As a consequence, an easy induction shows that $\zeta_G^n(f)(g_0, \dots, g_n) = f(g_0b_0, \dots, g_nb_0)$ for every $f \in C_{cbs}^n(\tilde{X})$, $(g_0, \dots, g_n) \in G^{n+1}$. By Lemma 4.4, this implies that $\zeta_G^n \circ \beta_G^n$ is the identity of $B^n(G)$, whence the conclusion. \square

4C. Ivanov’s contracting homotopy. In order to show that, under the hypothesis that (X, W) is good, bounded cochains provide a proper allowable pair of resolutions for $(G, A; \mathbb{R})$, we first recall Ivanov’s construction of a contracting homotopy for the resolution $C_b^*(\tilde{X})$.

It is shown in [Ivanov 1985] that one can construct an infinite tower of bundles

$$(7) \quad \dots \xrightarrow{p_m} X_m \xrightarrow{p_{m-1}} X_{m-1} \xrightarrow{p_{m-2}} \dots \xrightarrow{p_2} X_2 \xrightarrow{p_1} X_1,$$

where $X_1 = \tilde{X}$, $\pi_i(X_m) = 0$ for every $i \leq m$, $\pi_i(X_m) = \pi_i(X)$ for every $i > m$ and each map $p_m : X_{m+1} \rightarrow X_m$ is a principal H_m -bundle for some topological connected abelian group H_m , which has the homotopy type of a $K(\pi_{m+1}(X), m)$. Moreover, the induced chain maps $p_m^* : C_b^*(X_m) \rightarrow C_b^*(X_{m+1})$ admit left inverse chain maps $A_m^* : C_b^*(X_{m+1}) \rightarrow C_b^*(X_m)$ obtained by averaging cochains over the preimages in X_{m+1} of simplices in X_m , in such a way that the A_m ’s are norm-nonincreasing.

Denote by $W_m \subseteq X_m$ the preimage $p_{m-1}^{-1}(p_{m-2}^{-1}(\dots(p_1^{-1}(\tilde{W})))) \subseteq X_m$ (so W_{m+1} is a principal H_m -bundle over W_m for every $m \geq 1$). We denote simply by

$$p_m : W_{m+1} \rightarrow W_m$$

the restriction of p_m to W_{m+1} . It follows from Ivanov’s construction that each A_m^* induces a norm-nonincreasing chain map $C_b^*(W_{m+1}) \rightarrow C_b^*(W_m)$, which will still be denoted by A_m^* .

Lemma 4.6. *Suppose that (X, W) is good. Then $\pi_i(W_m) = 0$ for every $i \leq m$.*

Proof. Of course, it is sufficient to prove that $\pi_i(W_m) \cong \pi_i(X_m)$ for every $i \in \mathbb{N}$, $m \in \mathbb{N}$. Let us prove this last statement by induction on m . Since the inclusion map $W \hookrightarrow X$ is π_1 -injective we have $\pi_1(W_1) = \pi_1(X_1) = 0$. Therefore, since coverings induce isomorphisms on homotopy groups of order at least two, the case $m = 1$ follows from the fact that the pair (X, W) is good. The inductive step follows from an easy application of the Five Lemma to the following commutative diagram,

which descends in turn from the naturality of the homotopy exact sequences for the bundles $X_{m+1} \rightarrow X_m$, $W_{m+1} \rightarrow W_m$:

$$\begin{array}{ccccccccc} \pi_{i+1}(W_m) & \longrightarrow & \pi_i(H_m) & \longrightarrow & \pi_i(W_{m+1}) & \longrightarrow & \pi_i(W_m) & \longrightarrow & \pi_{i-1}(H_m) \\ \downarrow & & \parallel & & \downarrow & & \downarrow & & \parallel \\ \pi_{i+1}(X_m) & \longrightarrow & \pi_i(H_m) & \longrightarrow & \pi_i(X_{m+1}) & \longrightarrow & \pi_i(X_m) & \longrightarrow & \pi_{i-1}(H_m). \quad \square \end{array}$$

Now suppose that (X, W) is good. We choose basepoints $w_m \in W_m$ in such a way that $p_m(w_{m+1}) = w_m$ for every $m \geq 1$, and $w_1 \in W_1 = \tilde{W}$ coincides with the basepoint b_0 fixed above. Since X_m is m -connected, for every $n \leq m$ it is possible to construct a map $L_n^m : S_n(X_m) \rightarrow S_{n+1}(X_m)$ that associates to every $\sigma \in S_n(X_m)$ a cone of σ over w_m (see [Ivanov 1985]). We stress that, since W_m is also m -connected, if $\sigma \in S_n(W_m) \subseteq S_n(X_m)$, then $L_n^m(\sigma)$ can be chosen to belong to $S_{n+1}(W_m)$. The maps L_n^m , $n \leq m$, induce a (partial) homotopy between the identity and the null map of $C_*(X_m)$, which in turn induces a (partial) contracting homotopy $\{k_m^n\}_{n \leq m}$ for the (partial) complex $\{C_b^n(X_m)\}_{n \leq m}$. Since $L_n^m(S_n(W_m)) \subseteq S_{n+1}(W_m)$, this contracting homotopy induces a (partial) contracting homotopy for $\{C_b^n(W_m)\}_{n \leq m}$, which we still denote by k_m^n . Moreover, it is possible to choose these contracting homotopies in a compatible way, in the sense that the equality $A_m^{n-1} \circ k_{m+1}^n \circ p_m^n = k_m^n$ holds for every $n \leq m$ (see again [Ivanov 1985]). Thanks to this compatibility condition, one can finally define the contracting homotopy

$$k_G^* : C_b^*(\tilde{X}) \rightarrow C_b^{*-1}(\tilde{X}),$$

via the formula

$$k_G^n = A_1^{n-1} \circ \cdots \circ A_{m-1}^{n-1} \circ k_m^n \circ p_{m-1}^n \circ \cdots \circ p_2^n \circ p_1^n \quad \text{for any } m \geq n.$$

The very same formula defines a contracting homotopy for $C_b^*(\tilde{W})$. By construction, the restriction map $C_b^*(\tilde{X}) \rightarrow C_b^*(\tilde{W})$ commutes with these contracting homotopies, and it obviously restricts to a surjective map $C_b^*(\tilde{X})^G \rightarrow C_b^*(\tilde{W})^A$. Since $C_b^n(\tilde{X})$, $C_b^n(\tilde{W})$ are relatively injective for every $n \geq 0$ (see [Ivanov 1985]), we have finally proved the following:

Proposition 4.7. *The pair $(C_b^*(\tilde{X}), \delta^*)$, $(C_b^*(\tilde{W}), \delta^*)$ provides a proper pair of resolutions for $(G, A; \mathbb{R})$. If in addition (X, W) is good, then this pair of resolutions is also allowable.*

Corollary 4.8. *For every $n \in \mathbb{N}$, the map*

$$H^n(\eta_{G,A}^*) : H_{cbs}^n(X, W) \rightarrow H_b^n(X, W)$$

is a norm-nonincreasing isomorphism.

Proof. By Proposition 4.7, bounded cochains provide a proper pair of resolutions for $(G, A; \mathbb{R})$, so Proposition 3.8 implies that $H^n(\eta_{G,A}^*)$ is an isomorphism. That it is norm-nonincreasing is a direct consequence of its explicit description. \square

Remark 4.9. The fact that the pair of resolutions $(C_b^*(\tilde{X}), \delta^*), (C_b^*(\tilde{W}), \delta^*)$ is allowable is stated in [Park 2003, Lemma 4.2] under the only assumption that (X, W) is a pair of connected CW-pairs. However, at the moment we are not able to prove such a statement without the assumption that (X, W) is good. For example, let us suppose that X is simply connected and W is a point (so that $\pi_n(W)$ injects into $\pi_n(X)$ for every $n \in \mathbb{N}$, and $X_1 = \tilde{X} = X, W_1 = \tilde{W} = W$). Then for every $n \in \mathbb{N}$ there exists only one simplex in $S_n(W)$, namely the constant n -simplex σ_n^W . Therefore, the only possible contracting homotopy for W is given by the map which sends the cochain $\varphi \in C_b^n(W)$ to the cochain $k_A^n(\varphi)$ such that $k_A^n(\varphi)(\sigma_{n-1}^W) = \varphi(\sigma_n^W)$. On the other hand, it is not difficult to show that $\pi_i(W_m) = \pi_{i+1}(X)$ for every $i < m$, and $\pi_i(W_m) = 0$ for every $i \geq m$. Therefore, if $\pi_{i+1}(X) \neq 0$, then $\pi_i(W_m) \neq 0$ for every $m > i$. This readily implies that for $m > i$ one cannot construct *cone-like* operators $L_j^m : C_j(X_m) \rightarrow C_{j+1}(X_m), j \leq i$, such that $d_{j+1}L_j^m + L_{j-1}^m d_j = \text{Id}$ and $L_j^m(C_j(W_m)) \subseteq C_{j+1}(W_m)$ for every $j \leq i$, so it is not clear how to show that the pair of resolutions $C_b^*(\tilde{X}), C_b^*(\tilde{W})$ is allowable. This difficulty already arises for the pair (S^2, q) , where q is any point of the 2-dimensional sphere S^2 .

Some troubles arise also in the case when the inclusion induces surjective (but not bijective) maps between the homotopy groups of W and of X . For instance, if X is the Euclidean 3-space and $W = S^2$, then $X_m = X$ for every $m \in \mathbb{N}$, so $W_m = W$ for every $m \in \mathbb{N}$, and, if i is sufficiently high, the partial complex $\{C_j(X, W)\}_{j \leq i}$ does not support a relative cone-like operator. Also observe that, if $\{W'_m, m \in \mathbb{N}\}$ is the tower of bundles constructed starting from W just as X_m is constructed starting from X , then the only map $W'_m \rightarrow W_m = S^2 \subseteq \mathbb{R}^3 = X_m$ which commutes with the projections of W'_m and X_m onto $W_1 = S^2$ and $X_1 = \mathbb{R}^3$ is the projection $W'_m \rightarrow W_1 = S^2$. As a consequence, also in this case it is not clear why the pair of resolutions $C_b^*(\tilde{X}), C_b^*(\tilde{W})$ should be allowable.

4D. Proof of Theorem 4.1. We now come back to the proof of Theorem 4.1. By Proposition 4.5 and Corollary 4.8, we are only left to show that, under the assumption that (X, W) is good, the isomorphism

$$H^n(\eta_G^*) : H_{cbs}^n(X, W) \rightarrow H_b(X, W)$$

is isometric for every $n \in \mathbb{N}$.

So, suppose that (X, W) is good. By Proposition 4.7 bounded cochains provide a proper allowable pair of resolutions for $(G, A; \mathbb{R})$. Therefore, Proposition 3.10 provides a morphism of pairs of resolutions

$$\alpha_G^* : C_b^*(\tilde{X}) \rightarrow B^*(G), \quad \alpha_A^* : C_b^*(\tilde{W}) \rightarrow B^*(A),$$

such that the induced map $H^*(\alpha_{G,A}^*)$ is a norm-nonincreasing isomorphism.

We already know that all the maps in the diagram

$$\begin{array}{ccc}
 & H_b^*(G, A) & \\
 H^*(\beta_{G,A}^*) \swarrow & & \nwarrow H^*(\alpha_{G,A}^*) \\
 H_{cbs}^*(X, W) & \xrightarrow{H^*(\eta_{G,A}^*)} & H_b^*(X, W).
 \end{array}$$

are norm-nonincreasing isomorphisms, so in order to conclude it is sufficient to show that the diagram commutes. This fact is obviously implied by the following result, which concludes the proof of Theorem 4.1.

Proposition 4.10. *Suppose that (X, W) is good. Then, for every $n \in \mathbb{N}$ the composition*

$$\alpha_{G,A}^n \circ \eta_{G,A}^n \circ \beta_{G,A}^n : B^n(G, A) \rightarrow B^n(G, A)$$

is equal to the identity of $B^n(G, A)$.

Proof. Since the composition $\alpha_{G,A}^n \circ \eta_{G,A}^n \circ \beta_{G,A}^n$ coincides with the restriction of $\alpha_G^n \circ \eta_G^n \circ \beta_G^n$ to $B^n(G, A) \subseteq B^n(G)$, it is sufficient to show that $\alpha_G^n \circ \eta_G^n \circ \beta_G^n$ is the identity of $B^n(G)$.

Before going into the needed computations, let us stress that the definition of α_G^* involves the contracting homotopy for the resolution $C_b^*(\tilde{X})$ described in Section 4C. Being based on a non-explicit averaging procedure, this contracting homotopy cannot be described by an explicit formula, and the same is true for the chain map α_G^* . However, the explicit description of the composition $\alpha_G^* \circ \eta_G^*$ is sufficient to our purposes.

In fact, we already know from Lemma 4.4 that

$$\beta_G^n(f)(g_0 b_0, \dots, g_n b_0) = f(g_0, \dots, g_n)$$

for every $f \in B^n(G)$, $(g_0, \dots, g_n) \in G^{n+1}$. Therefore, in order to conclude it is sufficient to prove that

$$(8) \quad \alpha_G^n(\eta_G^n(f))(g_0, \dots, g_n) = f(g_0 b_0, \dots, g_n b_0)$$

for every $f \in C_{cbs}^n(\tilde{X})$. So, let t_G^* and k_G^* be the contracting homotopies for continuous bounded straight cochains and for bounded cochains, respectively; see (6) and (7). We first show that for every $n \in \mathbb{N}$ we have

$$(9) \quad k_G^n \circ \eta_G^n = \eta_G^{n-1} \circ t_G^n.$$

Fix $f \in C_{cbs}^n(\tilde{X})$ and $\sigma \in S_{n-1}(\tilde{X})$, and let us compute $k_G^n(\eta_G^n(f))(\sigma)$. With notation as in Section 4C, we choose $m \geq n$ and set

$$f_m = p_{m-1}^n(\dots p_1^n(\eta_G^n(f))) \in C_b^n(X_m).$$

Then, if σ_m is any lift of σ in X_m , we have $k_m^n(f_m)(\sigma_m) = f_m(\sigma'_m)$, where $\sigma'_m \in S_n(X_m)$ has vertices $w_m, \sigma_m(e_0), \dots, \sigma_m(e_{n-1})$. It readily follows that

$$k_m^n(f_m)(\sigma_m) = f(b_0, \sigma(e_0), \dots, \sigma(e_{n-1})).$$

We have thus shown that the cochain $k_m^n(f_m)$ is constant on all the lifts of σ in X_m . By definition, the value of $k_G^n(\eta_G^n(f))(\sigma)$ is obtained by suitably averaging the values taken by $k_m^n(f_m)$ on such lifts, so we finally get

$$k_G^n(\eta_G^n(f))(\sigma) = f(b_0, \sigma(e_0), \dots, \sigma(e_{n-1})),$$

whence (9).

Recall now that the map α_G^* is explicitly described (in terms of the contracting homotopy k_G^*) in Proposition 3.10; see (2). Therefore, (2) and (9) readily imply that the composition $\alpha_G^n \circ \eta_G^n$ can be described by the following inductive formula:

$$\alpha_G^n(\eta_G^n(f))(g_0, \dots, g_n) = \alpha_G^{n-1}(g_0(\eta_G^{n-1}(t_G^n(g_0^{-1}(f)))))(g_1, \dots, g_n).$$

An easy induction now implies (8), whence the conclusion. \square

4E. Proof of Theorem 1.7. We next describe how Theorem 1.7 can be deduced from Theorem 4.1. For every $n \in \mathbb{N}$ the module $C_{cb}^n(\tilde{X})$ (resp. $C_{cb}^n(\tilde{W})$) admits a natural structure of G -module (resp. A -module). Moreover, it is proved in [Frigerio 2011, Lemma 6.1] that the isometric isomorphism $C_b^*(X, W) \rightarrow C_b^*(\tilde{X}, \tilde{W})^G$ induced by the covering projection $p: \tilde{X} \rightarrow X$ restricts to an isometric isomorphism $C_{cb}^*(X, W) \rightarrow C_{cb}^*(\tilde{X}, \tilde{W})^G$, which induces in turn a natural identification

$$(10) \quad H_{cb}^*(X, W) \cong H^*(C_{cb}^*(\tilde{X}, \tilde{W})^G).$$

The G -chain map $v_G^*: C_{cbs}^*(\tilde{X}) \rightarrow C_{cb}^*(\tilde{X})$ defined by

$$v_G^n(f)(\sigma) = f(\sigma(e_0), \dots, \sigma(e_n)) \quad \text{for every } n \in \mathbb{N}, f \in C_{cbs}^n(\tilde{X}), \sigma \in S_n(\tilde{X}),$$

obviously restricts to a chain map $v_{G,A}^*: C_{cbs}^*(\tilde{X}, \tilde{W})^G \rightarrow C_{cb}^*(\tilde{X}, \tilde{W})^G$. Under the identifications described in (3) and (10), this chain map induces the norm-nonincreasing map

$$H^*(v_{G,A}^*) : H_{cbs}^*(X, W) \rightarrow H_{cb}^*(X, W)$$

(we cannot realize $H^*(v_{G,A}^*)$ as the map induced by a morphism of pairs of resolutions just because we are not able to prove that the pair $C_{cb}^*(\tilde{X}), C_{cb}^*(\tilde{W})$ provides a pair of resolutions for $(G, A; \mathbb{R})$; see Remark 4.11 below).

It readily follows from the definitions that the following diagram commutes:

$$\begin{array}{ccc}
 H_{cbs}^*(X, W) & \xrightarrow{H^*(\eta_{G,A}^*)} & H_b^*(X, W) \\
 \searrow^{H^*(\nu_{G,A}^*)} & & \nearrow^{H^*(\rho_b^*)} \\
 & H_{cb}^*(X, W) &
 \end{array}$$

where $H^*(\rho_b^*) : H_{cb}^*(X, W) \rightarrow H_b^*(X, W)$ is the map described in the Introduction.

Now suppose that (X, W) is good. Then Theorem 4.1 implies that the map $H^*(\eta_{G,A}^*)$ is an isometric isomorphism, so the map $H^*(\nu_{G,A}^*) \circ H^*(\eta_{G,A}^*)^{-1}$ provides a right inverse to $H^*(\rho_b^*)$. Since $H^*(\nu_{G,A}^*)$ is norm-nonincreasing, this map is an isometric embedding, and this concludes the proof of Theorem 1.7.

Remark 4.11. Suppose that (X, W) is good. If we were able to prove that the complexes $C_{cb}^*(\tilde{X})$, $C_{cb}^*(\tilde{W})$ provide a proper pair of resolutions for $(G, A; \mathbb{R})$, then we could prove that $H^*(\rho_b^*) : H_{cb}^*(X, W) \rightarrow H_b^*(X, W)$ is an isometric isomorphism for every good pair (X, W) . However, it is not clear why Ivanov's contracting homotopies should take continuous cochains into continuous cochains, thus restricting to contracting homotopies for $C_{cb}^*(\tilde{X})$, $C_{cb}^*(\tilde{W})$.

4F. (Unbounded) continuous cohomology of pairs. We conclude the section by proving Theorem 1.9, which asserts that, when (X, W) is a locally finite good CW-pair, the map

$$H^*(\rho^*) : H_c^*(X, W) \rightarrow H^*(X, W)$$

is an isometric isomorphism.

We first observe that, since W is closed in X , the subspace $S_n(W)$ is closed in $S_n(X)$ for every $n \in \mathbb{N}$. Moreover, since X is locally finite, it is metrizable, and this implies that $S_n(X)$ is also metrizable. Therefore, by Tietze's theorem, every continuous cochain on W extends to a continuous cochain on X ; i.e., the restriction map $C_c^*(X) \rightarrow C_c^*(W)$ is surjective. As a consequence, both rows of the following commutative diagram are exact:

$$\begin{array}{ccccccccc}
 H_c^{n+1}(X) & \longrightarrow & H_c^{n+1}(W) & \longrightarrow & H_c^n(X, W) & \longrightarrow & H_c^n(X) & \longrightarrow & H_c^n(W) \\
 \downarrow & & \downarrow & & \downarrow^{H^n(\rho^*)} & & \downarrow & & \downarrow \\
 H^{n+1}(X) & \longrightarrow & H^{n+1}(W) & \longrightarrow & H^n(X, W) & \longrightarrow & H^n(X) & \longrightarrow & H^n(W).
 \end{array}$$

We know from [Frigerio 2011, Theorem 1.1] that, in the absolute case, the vertical arrows are isomorphisms, and the Five Lemma implies now that $H^n(\rho^*)$ is an isomorphism. We are left to show that it is also an isometry.

The inclusions $C_b^*(X, W) \hookrightarrow C^*(X, W)$, $C_{cb}^*(X, W) \hookrightarrow C_c^*(X, W)$ induce the comparison maps $c^* : H_b^*(X, W) \rightarrow H^*(X, W)$, $c_c^* : H_{cb}^*(X, W) \rightarrow H_c^*(X, W)$ and

it follows from the very definitions that for every $\varphi \in H^n(X, W)$, $\varphi_c \in H_c^n(X, W)$ the following equalities hold:

$$\begin{aligned} \|\varphi\|_\infty &= \inf\{\|\psi\|_\infty \mid \psi \in H_b^n(X, W), c^n(\psi) = \varphi\}, \\ \|\varphi_c\|_\infty &= \inf\{\|\psi_c\|_\infty \mid \psi_c \in H_{cb}^n(X, W), c_c^n(\psi_c) = \varphi_c\}, \end{aligned}$$

where we understand that $\inf \emptyset = +\infty$. Moreover, since $H^*(\rho^*) \circ c_c^* = c^* \circ H^*(\rho_b^*)$, for every $\varphi_c \in H_c^*(X, W)$ we have

$$\begin{aligned} \|H^*(\rho^*)(\varphi_c)\|_\infty &= \inf\{\|\psi\|_\infty \mid \psi \in H_b^*(X, W), c^*(\psi) = H^*(\rho^*)(\varphi_c)\} \\ &= \inf\{\|\psi_c\|_\infty \mid \psi_c \in H_{cb}^*(X, W), c^*(H^*(\rho_b^*)(\psi_c)) = H^*(\rho^*)(\varphi_c)\} \\ &= \inf\{\|\psi_c\|_\infty \mid \psi_c \in H_{cb}^*(X, W), H^*(\rho^*)(c_c^*(\psi_c)) = H^*(\rho^*)(\varphi_c)\} \\ &= \inf\{\|\psi_c\|_\infty \mid \psi_c \in H_{cb}^*(X, W), c_c^*(\psi_c) = \varphi_c\} = \|\varphi_c\|_\infty, \end{aligned}$$

where the second equality is due to Theorem 1.7 (recall that locally finite CW-pairs are countable). The proof of Theorem 1.9 is now complete.

5. The duality principle

This section is mainly devoted to the proof of Theorem 1.3. As already mentioned in the Introduction, once a suitable duality pairing between measure homology and continuous bounded cohomology is established, Theorem 1.3 can be easily deduced from Theorem 1.7.

5A. Duality between singular homology and bounded cohomology. Let us begin by recalling the well-known duality between bounded cohomology and singular homology. Let (X, W) be any pair of topological spaces. By definition, $C^n(X, W)$ is the algebraic dual of $C_n(X, W)$, and it is readily seen that the L^∞ -norm on $C^n(X, W)$ is dual to the L^1 -norm on $C_n(X, W)$. As a consequence, $C_b^n(X, W)$ coincides with the topological dual of $C_n(X, W)$. This does *not* imply that $H_b^n(X, W)$ is the topological dual of $H_n(X, W)$, because taking duals of normed chain complexes does not commute in general with homology (see [Löh 2008] for a detailed discussion of this issue). However, if we denote by

$$\langle \cdot, \cdot \rangle : H_b^n(X, W) \times H_n(X, W) \rightarrow \mathbb{R}$$

the *Kronecker product* induced by the pairing $C_b^n(X, W) \times C_n(X, W) \rightarrow \mathbb{R}$, then an application of Hahn–Banach theorem (for details, see [Löh 2007, Theorem 3.8], for instance) gives the following:

Proposition 5.1. *For every $\alpha \in H_n(X, W)$ we have*

$$\|\alpha\|_1 = \sup \left\{ \frac{1}{\|\varphi\|_\infty} \mid \varphi \in H_b^n(X, W), \langle \varphi, \alpha \rangle = 1 \right\},$$

where we understand that $\sup \emptyset = 0$.

5B. Duality between measure homology and continuous bounded cohomology.

The topological dual of $\mathcal{C}_*(X, W)$ does not admit an easy description, so in order to compute seminorms in $\mathcal{H}_*(X, W)$ via duality more work is needed. We first observe that, if μ is any measure on $S_n(X)$ with compact determination set and f is any continuous function on $S_n(X)$, it makes sense to integrate f with respect to μ . Therefore, for every $n \in \mathbb{N}$ the bilinear pairing

$$\langle \cdot, \cdot \rangle : C_{cb}^n(X, W) \times \mathcal{C}_n(X, W) \rightarrow \mathbb{R}, \quad \langle f, \mu \rangle = \int_{S_n(X)} f(\sigma) d\mu(\sigma)$$

is well-defined. It readily follows from the definitions that $|\langle f, \mu \rangle| \leq \|f\|_\infty \cdot \|\mu\|_m$ for every $f \in C_{cb}^n(X, W)$, $\mu \in \mathcal{C}_n(X, W)$, so $C_{cb}^*(X, W)$ lies in the topological dual of $\mathcal{C}_*(X, W)$. Moreover, for every $i \in \mathbb{N}$, $f \in C_{cb}^i(X, W)$ and $\mu \in \mathcal{C}_{i+1}(X, W)$ we have $\langle \delta f, \mu \rangle = \langle f, \partial \mu \rangle$, so this pairing defines a Kronecker product

$$\langle \cdot, \cdot \rangle : H_{cb}^n(X, W) \times \mathcal{H}_n(X, W) \rightarrow \mathbb{R}$$

such that

$$(11) \quad |\langle \varphi_c, \alpha \rangle| \leq \|\varphi_c\|_\infty \cdot \|\alpha\|_{mh} \quad \text{for every } \varphi_c \in H_{cb}^n(X, W), \alpha \in \mathcal{H}_n(X, W).$$

The following proposition is an immediate consequence of inequality (11), and provides a sort of weak duality theorem for continuous bounded cohomology and measure homology. The term “weak” refers to the fact that while Proposition 5.1 allows to compute seminorms in homology in terms of seminorms in bounded cohomology, here only an inequality is established. However, this turns out to be sufficient to our purposes. Moreover, once Theorem 1.3 is proved, one could easily prove that (in the case of good CW-pairs) the inequality of Proposition 5.2 is in fact an equality, thus recovering a “full” duality between continuous bounded cohomology and measure homology.

Proposition 5.2. *For every $\alpha \in \mathcal{H}_n(X, W)$ we have*

$$\|\alpha\|_{mh} \geq \sup \left\{ \frac{1}{\|\varphi_c\|_\infty} \mid \varphi_c \in H_{cb}^n(X, W), \langle \varphi_c, \alpha \rangle = 1 \right\},$$

where we understand that $\sup \emptyset = 0$.

To conclude the proof of Theorem 1.3, we need one more result, which follows readily from the definitions and ensures that the Kronecker products introduced above are compatible with each other:

Proposition 5.3. *For every $\varphi_c \in H_{cb}^n(X, W)$, $\alpha \in H_n(X, W)$ we have*

$$\langle H^n(\rho_b^*)(\varphi_c), \alpha \rangle = \langle \varphi_c, H_n(\iota_*)(\alpha) \rangle.$$

Proof of Theorem 1.3. Suppose that (X, W) is a good CW-pair. We already know that the map $H_*(\iota_*) : H_*(X, W) \rightarrow \mathcal{H}_*(X, W)$ is a norm-nonincreasing isomorphism, so we are left to show that $\|H_*(\iota_*)(\alpha)\|_{\text{mh}} \geq \|\alpha\|_1$ for every $\alpha \in H_*(X, W)$.

However, for every $\alpha \in H_n(X, W)$ we have

$$\begin{aligned} \|H_n(\iota_*)(\alpha)\|_{\text{mh}} &\geq \sup \left\{ \frac{1}{\|\varphi_c\|_\infty} \mid \varphi_c \in H_{cb}^n(X, W), \langle \varphi_c, H_n(\iota_*)(\alpha) \rangle = 1 \right\} \\ &= \sup \left\{ \frac{1}{\|\varphi_c\|_\infty} \mid \varphi_c \in H_{cb}^n(X, W), \langle H^n(\rho_b^*)(\varphi_c), \alpha \rangle = 1 \right\} \\ &= \sup \left\{ \frac{1}{\|\varphi\|_\infty} \mid \varphi \in H_b^n(X, W), \langle \varphi, \alpha \rangle = 1 \right\} \\ &= \|\alpha\|_1, \end{aligned}$$

where the inequality is due to Proposition 5.2, the first equality to Proposition 5.3, the second equality to Theorem 1.7, and the last equality to Proposition 5.1. \square

Remark 5.4. Let (X, W) be any CW-pair. The arguments described in this section show that if $H^*(\rho_b^*) : H_{cb}^*(X, W) \rightarrow H_b^*(X, W)$ admits a norm-nonincreasing right inverse, then the map $H_*(\iota_*) : H_*(X, W) \rightarrow \mathcal{H}_*(X, W)$ is an isometric isomorphism.

6. A comparison with Park’s seminorms

Park [2003] describes an algebraic foundation of relative bounded cohomology of pairs, both in the case of a pair of groups (G, A) equipped with a homomorphism $A \rightarrow G$ and in the case of a pair of path-connected topological spaces (X, W) equipped with a continuous map $W \rightarrow X$. However, recall from the Introduction that the seminorms considered by Park are quite different from the ones considered in this paper, which go back to [Gromov 1982]. In this section we investigate the relationships between our seminorms and the seminorms introduced in [Park 2003], proving in particular that there exist examples for which they are *not* isometric to each other.

6A. Park’s mapping cone for homology. Let (X, W) be a countable CW-pair, where both X and W are connected, and let us suppose that the inclusion $i : W \hookrightarrow X$ induces an injective map on the fundamental groups (several considerations here below also hold without this last assumption, but this is not relevant to our purposes). We also denote by $i_* : C_*(W) \rightarrow C_*(X)$ the map induced by the inclusion i . The homology mapping cone complex of (X, W) is the complex

$$(C_*(W \rightarrow X), \bar{d}_*) = (C_*(X) \oplus C_{*-1}(W), \bar{d}_*),$$

where

$$\begin{aligned} \bar{d}_n : C_n(X) \oplus C_{n-1}(W) &\rightarrow C_{n-1}(X) \oplus C_{n-2}(W) \\ (u_n, v_{n-1}) &\mapsto (d_n u_n + i_{n-1}(v_{n-1}), -d_{n-1} v_{n-1}), \end{aligned}$$

and d_* denotes the usual differential both of $C_*(X)$ and of $C_*(W)$. The homology of the mapping cone $(C_*(W \rightarrow X), \bar{d}_*)$ is denoted by $H_*(W \rightarrow X)$. For every $\omega \in [0, \infty)$ one can endow $C_*(W \rightarrow X)$ with the L^1 -norm

$$\|(u, v)\|_1(\omega) = \|u\|_1 + (1 + \omega)\|v\|_1,$$

which induces in turn a seminorm (still denoted by $\|\cdot\|_1(\omega)$) on $H_*(W \rightarrow X)$ (in fact, in [Park 2004] the case $\omega = \infty$ is also considered, but this is not relevant to our purposes).

As observed in [Park 2004], the chain map

$$(12) \quad \beta_* : C_*(W \rightarrow X) \rightarrow C_*(X, W) = C_*(X)/C_*(W), \quad \beta_*(u, v) = [u]$$

induces an isomorphism

$$H_*(\beta_*) : H_*(W \rightarrow X) \rightarrow H_*(X, W).$$

The explicit description of β_* implies that

$$\|H_*(\beta_*)(\alpha)\|_1 \leq \|\alpha\|_1(0) \leq \|\alpha\|_1(\omega)$$

for every $\alpha \in H_*(W \rightarrow X)$, $\omega \in [0, \infty)$.

6B. Park's mapping cone for bounded cohomology. We define the mapping cone for bounded cohomology as the (topological) dual of the mapping cone for homology. More precisely, we fix $\omega \in [0, \infty)$, and endow $C_*(W \rightarrow X)$ with the norm $\|\cdot\|_1(\omega)$. It is readily seen that the topological dual of $C_n(W \rightarrow X) = C_n(X) \oplus C_{n-1}(W)$ is isometrically isomorphic to the space

$$C_b^n(W \rightarrow X) = C_b^n(X) \oplus C_b^{n-1}(W)$$

endowed with the L^∞ -norm $\|\cdot\|_\infty(\omega)$ defined by

$$\|(f, g)\|_\infty(\omega) = \max\{\|f\|_\infty, (1 + \omega)^{-1}\|g\|_\infty\}.$$

In other words, the pairing

$$C_b^*(W \rightarrow X) \times C_*(W \rightarrow X) \rightarrow \mathbb{R}, \quad ((f, f'), (a, a')) \mapsto f(a) - f'(a')$$

realizes $C_b^*(W \rightarrow X)$ as the topological dual of $C_*(W \rightarrow X)$, and an easy computation shows that the norm $\|\cdot\|_\infty(\omega)$ just introduced on $C_b^*(W \rightarrow X)$ coincides with the operator norm (with respect to the norm $\|\cdot\|_1(\omega)$ fixed on $C_*(W \rightarrow X)$). Therefore, if $i^* : C_b^*(X) \rightarrow C_b^*(W)$ is the cochain map induced by the inclusion, then the

cohomology mapping cone complex of (X, W) is the complex $(C_b^*(W \rightarrow X), \bar{\delta}^*)$, where $\bar{\delta}^*$ is defined as the dual map of \bar{d}_* , and admits therefore the following explicit description (see [Park 2003] for the details):

$$\begin{aligned} \bar{\delta}^n : C_b^n(X) \oplus C_b^{n-1}(W) &\rightarrow C_b^{n+1}(X) \oplus C_b^n(W) \\ (f_n, g_{n-1}) &\mapsto (\delta^n f_n, -i^n(f_n) - \delta^{n-1} g_{n-1}) \end{aligned}$$

(here δ^* denotes the usual differential both of $C_b^*(X)$ and of $C_b^*(W)$). The cohomology of the complex $(C_b^*(W \rightarrow X), \bar{\delta}^*)$ is denoted by $H_b^*(W \rightarrow X)$. Just as in the case of homology, the L^∞ -norm $\|\cdot\|_\infty(\omega)$ on $C_b^n(W \rightarrow X)$ descends to a seminorm (still denoted by $\|\cdot\|_\infty(\omega)$) on $H_b^*(W \rightarrow X)$.

The chain map

$$\beta^* : C_b^*(X, W) \rightarrow C_b^*(W \rightarrow X), \quad \beta^*(f) = (f, 0)$$

is the dual of the chain map β_* introduced in Equation (12) above, and induces an isomorphism

$$H^*(\beta^*) : H_b^*(X, W) \rightarrow H_b^*(W \rightarrow X)$$

such that

$$\|H^*(\beta^*)(\varphi)\|_\infty(\omega) \leq \|H^*(\beta^*)(\varphi)\|_\infty(0) \leq \|\varphi\|_\infty$$

for every $\varphi \in H_b^*(X, W)$, $\omega \in [0, \infty)$. More precisely:

Theorem 6.1 [Park 2003, Theorem 4.6]. *For every $n \in \mathbb{N}$, the isomorphism $H^n(\beta^*)$ is such that*

$$\frac{1}{n+2} \|\varphi\|_\infty \leq \|H^n(\beta^*)(\varphi)\|_\infty(0) \leq \|\varphi\|_\infty \quad \text{for every } \varphi \in H_b^n(X, W).$$

It is asked in [Park 2003] whether $H^*(\beta^*)$ is actually an isometry or not. We show in Proposition 6.4 below that there exist examples for which $H^*(\beta^*)$ is *not* an isometry.

6C. Mapping cones and duality. In the previous subsection we have seen that, for every $\omega \geq 0$, the normed space $(C_b^*(W \rightarrow X), \|\cdot\|_\infty(\omega))$ coincides with the topological dual of the normed space $(C_*(W \rightarrow X), \|\cdot\|_1(\omega))$. We may therefore apply the duality result proved in [Löh 2007, Theorem 3.14], and obtain the following:

Proposition 6.2. *If the map*

$$H^*(\beta^*) : (H_b^*(X, W), \|\cdot\|_\infty) \rightarrow (H_b^*(W \rightarrow X), \|\cdot\|_\infty(\omega))$$

is an isometric isomorphism, then

$$\|H_*(\beta_*)(\alpha)\|_1 = \|\alpha\|_1(\omega)$$

for every $\alpha \in H_(X, W)$.*

6D. An explicit example. Let M be a compact, connected, oriented manifold with connected boundary, and suppose that the inclusion $i : \partial M \rightarrow M$ induces an injective homomorphism $i_* : \pi_1(\partial M) \rightarrow \pi_1(M)$.

We denote by $[M, \partial M]$ the (real) fundamental class in $H_n(M, \partial M)$ and we set

$$[\partial M \rightarrow M] = H_n(\beta_*)^{-1}([M, \partial M]) \in H_n(\partial M \rightarrow M).$$

The L^1 -seminorm $\|[M, \partial M]\|_1$ of the real fundamental class of M is usually known as the *simplicial volume* of M , and it is denoted simply by $\|M\|$. Similarly, the L^1 -seminorm of the real fundamental class $[\partial M] \in H_{n-1}(\partial M)$ is the simplicial volume of ∂M , and it is denoted by $\|\partial M\|$.

Lemma 6.3. *We have*

$$\|[\partial M \rightarrow M]\|_1(\omega) \geq \|M\| + (1 + \omega)\|\partial M\|.$$

Proof. It is shown in [Park 2004] that, if $\alpha \in C_i(M)$ is such that $d_i\alpha \in C_{i-1}(\partial M)$ (so that α defines an element $[\alpha] \in H_i(M, \partial M)$), then

$$H_i(\beta_*)^{-1}([\alpha]) = [(\alpha, -d_i\alpha)].$$

Therefore, if $\alpha \in C_n(M)$ is a representative of the fundamental class $[M, \partial M] \in H_n(M, \partial M)$, then $(\alpha, -d_n\alpha)$ is a representative of $[\partial M \rightarrow M] \in H_n(\partial M \rightarrow M)$. If (α', γ) is any other representative of such a class, then by definition of mapping cone there exist $x \in C_{n+1}(M)$ and $y \in C_n(\partial M)$ such that:

$$\alpha - \alpha' = d_{n+1}x + i_n(y) \quad \text{and} \quad \gamma + d_n\alpha = -d_n y.$$

These equalities readily imply that $[\alpha'] = [\alpha]$ in $H_n(M, \partial M)$ and $[\gamma] = [-d_n\alpha]$ in $H_{n-1}(\partial M)$. As a consequence, since $d_n\alpha$ is a representative of the fundamental class of ∂M , we have $\|\alpha'\|_1 \geq \|[\alpha']\|_1 = \|M\|$ and $\|\gamma\|_1 \geq \|[\gamma]\|_1 = \|\partial M\|$, whence

$$\|(\alpha', \gamma)\|_1(\omega) \geq \|M\| + (1 + \omega)\|\partial M\|.$$

The conclusion follows from the fact that (α', γ) is an arbitrary representative of $[\partial M \rightarrow M]$. \square

Proposition 6.4. *Let M be a compact connected oriented hyperbolic n -manifold with connected geodesic boundary. Then, for every $\omega \in [0, \infty)$ the isomorphism*

$$H^n(\beta^*) : (H_b^n(M, \partial M), \|\cdot\|_\infty) \rightarrow (H_b^n(\partial M \rightarrow M), \|\cdot\|_\infty(\omega))$$

is not isometric.

Proof. It is well-known that the inclusion $\partial M \hookrightarrow M$ induces an injective map on fundamental groups. Moreover, since ∂M is a closed oriented hyperbolic $(n-1)$ -manifold, we also have $\|\partial M\| > 0$. By Proposition 6.2, if $H^n(\beta^*)$ were an isometry

we would have $\|[\partial M \rightarrow M]\|_1(\omega) = \|[M, \partial M]\|_1 = \|M\|$, and this contradicts Lemma 6.3. \square

Acknowledgement

The authors thank Maria Beatrice Pozzetti for several useful conversations about the contents of [Ivanov 1985].

References

- [Berlanga 2008] R. Berlanga, “A topologised measure homology”, *Glasg. Math. J.* **50**:3 (2008), 359–369. MR 2009i:55008 Zbl 1167.55002
- [Bieri and Eckmann 1978] R. Bieri and B. Eckmann, “Relative homology and Poincaré duality for group pairs”, *J. Pure Appl. Algebra* **13**:3 (1978), 277–319. MR 80k:20048 Zbl 0392.20032
- [Bourgin 1952] D. G. Bourgin, “The paracompactness of the weak simplicial complex”, *Proc. Nat. Acad. Sci. U. S. A.* **38** (1952), 305–313. MR 14,70g Zbl 0046.40305
- [Bridson and Haefliger 1999] M. R. Bridson and A. Haefliger, *Metric spaces of non-positive curvature*, Grundlehren der Math. Wiss. **319**, Springer, Berlin, 1999. MR 2000k:53038 Zbl 0988.53001
- [Frigerio 2011] R. Frigerio, “(Bounded) continuous cohomology and Gromov’s proportionality principle”, *Manuscripta Math.* **134**:3–4 (2011), 435–474. MR 2012f:55006 Zbl 1220.55003
- [Frigerio and Pagliantini 2010] R. Frigerio and C. Pagliantini, “The simplicial volume of hyperbolic manifolds with geodesic boundary”, *Algebr. Geom. Topol.* **10** (2010), 979–1001. MR 2011c:53082 Zbl 1206.53045
- [Fujiwara and Manning 2011] K. Fujiwara and J. F. Manning, “Simplicial volume and fillings of hyperbolic manifolds”, *Algebr. Geom. Topol.* **11** (2011), 2237–2264. MR 2012g:53062 Zbl 05959839
- [Gromov 1982] M. Gromov, “Volume and bounded cohomology”, *Inst. Hautes Études Sci. Publ. Math.* **56** (1982), 5–99. MR 84h:53053 Zbl 0516.53046
- [Hansen 1998] S. K. Hansen, “Measure homology”, *Math. Scand.* **83**:2 (1998), 205–219. MR 86h:55005 Zbl 0932.55003
- [Hatcher 2002] A. Hatcher, *Algebraic topology*, Cambridge University Press, Cambridge, 2002. MR 2002k:55001 Zbl 1044.55001
- [Ivanov 1985] N. V. Ivanov, “Foundations of the theory of bounded cohomology”, *Zap. Nauchn. Sem. Leningrad. Otdel. Mat. Inst. Steklov. (LOMI)* **143** (1985), 69–109. In Russian; translated in *J. Soviet Math.* **37** (1987), 1090–1114. MR 87b:53070
- [Löh 2006] C. Löh, “Measure homology and singular homology are isometrically isomorphic”, *Math. Z.* **253**:1 (2006), 197–218. MR 2006m:55021 Zbl 1093.55004
- [Löh 2007] C. Löh, *l^1 -homology and simplicial volume*, Ph.D. thesis, WWU Münster, 2007, available at <http://nbn-resolving.de/urn:nbn:de:hbz:6-37549578216>. Zbl 1152.57304
- [Löh 2008] C. Löh, “Isomorphisms in l^1 -homology”, *Münster J. Math.* **1** (2008), 237–265. MR 2010b:55007 Zbl 1158.55007
- [Löh and Sauer 2009] C. Löh and R. Sauer, “Degree theorems and Lipschitz simplicial volume for nonpositively curved manifolds of finite volume”, *J. Topol.* **2**:1 (2009), 193–225. MR 2010j:53065 Zbl 1187.53043
- [Mineyev and Yaman 2007] I. Mineyev and A. Yaman, “Relative hyperbolicity and bounded cohomology”, preprint, 2007, available at <http://www.math.uiuc.edu/~mineyev/math/art/rel-hyp.pdf>.

- [Miyazaki 1952] H. Miyazaki, “The paracompactness of CW -complexes”, *Tôhoku Math. J. (2)* **4** (1952), 309–313. MR 14,894c Zbl 0049.12502
- [Monod 2001] N. Monod, *Continuous bounded cohomology of locally compact groups*, Lecture Notes in Mathematics **1758**, Springer, Berlin, 2001. MR 2002h:46121 Zbl 0967.22006
- [Pagliantini 2012] C. Pagliantini, *Relative (continuous) bounded cohomology and simplicial volume of hyperbolic manifolds with geodesic boundary*, Ph.D. thesis, University of Pisa, 2012. In preparation.
- [Park 2003] H. Park, “Relative bounded cohomology”, *Topology Appl.* **131**:3 (2003), 203–234. MR 2004e:55008 Zbl 1042.55003
- [Park 2004] H. Park, “Foundations of the theory of l_1 homology”, *J. Korean Math. Soc.* **41**:4 (2004), 591–615. MR 2005c:55011 Zbl 1061.55004
- [Thurston 1979] W. P. Thurston, “The geometry and topology of three-manifolds”, lecture notes, Princeton University, 1979, available at <http://msri.org/publications/books/gt3m>.
- [Zastrow 1998] A. Zastrow, “On the (non)-coincidence of Milnor–Thurston homology theory with singular homology theory”, *Pacific J. Math.* **186**:2 (1998), 369–396. MR 2000a:55008 Zbl 0933.55008

Received June 8, 2011. Revised February 9, 2012.

ROBERTO FRIGERIO
DIPARTIMENTO DI MATEMATICA
UNIVERSITÀ DI PISA
LARGO B. PONTECORVO 5
56121 PISA
ITALY
frigerio@dm.unipi.it

CRISTINA PAGLIANTINI
DIPARTIMENTO DI MATEMATICA
UNIVERSITÀ DI PISA
LARGO B. PONTECORVO 5
56121 PISA
ITALY
pagliantini@mail.dm.unipi.it

NORMAL ENVELOPING ALGEBRAS

ALEXANDRE N. GRISHKOV, MARINA RASSKAZOVA
AND SALVATORE SICILIANO

A full characterization is given of ordinary and restricted enveloping algebras which are normal with respect to the principal involution.

1. Introduction

Let A be an algebra with involution $*$ over a field \mathbb{F} . We recall that A is said to be normal if $xx^* = x^*x$ for every $x \in A$. Over the decades, normal algebras with involutions have been extensively investigated on their own; see, for example, [Beidar et al. 1981; Bovdi et al. 1985; Bovdi 1990; 1997; Bovdi and Siciliano 2007; Brešar and Vukman 1989; Herstein 1976; Knus et al. 1998; Lim 1977; 1979; Maxwell 1972]. Moreover, they have several applications in linear algebra and functional analysis; see, for example, [Berberian 1959; Fuglede 1950; Maxwell 1972; Mosić and Djordjević 2009; Putnam 1951; Yood 1974]. It is well-known that any normal algebra with involution satisfies the standard polynomial identity of degree 4 [Herstein 1976, Section 5]. Moreover, Maxwell [1972] determined the structure of a normal simple algebra of matrices with entries in a field with involution. He also proved that a division algebra D with involution is normal if and only if D is either a field or a generalized quaternion algebra over its center. Furthermore, a characterization of group algebras which are normal under the standard involution was established by Bovdi, Gudivok, and Semirov [Bovdi et al. 1985]. Subsequently, such a result has been extended to twisted group algebras [Bovdi 1990; 1997] and to group algebras under a Novikov involution [Bovdi and Siciliano 2007].

On the other hand, it seems that the rather natural problems of characterizing ordinary and restricted enveloping algebras which are normal under their canonical involutions have not been settled yet. The present paper is just devoted to answering these questions.

For an arbitrary Lie algebra L we denote by $U(L)$ the universal enveloping algebra of L . Moreover, if L is restricted with a p -map $[p]$ over a field \mathbb{F} of

The first author was supported by FAPESP and CNPq (Brazil) and grant RFFI-10.01.00383a (Russia).

MSC2010: 16S30, 16W10, 17B50.

Keywords: restricted Lie algebra, enveloping algebra, normal ring, principal involution.

characteristic $p > 0$, then we denote by $u(L)$ the restricted enveloping algebra of L . We consider $U(L)$ and $u(L)$ with the *principal involution* $*$, namely, the unique \mathbb{F} -antiautomorphism such that $x^* = -x$ for every x in L ; see [Bourbaki 2007, Section 2] or [Dixmier 1974, Section 2]. Note that $*$ is just the antipode of the \mathbb{F} -Hopf algebras $U(L)$ or $u(L)$.

We use the symbols $Z(L)$ and L' for the center of L and the derived subalgebra of L , respectively. If $S \subseteq L$, we denote by $\langle S \rangle_{\mathbb{F}}$ the \mathbb{F} -vector space generated by S . Also, if L is restricted, $\langle S \rangle_p$ denotes the restricted subalgebra generated by S , and we put $S^{[p]} = \{x^{[p]} \mid x \in S\}$. In our first main result we completely settle the restricted case:

Theorem 1.1. *Let L be a restricted Lie algebra over a field \mathbb{F} of characteristic $p > 0$. Then $u(L)$ is normal if and only if either L is abelian or $p = 2$, L is nilpotent of class 2, and one of the following conditions holds:*

- (i) L contains an abelian restricted ideal I of codimension 1.
- (ii) $\dim_{\mathbb{F}} L/Z(L) = 3$.
- (iii) $\dim_{\mathbb{F}} L' = 1$ and $(L')^{[2]} = 0$.
- (iv) $L = \langle x, x_1, x_2, x_3 \rangle_p + Z(L)$ with

$$[x_1, x_2] = \xi[x, x_3],$$

$$[x_1, x_3] = \mu[x, x_2],$$

$$[x_2, x_3] = \lambda[x, x_1],$$

and

$$\lambda[x, x_1]^{[2]} + \mu[x, x_2]^{[2]} + \xi[x, x_3]^{[2]} = 0$$

for some $\lambda, \mu, \xi \in \mathbb{F}$.

Afterwards we apply Theorem 1.1 in order to solve the ordinary case:

Theorem 1.2. *Let L be a Lie algebra over an arbitrary field \mathbb{F} . Then $U(L)$ is normal if and only if either L is abelian or $p = 2$, L is nilpotent of class 2, and one of the following conditions holds:*

- (i) L contains an abelian ideal of codimension 1.
- (ii) $\dim_{\mathbb{F}} L/Z(L) = 3$.

2. Proofs

For any associative algebra A , we shall consider the Lie bracket on A defined by $[a, b] := ab - ba \in A$, $a, b \in A$. The symbol $Z(A)$ will denote the center of A . Moreover, for a subset S of a Lie algebra L we shall denote by $C_L(S)$ the centralizer of S in L .

It is easy to verify that a normal algebra with involution satisfies the $*$ -polynomial identity $[x, y] = [x^*, y^*]$. The converse is also true in characteristic different from 2, but in general it fails without such an assumption [Lim 1977]. However, for restricted Lie algebras we have the following:

Lemma 2.1. *Let L be a restricted Lie algebra over a field \mathbb{F} of characteristic 2 such that $[x, y] = [x^*, y^*]$ for every $x, y \in u(L)$. Then L is nilpotent of class at most 2 and $u(L)$ is normal.*

Proof. For every $a, b, c \in L$, we have

$$0 = [ab, c] + [(ab)^*, c^*] = [[a, b], c].$$

Hence L is nilpotent of class at most 2.

Let $(e_i)_{i \in I}$ be an ordered \mathbb{F} -basis of L . Then every element u of $u(L)$ is an \mathbb{F} -linear combination of elements $e_{i_1} \cdots e_{i_m}$, where $m \geq 0$ and the indices $i_1 < \cdots < i_m$ are in I . As L is nilpotent of class at most 2, for every $z \in L$ we have $z^{[2]} \in Z(L)$, and then

$$[e_{i_1} \cdots e_{i_m}, (e_{i_1} \cdots e_{i_m})^*] = 0.$$

Moreover, by hypothesis we clearly have $[x, y^*] = [x^*, y]$ for every $x, y \in u(L)$. We conclude that $[u, u^*] = 0$, so that $u(L)$ is normal. \square

Lemma 2.2. *Let L be a restricted Lie algebra over a field \mathbb{F} of characteristic $p > 0$ such that $u(L)$ is normal. Then either L is abelian, or $p = 2$ and L is nilpotent of class 2.*

Proof. As $u(L)$ satisfies the $*$ -polynomial identity $[x, y] = [x^*, y^*]$, if $p = 2$, Lemma 2.1 assures that L is nilpotent of class at most 2. Now suppose $p > 2$. For every $x, y \in L$, we have

$$0 = [x^2 + y, (x^2 + y)^*] = -4x[x, y] + 2[x, [x, y]].$$

Since $p > 2$, in view of the Poincaré–Birkhoff–Witt (PBW) theorem for restricted Lie algebras [Strade and Farnsteiner 1988, Section 2, Theorem 5.1], the previous relation is possible only when $[x, y] = 0$, so that L is abelian. This yields the claim. \square

Let L be a restricted Lie algebra over a field of characteristic 2. For every $a, b, c, d \in L$, we put

$$\Theta(a, b, c, d) := [a, b][c, d] + [a, c][b, d] + [a, d][b, c] \in u(L).$$

The following result will be extremely useful in the sequel.

Lemma 2.3. *Let L be a restricted Lie algebra over a field \mathbb{F} of characteristic 2, and suppose L to be nilpotent of class 2. Then $u(L)$ is normal if and only if $\Theta(a, b, c, d) = 0$ for all $a, b, c, d \in L$.*

Proof. If $u(L)$ is normal, for all $a, b, c, d \in L$ we have

$$\Theta(a, b, c, d) = [a, bcd] + [a, dcb] = [a, bcd] + [a, (bcd)^*] = 0.$$

Conversely, assume that $\Theta(a, b, c, d) = 0$ for all $a, b, c, d \in L$. Let $(e_j)_{j \in J}$ be an ordered \mathbb{F} -basis of L containing an \mathbb{F} -basis of $Z(L)$. Since $u(L)$ is a free $u(Z(L))$ -module, there exists a unique homomorphism of $u(Z(L))$ -modules

$$\phi : u(L) \rightarrow u(L),$$

which vanishes on 1 and L , and such that for every $n > 1$ and $j_1 < \dots < j_n$, one has

$$\phi(e_{j_1} \cdots e_{j_n}) = \sum_{1 \leq h < k \leq n} e_{j_1} \cdots \hat{e}_{j_h} \cdots \hat{e}_{j_k} \cdots e_{j_n} [e_{j_h}, e_{j_k}],$$

where the symbol \hat{e}_{i_h} indicates that e_{i_h} is to be omitted.

We claim that

$$\text{Im}(\phi) \subseteq Z(u(L)).$$

For this purpose it is enough to prove that $[x, \phi(e_{j_1} \cdots e_{j_n})] = 0$ for every $x \in L$, $n > 1$, and $j_1, \dots, j_n \in J$ with $j_1 < \dots < j_n$. Indeed, by the hypothesis we have

$$\begin{aligned} [x, \phi(e_{j_1} \cdots e_{j_n})] &= \left[x, \sum_{1 \leq h < k \leq n} e_{j_1} \cdots \hat{e}_{j_h} \cdots \hat{e}_{j_k} \cdots e_{j_n} [e_{j_h}, e_{j_k}] \right] \\ &= \sum_{1 \leq h < k \leq n} \sum_{\substack{1 \leq s \leq n \\ s \neq h, k}} e_{j_1} \cdots \hat{e}_{j_h} \cdots \hat{e}_{j_s} \cdots \hat{e}_{j_k} \cdots e_{j_n} [e_{j_h}, e_{j_k}] [x, e_{j_s}] \\ &= \sum_{1 \leq h < k < s \leq n} e_{j_1} \cdots \hat{e}_{j_h} \cdots \hat{e}_{j_k} \cdots \hat{e}_{j_s} \cdots e_{j_n} ([e_{j_h}, e_{i_k}] [x, e_{j_s}] \\ &\quad + [e_{j_h}, e_{i_s}] [x, e_{j_k}] + [e_{j_k}, e_{i_s}] [x, e_{j_h}]) = 0, \end{aligned}$$

yielding the claim.

Now we shall prove that

$$a = a^* + \phi(a)$$

for every $a \in u(L)$. For this purpose it is enough to show that for all $n \geq 0$ and $j_1, \dots, j_n \in J$ with $j_1 < \dots < j_n$, one has

$$e_{j_1} \cdots e_{j_n} = e_{j_n} \cdots e_{j_1} + \phi(e_{j_1} \cdots e_{j_n}).$$

Let us proceed by induction on n . By the proved claim and the inductive assumption, we have, for $n > 0$,

$$\begin{aligned} e_{j_1} \cdots e_{j_n} &= (e_{j_{n-1}} \cdots e_{j_1})e_{j_n} + \phi(e_{j_1} \cdots e_{j_{n-1}})e_{j_n} \\ &= e_{j_n}e_{j_{n-1}} \cdots e_{j_1} + [e_{j_{n-1}} \cdots e_{j_1}, e_{j_n}] + \phi(e_{j_1} \cdots e_{j_{n-1}})e_{j_n} \\ &= e_{j_n}e_{j_{n-1}} \cdots e_{j_1} + [e_{j_1} \cdots e_{j_{n-1}}, e_{j_n}] + [\phi(e_{j_1} \cdots e_{j_{n-1}}), e_{j_n}] + \phi(e_{j_1} \cdots e_{j_{n-1}})e_{j_n} \\ &= e_{j_n} \cdots e_{j_1} + \phi(e_{j_1} \cdots e_{j_n}), \end{aligned}$$

completing the inductive step.

Finally, by applying the properties proved above, for all $a, b \in u(L)$, we have

$$[a, b] = [a^* + \phi(a), b^* + \phi(b)] = [a^*, b^*].$$

Hence $u(L)$ is normal by Lemma 2.1, as required. \square

Remark 2.4. Since Θ is an alternating \mathbb{F} -multilinear function, by Lemma 2.3 it is clear that in order to conclude that $u(L)$ is normal, it suffices to check that $\Theta(a, b, c, d) = 0$ for all pairwise distinct noncentral elements a, b, c, d in a fixed \mathbb{F} -basis of L .

We are now in position to prove Theorem 1.1:

Proof of Theorem 1.1. Assume that $u(L)$ is normal and L is not abelian. Then, by Lemma 2.3, we know that \mathbb{F} has characteristic 2 and L is nilpotent of class 2. Let us proceed with a case-by-case analysis.

Case 1. $\max\{\dim_{\mathbb{F}}[L, x] \mid x \in L\} = 1$. Let x_1 and y_1 be two noncommuting element of L and put $z_1 := [x_1, y_1]$. By assumption we have $[L, x_1] = [L, y_1] = \mathbb{F} z_1$ and $L = \mathbb{F} y_1 \oplus C_L(x_1)$. Now, if $C_L(x_1)$ is abelian, L satisfies alternative (i) of the statement. Suppose then that there exist $x_2, y_2 \in C_L(x_1)$ such that $[x_2, y_2] := z_2 \neq 0$. From Lemma 2.3 it follows that

$$(1) \quad z_1 z_2 = \Theta(x_1, y_1, x_2, y_2) = 0.$$

Therefore the PBW theorem for restricted Lie algebras entails that $z_1 = \lambda z_2$ for some $\lambda \in \mathbb{F}$, which shows that $L' = \mathbb{F} z_1$. Also, as $\lambda \neq 0$, by (1), we have $z_1^{[2]} = 0$. Thus $(L')^{[2]} = 0$, and alternative (iii) of the statement holds.

Case 2. $\max\{\dim_{\mathbb{F}}[L, x] \mid x \in L\} = 2$. Let $x, x_1, x_2 \in L$ such that $z_1 := [x, x_1]$ and $z_2 := [x, x_2]$ are \mathbb{F} -linearly independent. We clearly have $L = \langle x_1, x_2 \rangle_{\mathbb{F}} \oplus C_L(x)$. Furthermore, by Lemma 2.3, we have, for all $y_1, y_2 \in C_L(x)$,

$$0 = \Theta(x, x_1, y_1, y_2) = z_1[y_1, y_2] \quad \text{and} \quad 0 = \Theta(x, x_2, y_1, y_2) = z_2[y_1, y_2].$$

Since z_1 and z_2 are \mathbb{F} -linearly independent, the PBW theorem forces $[y_1, y_2] = 0$. Hence $C_L(x)$ is abelian. Again by Lemma 2.3, for every $y \in C_L(x)$, we have

$$(2) \quad 0 = \Theta(x, x_1, x_2, y) = z_1[x_2, y] + z_2[x_1, y].$$

At this stage, a straightforward application of the PBW theorem yields

$$[x_1, y] = \lambda_{11}(y)z_1 + \lambda_{12}(y)z_2 \quad \text{and} \quad [x_2, y] = \lambda_{21}(y)z_1 + \lambda_{22}(y)z_2$$

for some $\lambda_{11}(y), \lambda_{12}(y), \lambda_{21}(y), \lambda_{22}(y) \in \mathbb{F}$. From (2) it follows that

$$(\lambda_{11}(y) + \lambda_{22}(y))z_1z_2 = \lambda_{21}(y)z_1^2 + \lambda_{12}(y)z_2^2 \in L,$$

and, again by the PBW theorem, the preceding relation is possible only when $\lambda_{11}(y) = \lambda_{22}(y) := \lambda(y)$. With the notation just introduced, we consider the following subcases.

Subcase 2.1. For every $u \in C_L(x)$, one has $\lambda_{12}(u) = \lambda_{21}(u) = 0$. Let $y \in C_L(x)$ and put $\bar{y} := \lambda(y)x + y$. Then we have $[\bar{y}, x] = [\bar{y}, x_1] = [\bar{y}, x_2] = 0$. As $C_L(x)$ is abelian, it follows that $\bar{y} \in Z(L)$ and then $C_L(x) = \mathbb{F}x \oplus Z(L)$. Thus $\dim_{\mathbb{F}} L/Z(L) = 3$, and alternative (ii) of the statement holds.

Subcase 2.2. There exists $u \in C_L(x)$ such that $\lambda_{12}(u) \neq 0$ and $\lambda_{21}(u) = 0$. By replacing u by $\lambda_{12}^{-1}(u)u$, we can suppose that $\lambda_{12}(u) = 1$. Put $y := \lambda(u)x + u$. Then we have

$$[x_1, y] = z_2 \quad \text{and} \quad [x_2, y] = 0.$$

Let $y_1 \in C_L(x)$. Since $C_L(x)$ is abelian, by Lemma 2.3 we have

$$(3) \quad 0 = \Theta(x_1, x_2, y, y_1) = z_2[x_2, y_1] = z_2(\lambda_{21}(y_1)z_1 + \lambda(y_1)z_2).$$

Consequently, as z_1 and z_2 are \mathbb{F} -linearly independent, the PBW theorem forces $\lambda_{21}(y_1) = 0$. Also, from relation (3) (applied for $y_1 = x$), we infer that $z_2^2 = 0$. Now put $\bar{y}_1 := \lambda(y_1)x + \lambda_{12}(y_1)y + y_1$. Then $\bar{y}_1 \in Z(L)$, and $C_L(x) = \mathbb{F}x \oplus \mathbb{F}y \oplus Z(L)$. We conclude that $L = \langle x, x_1, x_2, y \rangle_p + Z(L)$, and it is clear that L is a restricted Lie algebra satisfying alternative (iv) of the statement.

Subcase 2.3. There exists $u \in C_L(x)$ such that $\lambda_{12}(u) = 0$ and $\lambda_{21}(u) \neq 0$. This is analogous to Subcase 2.2.

Subcase 2.4. There exists $u \in C_L(x)$ such that $\lambda_{12}(u) \neq 0$ and $\lambda_{21}(u) \neq 0$. By replacing u by $\lambda_{12}^{-1}(u)u$, we can suppose that $\lambda_{12}(u) = 1$. Put $y := \lambda(u)x + u$. Then we have

$$[x_1, y] = z_2 \quad \text{and} \quad [x_2, y] = \lambda_{21}(u)z_1.$$

Moreover, Lemma 2.3 yields

$$0 = \Theta(x, x_1, x_2, y) = \lambda_{21}(u)z_1^2 + z_2^2.$$

Let $y_1 \in C_L(x)$ and put $\bar{y}_1 := \lambda(y_1)x + y_1$. As $C_L(x)$ is abelian, Lemma 2.3 yields $0 = \Theta(x_1, x_2, y, \bar{y}_1) = z_2[x_2, \bar{y}_1] + \lambda_{21}(u)z_1[x_1, \bar{y}_1] = (\lambda_{21}(\bar{y}_1) + \lambda_{21}(u)\lambda_{12}(\bar{y}_1))z_1z_2$, so that $\lambda_{21}(\bar{y}_1) = \lambda_{21}(u)\lambda_{12}(\bar{y}_1)$. Put $\hat{y}_1 := \bar{y}_1 + \lambda_{12}(\bar{y}_1)y$. Then we have $[x_1, \hat{y}_1] = 0$. Now, if for some $y_1 \in C_L(x)$ one has $[x_2, \hat{y}_1] = \lambda_{21}(\hat{y}_1)z_1 \neq 0$ then we can replace y by \hat{y}_1 and conclude by Subcase 2.3 that alternative (iv) holds. On the other hand, if $[x_2, \hat{y}_1] = 0$ for every $y_1 \in C_L(x)$ then $L = \langle x, x_1, x_2, y \rangle_p + Z(L)$, and it is clear that, also in this case, L is a restricted Lie algebra satisfying alternative (iv).

Case 3. $\max\{\dim_{\mathbb{F}}[L, x] \mid x \in L\} = 3$. Let $x, u_1, u_2, u_3 \in L$ such that $z_1 := [x, u_1]$, $z_2 := [x, u_2]$, and $z_3 := [x, u_3]$ are \mathbb{F} -linearly independent. We clearly have $L = \langle u_1, u_2, u_3 \rangle_{\mathbb{F}} \oplus C_L(x)$, and one can show that $C_L(x)$ is abelian in the same way as in Case 2. Moreover, in view of Lemma 2.3, we have

$$(4) \quad 0 = \Theta(x, u_1, u_2, u_3) = z_1[u_2, u_3] + z_2[u_1, u_3] + z_3[u_1, u_2].$$

Thus, for every $1 \leq i < j \leq 3$, by the PBW theorem, we see that

$$(5) \quad [u_i, u_j] = \sum_{k=1}^3 \alpha_{ij}^{(k)} z_k,$$

where $\alpha_{ij}^{(k)} \in \mathbb{F}$, $k = 1, 2, 3$. By (4) and (5), another application of the PBW theorem yields

$$\alpha_{12}^{(1)} = \alpha_{23}^{(3)}, \quad \alpha_{12}^{(2)} = \alpha_{13}^{(3)}, \quad \alpha_{13}^{(1)} = \alpha_{23}^{(2)}.$$

Put

$$x_1 := u_1 + \alpha_{12}^{(2)}x, \quad x_2 := u_2 + \alpha_{12}^{(1)}x, \quad x_3 := u_3 + \alpha_{13}^{(1)}x,$$

and, moreover, $\alpha_{23}^{(1)} := \lambda$, $\alpha_{13}^{(2)} := \mu$, and $\alpha_{12}^{(3)} := \xi$. Then we have

$$[x_1, x_2] = \xi z_3, \quad [x_1, x_3] = \mu z_2, \quad [x_2, x_3] = \lambda z_1.$$

From Lemma 2.3 it follows that

$$\lambda z_1^{[2]} + \mu z_2^{[2]} + \xi z_3^{[2]} = \Theta(x, x_1, x_2, x_3) = 0.$$

Now, let $y \in C_L(x)$. By Lemma 2.3 we obtain

$$\begin{aligned} \Theta(x, x_1, x_2, y) &= z_1[x_2, y] + z_2[x_1, y] = 0, \\ \Theta(x, x_1, x_3, y) &= z_1[x_3, y] + z_3[x_1, y] = 0, \\ \Theta(x, x_2, x_3, y) &= z_2[x_3, y] + z_3[x_2, y] = 0. \end{aligned}$$

Consequently, by the PBW theorem there exists $\beta \in \mathbb{F}$ such that $[x_i, y] = \beta z_i$ for every $i = 1, 2, 3$. Put $\bar{y} := y + \beta x$. Then $\bar{y} \in Z(L)$ and $C_L(x) = \mathbb{F}x \oplus Z(L)$. We conclude that $L = \langle x, x_1, x_2, x_3 \rangle_p + Z(L)$, and alternative (iv) is satisfied.

Case 4. $\max\{\dim_{\mathbb{F}}[L, x] \mid x \in L\} > 3$. Let $S := (u_i)_{i \in I}$ be a subset of L such that the elements $z_i := [x, u_i]$, $i \in I$, are \mathbb{F} -linearly independent, and $[S, x] = [L, x]$. We clearly have $L = \langle S \rangle_{\mathbb{F}} \oplus C_L(x)$, and one can show that $C_L(x)$ is abelian by proceeding in a similar way as in Case 2. Let $i, j \in I$, $i \neq j$. In view of Lemma 2.3, for every $k \in I \setminus \{i, j\}$, we have

$$0 = \Theta(x, u_i, u_j, u_k) = z_i[u_j, u_k] + z_j[u_i, u_k] + z_k[u_i, u_j].$$

At this stage, by arguing as in the first case of Case 3, we have that $[u_i, u_j] \in \mathbb{F}z_k$. As $|I| > 3$, we conclude that $[u_i, u_j] = 0$. Finally, let $y \in C_L(x)$. By Lemma 2.3, for all pairwise distinct elements i, j, k of I , we have

$$\Theta(x, u_i, u_j, y) = z_i[u_j, y] + z_j[u_i, y] = 0,$$

$$\Theta(x, u_i, u_k, y) = z_i[x_k, y] + z_k[u_i, y] = 0.$$

Therefore, an application of the PBW theorem shows that there exists $\beta \in \mathbb{F}$ such that $[u_i, y] = \beta z_i$ for every $i \in I$. Put $\bar{y} := y + \beta x$. Then $\bar{y} \in Z(L)$, so that $C_L(x) = \mathbb{F}x \oplus Z(L)$. Therefore, as $L^{[2]} \subseteq Z(L)$, we conclude that $Z(L) + \langle S \rangle_{\mathbb{F}}$ is an abelian restricted ideal of codimension 1 in L , and the proof of the necessity part is finished.

Now let us prove sufficiency. The claim is trivial if L is abelian. Then assume that the ground field has characteristic 2 and L is nilpotent of class 2. If L has an abelian restricted ideal of codimension 1, it is clear that $\Theta(a, b, c, d) = 0$ for any $a, b, c, d \in L$, and so, by Lemma 2.3, $u(L)$ is normal. Also, if $\dim_{\mathbb{F}} L/Z(L) = 3$ then $u(L)$ is normal by Lemma 2.3 and Remark 2.4. Furthermore, the claim is clear whenever $L' = \mathbb{F}z$ for some $0 \neq z \in L$ with $z^{[2]} = 0$. Finally suppose that alternative (iv) holds. We can assume that x, x_1, x_2 , and x_3 are \mathbb{F} -linearly independent (otherwise alternative (i) or (ii) holds). Extend the set $\{x, x_1, x_2, x_3\}$ by central elements in order to form an \mathbb{F} -basis of L . We have

$$\begin{aligned} \Theta(x, x_1, x_2, x_3) &= [x, x_1][x_2, x_3] + [x, x_2][x_1, x_3] + [x, x_3][x_1, x_2] \\ &= \lambda[x, x_1]^{[2]} + \mu[x, x_2]^{[2]} + \xi[x, x_3]^{[2]} = 0. \end{aligned}$$

From Lemma 2.3 and Remark 2.4 it follows that $u(L)$ is normal. □

Finally, we deal with ordinary universal enveloping algebras of arbitrary Lie algebras. Indeed, we shall prove Theorem 1.2 as a consequence of Theorem 1.1.

Proof of Theorem 1.2. Suppose first that ground field \mathbb{F} has characteristic zero. If L is abelian then $U(L)$ is obviously normal. On the other hand, if $U(L)$ is normal then it satisfies the standard polynomial identity of degree 4 [Herstein 1976, Section 5]. Therefore, in view of a theorem of Latyshev [Bahturin 1987, Section 6.7,

Theorem 25], L is necessarily abelian. Now suppose $p > 0$. Put

$$\hat{L} := \sum_{k \geq 0} L^{p^k} \subseteq U(L),$$

where L^{p^k} is the \mathbb{F} -vector space spanned by the set $\{l^{p^k} \mid l \in L\}$. Then \hat{L} is a restricted Lie algebra with $h^{[p]} = h^p$ for all $h \in \hat{L}$. Moreover, by [Strade 2004, Section 1, Corollary 1.1.4], we have $U(L) = u(\hat{L})$, and then Theorem 1.1 applies. Suppose first that $U(L)$ is normal. If $p > 2$, Theorem 1.1 forces \hat{L} (and so L) to be abelian. Now assume that $p = 2$ and L is not abelian. Then \hat{L} satisfies one of the alternatives (i)–(iv) in the statement of Theorem 1.1. If \hat{L} contains an abelian restricted ideal of codimension 1 then L contains an abelian ideal of codimension 1. Likewise, if $\dim_{\mathbb{F}} \hat{L}/Z(\hat{L}) = 3$, $\dim_{\mathbb{F}} L/Z(L) = 3$. Observe that, as $u(\hat{L}) = U(L)$ is a domain, alternative (iii) in the statement of Theorem 1.1 cannot occur. Finally, suppose that $\hat{L} = \langle x, x_1, x_2, x_3 \rangle_p + Z(\hat{L})$, where x, x_1, x_2 , and x_3 are elements of L with $[x_1, x_2] = \xi[x, x_3]$, $[x_1, x_3] = \mu[x, x_2]$, $[x_2, x_3] = \lambda[x, x_1]$, and

$$\lambda[x, x_1]^{[2]} + \mu[x, x_2]^{[2]} + \xi[x, x_3]^{[2]} = 0$$

for some $\lambda, \mu, \xi \in \mathbb{F}$. Now, if $\dim_{\mathbb{F}} L' = 3$, the PBW theorem for ordinary enveloping algebras forces $\lambda = \mu = \xi = 0$. Hence L contains an abelian ideal of codimension 1. If $\dim_{\mathbb{F}} L' = 2$, we can suppose without loss of generality that $[x, x_1]$ and $[x, x_2]$ are \mathbb{F} -linearly independent and $[x, x_3] = \alpha[x, x_1] + \beta[x, x_2]$ for suitable $\alpha, \beta \in \mathbb{F}$. Consequently, we have

$$\alpha^2 \xi [x, x_1]^2 + \beta^2 \xi [x, x_2]^2 = \xi [x, x_3]^2 = \lambda [x, x_1]^2 + \mu [x, x_2]^2,$$

and the PBW theorem gets $\lambda = \alpha^2 \xi$ and $\mu = \beta^2 \xi$. Put

$$y := \alpha \beta \xi x + \alpha x_1 + \beta x_2 + x_3.$$

Then $y \in Z(\hat{L})$ and $\hat{L} = \langle x, x_1, x_2, y \rangle_p + Z(\hat{L})$. It follows that $\dim_{\mathbb{F}} \hat{L}/Z(\hat{L}) = 3$ and then $\dim_{\mathbb{F}} L/Z(L) = 3$ as well. Finally, if $\dim_{\mathbb{F}} L' = 1$ then it is easy to see that L contains an abelian ideal of codimension 1, and the necessity part is proved. Sufficiency easily follows from Theorem 1.1 and the fact that $U(L) = u(\hat{L})$. \square

Acknowledgement

We thank W. de Graaf, S. Cicalò, and the referee for useful comments.

References

- [Bahturin 1987] Y. A. Bahturin, *Identical relations in Lie algebras*, VNU Science Press, Utrecht, 1987. MR 88f:17032 Zbl 0691.17001
- [Beidar et al. 1981] C. I. Beidar, A. V. Mikhalev, and C. Salavova, “Generalized identities and semiprime rings with involution”, *Math. Z.* **178**:1 (1981), 37–62. MR 83b:16012 Zbl 0471.16008

- [Berberian 1959] S. K. Berberian, “Note on a theorem of Fuglede and Putnam”, *Proc. Amer. Math. Soc.* **10**:2 (1959), 175–182. MR 21 #6548 Zbl 0092.32004
- [Bourbaki 2007] N. Bourbaki, *Groupes et algèbres de Lie: Chapitre 1*, Springer, Berlin, 2007. MR 42 #6159 Zbl 1120.17001
- [Bovdi 1990] V. A. Bovdi, “Normal twisted group rings”, *Dokl. Akad. Nauk Ukrain. SSR Ser. A* **7** (1990), 6–8. In Russian. MR 91i:20004 Zbl 0718.16020
- [Bovdi 1997] V. A. Bovdi, “Structure of normal twisted group rings”, *Publ. Math. Debrecen* **51**:3-4 (1997), 279–293. MR 98j:16016 Zbl 0906.16013
- [Bovdi and Siciliano 2007] V. A. Bovdi and S. Siciliano, “Normality in group rings”, *Algebra i Analiz* **19**:2 (2007), 1–9. In Russian; translated in *St. Petersburg Math. J.* **19**:2 (2008), 159–165. MR 2008d:16040 Zbl 1200.16036
- [Bovdi et al. 1985] A. A. Bovdi, P. M. Gudivok, and M. S. Semirov, “Normal group rings”, *Ukrain. Mat. Zh.* **37**:1 (1985), 3–8. In Russian; translated in *Ukrain. Math. J.* **37**:1 (1985), 1–5. MR 86h:16013 Zbl 0572.16005
- [Brešar and Vukman 1989] M. Brešar and J. Vukman, “On some additive mappings in rings with involution”, *Aequationes Math.* **38**:2-3 (1989), 178–185. MR 90j:16076 Zbl 0691.16041
- [Dixmier 1974] J. Dixmier, *Algèbres enveloppantes*, Cahiers Scientifiques **37**, Gauthier-Villars, Paris, 1974. Translated as *Enveloping algebras*, Graduate Studies in Mathematics **11**, Amer. Math. Soc., Providence, RI, 1996. MR 58 #16803a Zbl 0308.17007
- [Fuglede 1950] B. Fuglede, “A commutativity theorem for normal operators”, *Proc. Nat. Acad. Sci. USA* **36**:1 (1950), 35–40. MR 11,371c Zbl 0035.35804
- [Herstein 1976] I. N. Herstein, *Rings with involution*, University of Chicago Press, Chicago, IL, 1976. MR 56 #406 Zbl 0343.16011
- [Knus et al. 1998] M.-A. Knus, A. Merkurjev, M. Rost, and J.-P. Tignol, *The book of involutions*, Amer. Math. Soc. Colloq. Publ. **44**, Amer. Math. Soc., Providence, RI, 1998. MR 2000a:16031 Zbl 0955.16001
- [Lim 1977] T. P. Lim, “Conjugacy of elements in a normal ring”, *Canad. Math. Bull.* **20** (1977), 113–115. MR 56 #3049 Zbl 0358.16008
- [Lim 1979] T. P. Lim, “Some classes of rings with involution satisfying the standard polynomial of degree 4”, *Pacific J. Math.* **85**:1 (1979), 125–130. MR 81j:16018 Zbl 0444.16004
- [Maxwell 1972] G. Maxwell, “Algebras of normal matrices”, *Pacific J. Math.* **43**:2 (1972), 421–428. MR 47 #6736 Zbl 0234.15028
- [Mosić and Djordjević 2009] D. Mosić and D. S. Djordjević, “Moore–Penrose-invertible normal and Hermitian elements in rings”, *Linear Algebra Appl.* **431**:5-7 (2009), 732–745. MR 2011a:16058 Zbl 1186.16046
- [Putnam 1951] C. R. Putnam, “On normal operators in Hilbert space”, *Amer. J. Math.* **73**:2 (1951), 357–362. MR 12,717f Zbl 0042.34501
- [Strade 2004] H. Strade, *Simple Lie algebras over fields of positive characteristic, I: Structure theory*, Expositions in Math. **38**, De Gruyter, Berlin, 2004. MR 2005c:17025 Zbl 1074.17005
- [Strade and Farnsteiner 1988] H. Strade and R. Farnsteiner, *Modular Lie algebras and their representations*, Pure and Appl. Math. **116**, Dekker, New York, 1988. MR 89h:17021 Zbl 0648.17003
- [Yood 1974] B. Yood, “Commutativity properties in Banach $*$ -algebras”, *Pacific J. Math.* **53**:1 (1974), 307–317. MR 50 #14247 Zbl 0268.46050

Received June 14, 2011. Revised November 11, 2011.

ALEXANDRE N. GRISHKOV
INSTITUTO DE MATEMÁTICA E ESTATÍSTICA
UNIVERSIDADE DE SÃO PAULO
RUA DO MATÃO, 1010 CIDADE UNIVERSITÁRIA
CEP 05508-090 SÃO PAULO
BRAZIL
grishkov@ime.usp.br

MARINA RASSKAZOVA
OMSK INSTITUTE OF CONSUMER SERVICE TECHNOLOGY
UL. PEVZOVA 13, OMSK 644099
RUSSIA
marras123@gmail.com

SALVATORE SICILIANO
DIPARTIMENTO DI MATEMATICA E FISICA “ENNIO DE GIORGI”
UNIVERSITÀ DEL SALENTO
VIA PROVINCIALE LECCE-ARNESANO, 73100-LECCE
ITALY
salvatore.siciliano@unisalento.it

BOUNDED AND UNBOUNDED CAPILLARY SURFACES IN A CUSP DOMAIN

YASUNORI AOKI AND DAVID SIEGEL

We study asymptotic behavior of the height of a static liquid surface in a cusp domain as modelled by the Laplace–Young capillary surface equation. We introduce a new form of an asymptotic expansion in terms of the functions defining the boundary curves forming a cusp. We are able to address the asymptotic behavior of the capillary surface in cusp domains not previously considered, such as an exponential cusp. In addition, we have shown that the capillary surface in a cusp domain is bounded if the contact angles of the boundary walls forming a cusp are supplementary angles, which implies the continuity of the capillary surface at the cusp.

1. Introduction

Background. In everyday life, it is often safe to assume that the surface of water at rest is almost flat; however, careful observation shows that the surface of water in a container can exhibit complicated geometry near the interface where the water meets the container. One of the most extreme examples is when the container has a sharp (cusped) boundary. As seen in the photo, the static liquid surface (capillary surface) rises very steeply near a cusp—formed in the case illustrated here by the tangency between a circular cylinder and a straight wall. This behavior can be understood through a singular solution of the Laplace–Young capillary surface equation.

As noted in [Finn 1986], the study of a singular capillary surface can be traced back to Brook Taylor in 1712. Later contributions to the study of singular capillary surfaces by Concus and Finn [1969] and Miersemann [1993] spurred considerable interest in the field; see, for example [King et al. 1999; Scholz 2001; 2004; Norbury et al. 2005; Aoki 2007]. In particular, Scholz’s work on capillary surfaces in a



MSC2010: 35A20, 35C20, 35J60, 76B45.

Keywords: singularity, asymptotic analysis, nonlinear elliptic PDE.

domain containing a cusp where the boundaries can be approximated by power series (including fractional powers) led him to conclude that “[the capillary surface] rises with the same order [as] the order of contact of the two arcs, which form the cusp” [Scholz 2004]. Since this is a very intuitive statement, our curiosity led us to ask whether this statement holds for cases that Scholz did not consider in his paper [2004].

In this paper we extend Scholz’s results in two directions. We first consider cusp domains not limited to the power-law cusp. Instead of approximating the boundary by power series, we directly use the distance between two arcs forming a cusp in the asymptotic expansion. Although one may argue that most of the shapes used in real life applications can be approximated by power series, our main focus was to justify the above statement in a more direct and intuitive manner, by avoiding the extra approximation step. The second direction of extension is to include cases in which the contact angles of the boundary walls forming a cusp are supplementary angles. Although all the known results suggest that a capillary surface in a domain with a cusp is unbounded, we have shown that a capillary surface can be bounded, and hence continuous, if the contact angles are supplementary angles.

Statement of the problems. Here we state the problems we are going to consider in this paper. We first define a cusp domain. Without loss of generality, and for simplicity of writing, we consider the following domain (see Figure 1):

$$(1-1) \quad \Omega = \{(x, y) : x > 0, f_2(x) < y < f_1(x)\},$$

where

$$(1-2) \quad \begin{aligned} & f_1(x), f_2(x) \in C^3(0, \infty), \quad f_1(x) > f_2(x) \quad \text{for } x > 0, \\ & \lim_{x \rightarrow 0^+} f_1(x) = \lim_{x \rightarrow 0^+} f_2(x) = 0, \quad \lim_{x \rightarrow 0^+} f_1'(x) = \lim_{x \rightarrow 0^+} f_2'(x) = 0. \end{aligned}$$

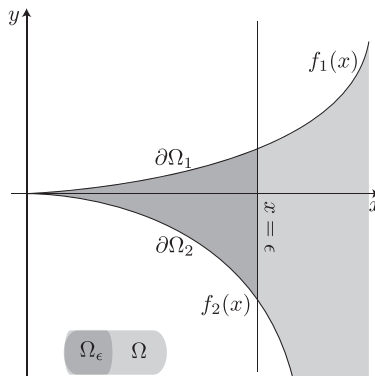


Figure 1. The cusped domain Ω and its boundary.

Also we denote the boundaries as follows:

$$\partial\Omega_1 = \{(x, y) : x > 0, y = f_1(x)\}, \quad \partial\Omega_2 = \{(x, y) : x > 0, y = f_2(x)\}.$$

Although we base our discussion on this infinite domain, all of the results presented in this paper only depend locally on a domain sufficiently close to the cusp, so the results hold for any domain that coincides with Ω in a neighborhood of the origin.

We now state the partial differential equation that interests us, the Laplace–Young capillary surface equation. Let $u(x, y)$ be the height of a capillary surface in domain Ω . It satisfies the following boundary value problem (see [Finn 1986] for a derivation):

$$(1-3) \quad \nabla \cdot Tu = \kappa u \quad \text{in } \Omega,$$

$$(1-4) \quad \vec{v}_1 \cdot Tu = \cos \gamma_1 \quad \text{on } \partial\Omega_1,$$

$$(1-5) \quad \vec{v}_2 \cdot Tu = \cos \gamma_2 \quad \text{on } \partial\Omega_2,$$

where

$$(1-6) \quad Tu = \frac{\nabla u}{\sqrt{1 + |\nabla u|^2}},$$

κ is the capillarity constant, \vec{v}_1 and \vec{v}_2 are exterior unit normal vectors on the boundaries $\partial\Omega_1$ and $\partial\Omega_2$, and γ_1, γ_2 are the contact angles. The capillarity constant κ can be normalized by rescaling x, y , and u . In the sequel we let $\kappa = 1$.

Here we introduce the big theta notation to replace the statement “is of the same order as”, to make this expression more precise. If $f(x) = \Theta(g(x))$, there exist constants $k_1, k_2 > 0$ and $x_0 > 0$ such that

$$(1-7) \quad k_1|g(x)| < |f(x)| < k_2|g(x)| \quad \text{for all } x < x_0.$$

We note that Θ is a more strict order relation than that of O , i.e., if $f(x) = \Theta(g(x))$ then $f(x) = O(g(x))$; however the converse is not true.

We can now write our core research questions as follows:

- Suppose $\gamma_1 + \gamma_2 \neq \pi$. Does $u(x, y) = \Theta\left(\frac{1}{f_1(x) - f_2(x)}\right)$ hold for any $f_1(x)$ and $f_2(x)$ satisfying (1-2)?
- How does $u(x, y)$ behave asymptotically as $x \rightarrow 0^+$ when $\gamma_1 + \gamma_2 = \pi$?

Structure of the paper. As the title of this paper suggests, there are two main parts: unbounded and bounded cases.

In Section 2 we consider unbounded capillary surfaces in cusp domains. We first prove in Section 2A that capillary surfaces are unbounded if $\gamma_1 + \gamma_2 \neq \pi$. Then in Section 2B the formal asymptotic expansion is presented. Using the formal asymptotic expansion, in Section 2C we prove the asymptotic behavior of the

solution. In Section 2D we give examples of power-law and non-power-law cusps with the intention of comparing our findings with the results in [Scholz 2004].

In Section 3 we consider bounded capillary surfaces in cusp domains. We first prove in Section 3A that capillary surfaces are bounded if $\gamma_1 + \gamma_2 = \pi$ and the curvature of the boundaries is finite. In Section 3B we show that if a capillary surface is bounded at the cusp, then it is continuous at the cusp. Section 4 contains concluding remarks summarizing our findings and suggesting some future extensions of our results. In addition, an Appendix we have included the Concus–Finn comparison principle and its Corollary used in Sections 2C and 3A.

2. Unbounded capillary surfaces

In this section, we assume $\gamma_1 + \gamma_2 \neq \pi$ and aim to prove that

$$u(x, y) = \Theta \left(\frac{1}{f_1(x) - f_2(x)} \right) \quad \text{as } x \rightarrow 0^+,$$

with as few restrictions on $f_1(x)$ and $f_2(x)$ as possible.

2A. Unboundedness of the capillary surface when $\gamma_1 + \gamma_2 \neq \pi$. We show that $u(x, y) \neq O(1)$. This is intuitively obvious from the remarkable result of Concus and Finn [1969], as a cusp can be considered as a corner with zero opening angle.

Lemma 2.1 (unboundedness of $u(x, y)$ when $\gamma_1 + \gamma_2 \neq \pi$). *Let $u(x, y)$ be the solution of the boundary value problem (1-3)–(1-5).*

If $\cos \gamma_1 + \cos \gamma_2 > 0$, then $u(x, y)$ cannot be bounded from above.

If $\cos \gamma_1 + \cos \gamma_2 < 0$, then $u(x, y)$ cannot be bounded from below.

Proof. Similar to the proof in [Concus and Finn 1969], we work by contradiction. First consider the case $\cos \gamma_1 + \cos \gamma_2 > 0$, and assume there exists a constant $M > 0$ such that $u(x, y) < M$ in Ω . Integrate the PDE (1-3) in a subdomain Ω_ϵ given by

$$\Omega_\epsilon = \{(x, y) : 0 < x < \epsilon, f_2(x) < y < f_1(x)\}.$$

By applying the divergence theorem and the boundary conditions (1-4) and (1-5), we obtain after some calculation the equation

$$\begin{aligned} (2-1) \quad & \int_{x=0}^\epsilon \int_{y=f_2(x)}^{f_1(x)} u \, dy \, dx \\ & = \int_{x=0}^\epsilon (\cos \gamma_1 \sqrt{1+f_1'^2} + \cos \gamma_2 \sqrt{1+f_2'^2}) \, dx + \int_{y=f_2(\epsilon)}^{f_1(\epsilon)} \frac{u_x}{\sqrt{1+u_x^2+u_y^2}} \Big|_{x=\epsilon} \, dx. \end{aligned}$$

The trick is to realize that the last term of (2-1) can be bounded from below, i.e.,

$$\frac{u_x}{\sqrt{1+u_x^2+u_y^2}} > -1,$$

which implies

$$\int_{y=f_2(\epsilon)}^{f_1(\epsilon)} \frac{u_x}{\sqrt{1+u_x^2+u_y^2}} \Big|_{x=\epsilon} dx > -(f_1(\epsilon) - f_2(\epsilon)).$$

We now apply the assumption $u(x, y) < M$ and the preceding inequality to (2-1) and obtain the inequality

$$\begin{aligned} \epsilon M \max_{0 < x \leq \epsilon} (f_1(x) - f_2(x)) + (f_1(\epsilon) - f_2(\epsilon)) \\ > \int_{x=0}^{\epsilon} (\cos \gamma_1 \sqrt{1+f_1'^2} + \cos \gamma_2 \sqrt{1+f_2'^2}) dx. \end{aligned}$$

Dividing both sides by $\epsilon > 0$ and taking the limit as ϵ approaches 0 gives

$$\begin{aligned} \lim_{\epsilon \rightarrow 0^+} M \max_{0 < x \leq \epsilon} (f_1(x) - f_2(x)) + \lim_{\epsilon \rightarrow 0^+} \frac{f_1(\epsilon) - f_2(\epsilon)}{\epsilon} \\ \geq \lim_{\epsilon \rightarrow 0^+} \frac{\int_{x=0}^{\epsilon} (\cos \gamma_1 \sqrt{1+f_1'^2} + \cos \gamma_2 \sqrt{1+f_2'^2}) dx}{\epsilon}. \end{aligned}$$

Applying the definition of the derivative together with (1-2) then gives

$$f_1'(0) - f_2'(0) \geq (\cos \gamma_1 \sqrt{1+f_1'(0)^2} + \cos \gamma_2 \sqrt{1+f_2'(0)^2}),$$

which implies $0 \geq \cos \gamma_1 + \cos \gamma_2$. Hence we obtain a contradiction. The proof for the case where $\cos \gamma_1 + \cos \gamma_2 < 0$ can be constructed similarly. \square

Lemma 2.1 and Corollary A.1 together imply that $u(x, y)$ is unbounded at the cusp and bounded away from the cusp.

2B. Formal asymptotic expansion of the boundary value problem (1-3)–(1-5).

The main idea is to consider an asymptotic expansion of the form

$$(2-2) \quad v(x, y) = \frac{A}{f_1(x) - f_2(x)} + g(x, y) \frac{f_1'(x) - f_2'(x)}{f_1(x) - f_2(x)} + h(x, y) \frac{(f_1'(x) - f_2'(x))^2}{f_1(x) - f_2(x)},$$

where $g(x, y), h(x, y) \in O(1)$ as $x \rightarrow 0^+$. Recalling that $\lim_{x \rightarrow 0^+} f_1(x) = 0$ and $\lim_{x \rightarrow 0^+} f_2(x) = 0$, we have the first term significantly larger than the second term near the cusp. Also note that the leading order term is of the same order as the reciprocal of the distance between two boundaries measured in \vec{y} direction.

The aim of this subsection is to find $g(x, y)$ and $h(x, y)$ such that (2-2) satisfies asymptotically the PDE (1-3) and the boundary conditions (1-4) and (1-5).

For simplicity of computation, we introduce coordinate variables s and t as follows:

$$s := x, \quad t := \frac{2y - (f_1(x) + f_2(x))}{f_1(x) - f_2(x)}.$$

We have chosen t so that $y = f_1(x)$ when $t = 1$, and $y = f_2(x)$ when $t = -1$.

Lemma 2.2 (first two terms of the formal asymptotic expansion). *In (2-2), let $A = \cos \gamma_1 + \cos \gamma_2$, and*

$$g(s, t) = -\sqrt{1 - \left(\frac{\cos \gamma_1(t+1) + \cos \gamma_2(t-1)}{2} \right)^2} + C_1$$

(where C_1 is an arbitrary constant), and $h(s, t) = 0$. If $f_1(s)$ and $f_2(s)$ satisfy

$$(2-3) \quad \begin{aligned} f_1(s) - f_2(s) &= o(f_1'(s) - f_2'(s)), & \frac{f_1''(s) - f_2''(s)}{f_1(s) - f_2(s)} &= o\left(\frac{f_1'(s) - f_2'(s)}{(f_1(s) - f_2(s))^2}\right), \\ \frac{f_1'''(s) - f_2'''(s)}{f_1'(s) - f_2'(s)} &= o\left(\frac{1}{(f_1(s) - f_2(s))^2}\right), \end{aligned}$$

as $s \rightarrow 0^+$, then

$$(2-4) \quad \begin{aligned} \vec{v}_1 \cdot Tv|_{t=1} &= \cos \gamma_1 + o(1), & \vec{v}_2 \cdot Tv|_{t=-1} &= \cos \gamma_2 + o(1), \\ \nabla \cdot Tv - v &= o\left(\frac{1}{f_1(s) - f_2(s)}\right) \end{aligned}$$

as $s \rightarrow 0^+$.

A tedious but straightforward calculation will verify this lemma. Instead of showing this calculation, we briefly explain here how the expressions for A , g , and h in the statement of the lemma were deduced. We first let

$$v(s, t) = \frac{A}{f_1(s) - f_2(s)} + g(t) \frac{f_1'(s) - f_2'(s)}{f_1(s) - f_2(s)}.$$

(It is desirable—and, as it turns out, possible—to make the function g depend only on t , so we will suppress the dependence of g on s ; the same applies to the function h .) After some lengthy calculations with assumptions (2-3) we obtain

$$\begin{aligned} \vec{v}_1 \cdot Tv|_{t=1} &= \frac{2g'(1)}{\sqrt{A^2 + 4g'^2(1)}} + o(1), & \vec{v}_2 \cdot Tv|_{t=-1} &= -\frac{2g'(-1)}{\sqrt{A^2 + 4g'^2(-1)}} + o(1), \\ \nabla \cdot Tv - v &= \left(\frac{4g''(t)A^2}{(A^2 + 4g'^2(t))^{3/2}} - A \right) \frac{1}{f_1(s) - f_2(s)} + o\left(\frac{1}{f_1(s) - f_2(s)}\right). \end{aligned}$$

We now impose the desired equalities (2-4) and obtain a nonlinear ordinary differential equation of the first order in $g'(t)$,

$$(2-5) \quad \frac{4g''(t)A^2}{(A^2 + 4g'^2(t))^{3/2}} = A \quad \text{for } -1 < t < 1,$$

with boundary conditions

$$(2-6) \quad \frac{2g'(1)}{\sqrt{A^2 + 4g'^2(1)}} = \cos \gamma_1, \quad -\frac{2g'(-1)}{\sqrt{A^2 + 4g'^2(-1)}} = \cos \gamma_2.$$

Though there are two boundary conditions for this first-order ODE, note that A is an indeterminate constant. Both $g'(t)$ and A are determined by first integrating (2-5) under the boundary conditions (2-6). One essential observation from this derivation is that the coefficient A of the leading-order term was found together with that of the second-order term, $g(t)$. In fact this pattern continues; the constant on the second-order term C_1 will be determined (it vanishes) at the same time as the third-order term of the formal asymptotic expansion is found.

Lemma 2.3 (first three terms of the formal asymptotic expansion). *In (2-2), let $A = \cos \gamma_1 + \cos \gamma_2$,*

$$g(t) = -\sqrt{1 - \left(\frac{\cos \gamma_1(t+1) + \cos \gamma_2(t-1)}{2} \right)^2},$$

and

$$h(t) = -\frac{A}{4} \left(\delta t + \frac{t^2}{2} \right) + \frac{1-\alpha}{2A} g(t)^2 + C_2,$$

where C_2 is an arbitrary constant. If $f_1(s)$ and $f_2(s)$ satisfy the conditions

$$(2-7) \quad f_1'(s) > f_2'(s) \quad \text{for } s > 0,$$

$$(2-8) \quad f_1(s) - f_2(s) = o(f_1'(s) - f_2'(s)),$$

$$(2-9) \quad \frac{f_1''(s) - f_2''(s)}{f_1(s) - f_2(s)} = \alpha \frac{(f_1'(s) - f_2'(s))^2}{(f_1(s) - f_2(s))^2} + o\left(\frac{(f_1'(s) - f_2'(s))^2}{(f_1(s) - f_2(s))^2}\right),$$

$$(2-10) \quad \frac{f_1'''(s) - f_2'''(s)}{f_1'(s) - f_2'(s)} = O\left(\frac{(f_1'(s) - f_2'(s))^2}{(f_1(s) - f_2(s))^2}\right),$$

$$(2-11) \quad f_1'(s) + f_2'(s) = \delta(f_1'(s) - f_2'(s)) + o(f_1'(s) - f_2'(s)),$$

$$(2-12) \quad f_1''(s) + f_2''(s) = O(f_1''(s) - f_2''(s)),$$

as $s \rightarrow 0^+$, where $\alpha, \delta \in \mathbb{R}$, then

$$\vec{v}_1 \cdot Tv|_{t=1} = \cos \gamma_1 + o(f_1'(s) - f_2'(s)), \quad \vec{v}_2 \cdot Tv|_{t=-1} = \cos \gamma_2 + o(f_1'(s) - f_2'(s)),$$

$$\nabla \cdot Tv - v = o\left(\frac{f_1'(s) - f_2'(s)}{f_1(s) - f_2(s)}\right)$$

as $s \rightarrow 0^+$.

Again, a long tedious calculation will prove this lemma. We followed similar steps to determine $h(t)$, although solving the differential equation for $h(t)$ was not nearly as straightforward as for $g(t)$. The constant C_1 was determined to be 0 when $h(t)$ was determined and a new unknown constant C_2 appeared in the third-order term.

Comparing assumptions (2-3) with assumptions (2-8)–(2-12), we can see that the restrictions on f_1 and f_2 increase as the number of terms in the formal asymptotic expansion increases from two terms to three terms. Although these assumptions are not proven to be necessary conditions for these lemmas to hold, it is our suspicion that as the number of the terms in the asymptotic expansion increases, the restrictions on f_1 and f_2 do become more strict.

2C. Asymptotic behavior of the capillary surface. The main result of Section 2 is stated and proven in this subsection. We first show that the asymptotic growth order of the solution is the same order as the reciprocal of the distance between two arcs forming a cusp.

Theorem 2.1 (growth order of $u(x, y)$). *Let $u(x, y)$ be the solution of the boundary value problem (1-3)–(1-5). If $f_1(s)$ and $f_2(s)$ satisfy the conditions (2-3) and $|\cos \gamma_1| \neq 1$ and $|\cos \gamma_2| \neq 1$, then there exist positive constants s_0, k_1 and k_2 such that*

$$(2-13) \quad k_2 \left(\frac{1}{f_1(s) - f_2(s)} \right) < |u(s, t)| < k_1 \left(\frac{1}{f_1(s) - f_2(s)} \right), \quad \text{for } s < s_0.$$

Proof. The main idea of our proof is to construct a supersolution and a subsolution by modifying the formal asymptotic expansion given in Lemma 2.2. We prove these modified equations are in fact supersolution and subsolution by applying the Concus–Finn comparison principle (Theorem A.1). Let

$$v(s, t; K_1, K_2) = \frac{A(K_1)}{f_1(s) - f_2(s)} + g(t; K_1) \frac{f'_1(s) - f'_2(s)}{f_1(s) - f_2(s)} + K_2,$$

where

$$(2-14) \quad g(t; K_1) = -\frac{A}{A - \frac{1}{3}K_1} \sqrt{1 - \left(\frac{\cos \gamma_1(t+1) + \cos \gamma_2(t-1)}{2} - \frac{K_1}{6}t \right)^2};$$

here we choose K_1 and K_2 appropriately to construct the supersolution and the subsolution. The trick of this proof is to realize that A and $g(t)$, the first and second terms of the formal asymptotic expansion, need to be modified to obtain a supersolution and a subsolution. We first impose the following conditions on K_1

so that the quantities in (2-14) behave reasonably:

$$(2-15) \quad |K_1| < |\cos \gamma_1 + \cos \gamma_2|,$$

$$(2-16) \quad |K_1| < 6(1 - |\cos \gamma_1|),$$

$$(2-17) \quad |K_1| < 6(1 - |\cos \gamma_2|).$$

We restrict the choice of K_1 so that the sign of $A(K_1)$ only depends on the sign of $\cos \gamma_1 + \cos \gamma_2$. Also, if K_1 is chosen to satisfy (2-15)–(2-17), then $g(t, K_1)$ is real and bounded. After some calculations assuming (2-3), we obtain

$$(2-18) \quad \vec{v}_1 \cdot Tv|_{t=1} = \cos \gamma_1 + \frac{1}{3}K_1 + o(1), \quad \vec{v}_2 \cdot Tv|_{t=-1} = \cos \gamma_2 + \frac{1}{3}K_1 + o(1),$$

$$(2-19) \quad \nabla \cdot Tv - v = -\frac{1}{3}K_1 \frac{1}{f_1(s) - f_2(s)} - K_2 + o\left(\frac{f'_1(s) - f'_2(s)}{f_1(s) - f_2(s)}\right),$$

as $s \rightarrow 0^+$. The essential observation in this step of the proof is that the expressions in (2-18) do not depend on K_2 including the “small o” terms. Similarly, (2-19) has K_2 dependence only at the second term and not in the “small o” term.

We now construct a function v^+ that satisfies inequalities (A-1)–(A-4) in the Appendix, and is therefore a supersolution. We denote the associated constants by K_1^+ and K_2^+ ; i.e., $v^+ = v(s, t; K_1^+, K_2^+)$. Firstly, K_1^+ are chosen to be a small enough positive real number so as to satisfy (2-15)–(2-17). Then we choose a constant $s_0^+ > 0$ so that for all $s < s_0^+$ the inequalities

$$(2-20) \quad \vec{v}_1 \cdot Tv^+|_{t=1} - \cos \gamma_1 > 0, \quad \vec{v}_2 \cdot Tv^+|_{t=-1} - \cos \gamma_2 > 0,$$

$$(2-21) \quad \nabla \cdot Tv^+ - v^+ + K_2^+ < 0.$$

are satisfied. Based on our previous observation we note that the choice of s_0^+ is independent of K_2^+ . Let Ω_0^+ be the subdomain of Ω such that $s < s_0^+$. By adding a restriction on K_2^+ to be a positive real number, it follows from (2-21) that

$$\nabla \cdot Tv^+ - v^+ < 0 \quad \text{in } \Omega_0^+.$$

Note that v^+ now satisfies conditions (A-1)–(A-3) of the Concus–Finn comparison principle (Theorem A.1). It remains to choose K_2^+ so as to satisfy condition (A-4). According to Corollary A.1, $u(s, t)$ is bounded at $s = s_0^+$. Hence there exists K_2^+ such that

$$v^+ > u \quad \text{on } s = s_0^+.$$

Thus by Theorem A.1 we have shown that there exists Ω_0^+, K_1^+, K_2^+ such that

$$v^+(s, t; K_1^+, K_2^+) > u(s, t) \quad \text{in } \Omega_0^+.$$

Similarly we can construct a subsolution $v^-(s, t; K_1^-, K_2^-)$ such that

$$v^-(s, t; K_1^-, k_2^-) < u(s, t) \quad \text{in } \Omega_0^-.$$

Hence in $\Omega_0^+ \cap \Omega_0^-$ we have $v^- < u < v^+$, i.e.,

$$\frac{A(K_1^-)}{f_1(s) - f_2(s)} + g(t; K_1^-) \frac{f_1'(s) - f_2'(s)}{f_1(s) - f_2(s)} + K_2^- < u$$

and

$$u < \frac{A(K_1^+)}{f_1(s) - f_2(s)} + g(t; K_1^+) \frac{f_1'(s) - f_2'(s)}{f_1(s) - f_2(s)} + K_2^+.$$

Since K_1^+ and K_1^- were chosen to satisfy (2-15), $A(K_1^+)$ and $A(K_1^-)$ have the same sign. Without loss of generality assume $A(K_1^+) > 0$. Let

$$m_1(s) = A(K_1^+) + \left(\max_{-1 < t < 1} \{g(t; K_1^+)(f_1'(s) - f_2'(s))\} + K_2^+(f_1(s) - f_2(s)) \right),$$

$$m_2(s) = A(K_1^-) + \left(\min_{-1 < t < 1} \{g(t; K_1^-)(f_1'(s) - f_2'(s))\} + K_2^-(f_1(s) - f_2(s)) \right).$$

Since $f_1'(s) - f_2'(s)$ and $f_1(s) - f_2(s)$ are $o(1)$ and continuous, there exists $s_0 > 0$ so that $m_1(s), m_2(s) > 0$ for $s < s_0$. By choosing

$$(2-22) \quad k_1 = \max_{0 < s < s_0} m_1(s), \quad k_2 = \min_{0 < s < s_0} m_2(s),$$

we obtain (2-13). □

Note that the proof holds for arbitrarily small $|K_1^\pm|$. Hence it is natural to guess that $(\cos \gamma_1 + \cos \gamma_2)/(f_1(s) - f_2(s))$ is the correct leading-order term of the asymptotic expansion. We now show that the leading-order term of the formal asymptotic expansion is in fact the first-order term of the asymptotic expansion of $u(s, t)$.

Theorem 2.2 (leading-order behavior of $u(x, y)$). *Let $u(x, y)$ be the solution of the boundary value problem (1-3)–(1-5). Assume that $f_1(s)$ and $f_2(s)$ satisfy the conditions (2-8)–(2-12). Then*

$$(2-23) \quad u(s, t) = \frac{\cos \gamma_1 + \cos \gamma_2}{f_1(s) - f_2(s)} + O\left(\frac{f_1'(s) - f_2'(s)}{f_1(s) - f_2(s)}\right) \quad \text{as } s \rightarrow 0^+.$$

Proof. We let

$$v(s, t; K_3, K_4, K_5) = \frac{A}{f_1(s) - f_2(s)} + g(t, K_3) \frac{f_1'(s) - f_2'(s)}{f_1(s) - f_2(s)} + h(t; K_4) \frac{(f_1'(s) - f_2'(s))^2}{f_1(s) - f_2(s)} + K_5,$$

where

$$A = \cos \gamma_1 + \cos \gamma_2,$$

$$g(t; K_3) = -\sqrt{1 - \left(\frac{\cos \gamma_1(t+1) + \cos \gamma_2(t-1)}{2} \right)^2} + K_3,$$

$$h(t; K_4) = -\frac{A}{4} \left(\delta t + \frac{t^2}{2} \right) + \frac{1-\alpha}{2A} \left\{ 1 - \left(\frac{\cos \gamma_1(t+1) + \cos \gamma_2(t-1)}{2} \right)^2 \right\} + \frac{K_4}{2} t^2.$$

Unlike the proof of Theorem 2.1, we can choose K_3 and K_4 as any real numbers.

After some calculations assuming (2-8)–(2-12), we obtain

$$(2-24) \quad \vec{v}_1 \cdot Tv|_{t=1} = \cos \gamma_1 + K_4 \frac{(f'_1(s) - f'_2(s))}{(A^2 + 4(g'(t))^2)^{3/2}} + o(f'_1(s) - f'_2(s)),$$

$$(2-25) \quad \vec{v}_2 \cdot Tv|_{t=-1} = \cos \gamma_2 + K_4 \frac{(f'_1(s) - f'_2(s))}{(A^2 + 4(g'(t))^2)^{3/2}} + o(f'_1(s) - f'_2(s)),$$

$$(2-26) \quad \begin{aligned} \nabla \cdot Tv - v = & \left\{ \left(-\frac{12g'(t)t}{A^2 + 4(g'(t))^2} + \frac{4A^2}{(A^2 + 4(g'(t))^2)^{3/2}} \right) K_4 - K_3 \right\} \frac{f'_1(s) - f'_2(s)}{f_1(s) - f_2(s)} \\ & - K_5 + o\left(\frac{f'_1(s) - f'_2(s)}{f_1(s) - f_2(s)} \right), \end{aligned}$$

as $s \rightarrow 0^+$.

We now construct a supersolution. Let v^+ denote the supersolution, with associate constants K_3^+ , K_4^+ , K_5^+ ; i.e., $v^+ = v(s, t; K_3^+, K_4^+, K_5^+)$. We first choose the positive constant K_4^+ arbitrarily. Then we choose K_3^+ big enough so that

$$\left\{ \left(-\frac{12g'(t)t}{A^2 + 4(g'(t))^2} + \frac{4A^2}{(A^2 + 4(g'(t))^2)^{3/2}} \right) K_4^+ - K_3^+ \right\} < 0 \quad \text{for } -1 < t < 1.$$

We now choose $s_2^+ > 0$ so that

$$\vec{v}_1 \cdot Tv|_{t=1} - \cos \gamma_1 > 0, \quad \vec{v}_2 \cdot Tv|_{t=-1} - \cos \gamma_2 > 0, \quad \nabla \cdot Tv - v + K_5^+ < 0$$

for $0 < s < s_2^+$. Let Ω_2^+ be the subdomain of Ω such that $s < s_2^+$. By Corollary A.1, we know that $u(s_2^+, t)$ is bounded; hence there exists a large enough positive constant K_5^+ so that

$$v^+ > u \quad \text{on } s = s_2^+.$$

Thus by the Concus–Finn comparison principle (Theorem A.1) we have

$$v^+ > u \quad \text{in } \Omega_2^+.$$

Similarly we can construct a subsolution v^- by choosing suitable K_3^-, K_4^-, K_5^- and s_2^- . Thus we can bound the solution $u(s, t)$ by v^- and v^+ ; i.e.,

$$v^- < u < v^+ \quad \text{in } \Omega_2^+ \cap \Omega_2^-,$$

and (2-23) holds. □

From this section, we conclude that the height of a capillary surface near a cusp is proportional to the reciprocal of the distance between the two arcs forming the cusp, assuming these arcs satisfy (2-3).

2D. Examples of cusp domains. In the previous subsection, we have shown the behavior of the capillary surface near a cusp under certain assumptions $f_1(x)$ and $f_2(x)$ giving the shape of the boundaries. Those assumptions, expressed by (2-3) or (2-8)–(2-12), are left in these forms in order to make the theorem as general as possible. On the other hand, it is hard to grasp what kind of cusps are allowed or not. In this subsection, we will show through examples when the theorem is applicable and when it is not.

It is easy to show that if the difference between f_1 and f_2 can be written in the following form, these functions satisfy (2-8)–(2-10):

$$(2-27) \quad f_1(x) - f_2(x) = c x^{a_0} \exp\left(\sum_{i=1}^{\infty} a_i x^{b_i}\right),$$

where $c > 0, a_1 < 0, b_1 < 0, b_{i+1} > b_i$. An alternative way to write this is

$$(2-28) \quad f_1(x) - f_2(x) = \exp\left(\int_c^x \frac{\sum_{i=0}^{\infty} \tilde{a}_i \zeta^{\tilde{b}_i}}{\sum_{i=0}^{\infty} a_i \zeta^{b_i}} d\zeta\right),$$

where $c > 0, b_0 - \tilde{b}_0 \geq 1, b_{i+1} > b_i, a_0 > 0$ and $\tilde{a}_0 > 0$. As (2-8)–(2-10) are stricter requirements for $f_1(x)$ and $f_2(x)$ than (2-3), if $f_1(x) - f_2(x)$ can be written as (2-27) or (2-28), then f_1 and f_2 satisfy (2-3).

Note that (2-11) and (2-12) can be interpreted as saying that some osculating cusps (cusps with boundaries tangent to second order) are not allowed, and Equation (2-7) can be interpreted as saying that infinitely oscillating cusp boundaries are not allowed.

Example 1 (fractional power cusp). We now consider a cusp that can be analyzed through the result of Scholz. Consider (2-28) and let $b_0 > 1, \tilde{a}_i = a_i b_i$ and $\tilde{b}_i = b_i - 1$. Then we have

$$(2-29) \quad f_1(x) - f_2(x) = \tilde{c} \sum_{i=0}^{\infty} a_i x^{b_i}.$$

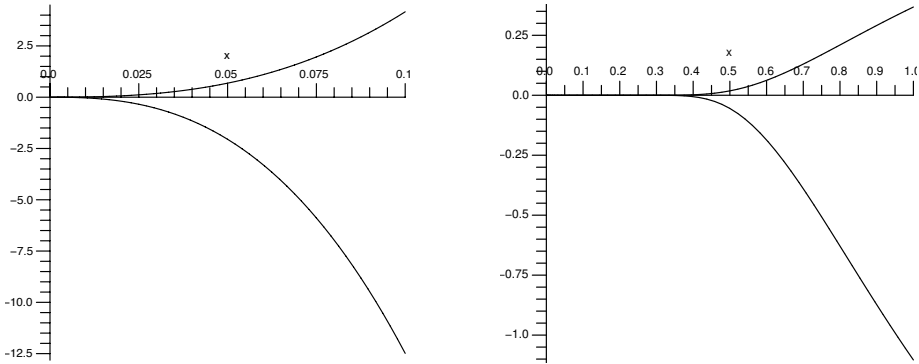


Figure 2. Left: fractional power cusp (Example 1). Right: exponential cusp (Example 2). In both cases, $p = 1$ and $q = -3$.

To be more specific, we consider the cusp boundaries

$$(2-30) \quad f_1(x) = p(x^{5/2} + x^3), \quad f_2(x) = q(x^{5/2} + x^3),$$

with constants $p > q$ (see Figure 2, left). According to Theorem 2.2, we obtain the asymptotic expansion

$$\begin{aligned} u(x, y) &= \frac{\cos \gamma_1 + \cos \gamma_2}{(p - q)(x^{5/2} + x^3)} + O(x^{-1}) \\ &= \frac{\cos \gamma_1 + \cos \gamma_2}{p - q} \left(\frac{1}{x^{5/2}} - \frac{1}{x^2} + \frac{1}{x^{3/2}} \right) + O(x^{-1}) \end{aligned}$$

as $x \rightarrow 0^+$. We note that this result is consistent with that of Scholz. It is noteworthy that by finding the first order term of our asymptotic expansion we find the first three terms of the asymptotic series solution in power series.

Example 2 (exponential cusp). We now consider cusps to which the results of Scholz do not apply. Equation (2-27) implies that $f_1(x)$ and $f_2(x)$ can contain exponential terms. We now consider a very sharp cusp, an “exponential cusp”, where

$$f_1(x) = p e^{-1/x^2}, \quad f_2(x) = q e^{-1/x^2}.$$

with constants $p > q$ (see Figure 2, right). According to Theorem 2.2, we obtain the asymptotic expansion

$$u(x, y) = \frac{\cos \gamma_1 + \cos \gamma_2}{p - q} e^{1/x^2} + O(x^{-3}) \quad \text{as } x \rightarrow 0^+.$$

This example shows that our result has extended the result of Scholz on the leading order behavior of a capillary surface in a cusp domain.

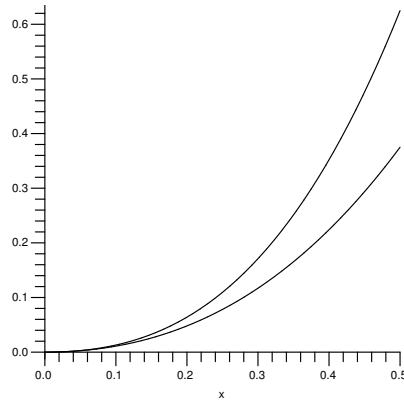


Figure 3. Osculatory cusp ($p = 3, q = 1$).

Example 3 (osculatory cusp). We now consider a case where Theorem 2.2 cannot be applied. Consider the cusp boundaries

$$(2-31) \quad f_1(x) = x^2 + px^3, \quad f_2(x) = x^2 + qx^3,$$

with constants $p > q$ (see Figure 3).

These functions do not satisfy (2-11)–(2-12); hence Theorem 2.2 does not apply. On the other hand, if $|\cos \gamma_1| \neq 1$ and $|\cos \gamma_2| \neq 1$, Theorem 2.1 applies, as this f_1 and f_2 satisfy (2-3). Hence even the case of the osculating cusp, we have shown that the height of the capillary surface rises as the same order as the reciprocal of the distance of two arcs forming a cusp, i.e.,

$$(2-32) \quad u(x, y) = \Theta \left(\frac{1}{x^3} \right).$$

As the two functions f_1 and f_2 forming a cusp only appear as $(f_1(x) - f_2(x))$ or $(f_1'(x) - f_2'(x))$ in the asymptotic expansion (2-2), it is not immediately obvious as to why we cannot conduct the asymptotic analysis of this problem similarly to the case where $f_1(x) = px^3, f_2(x) = qx^3$. However, the difference in asymptotic order between $f_1(x) - f_2(x)$ on the one hand and $f_1(x)$ or $f_2(x)$ on the other becomes crucial in calculating the asymptotic relations (2-24)–(2-26) of the boundary conditions and the PDE. For example, for the calculation of (2-24), since

$$\vec{v}_1 = \frac{(-f_1'(x), 1)}{\sqrt{1 + (f_1'(x))^2}},$$

the function $f_1(x)$ appears without subtracting $f_2(x)$. As a result, the asymptotic relation (2-24) does not hold for the case of osculatory cusp. Thus for the osculatory cusps, we cannot use the asymptotic expansion (2-2) to prove the leading order behavior.

3. Bounded capillary surfaces

In this section we assume $\gamma_1 + \gamma_2 = \pi$ and prove that $u(x, y)$ is bounded.

3A. Proof of the boundedness of the capillary surface when $\gamma_1 + \gamma_2 = \pi$.

Theorem 3.1 (boundedness of $u(x, y)$ when $\gamma_1 + \gamma_2 = \pi$). *Let $u(x, y)$ be the solution of the boundary value problem (1-3)–(1-4) with $\gamma_1 = \gamma$ and $\gamma_2 = \pi - \gamma$. If the boundaries $\partial\Omega_1$ and $\partial\Omega_2$ have finite curvatures in the neighborhood of the cusp, in other words, if there exists ϵ_o such that*

$$(3-1) \quad f_1(x), f_2(x) \in C^2([0, \epsilon_o]),$$

then $u(x, y)$ is bounded.

Proof. It follows immediately from Corollary A.1 that $u(x, y)$ is bounded in the domain away from the origin. Hence our problem reduces to show that $u(x, y)$ is bounded in the neighborhood of the origin.

First we show that $u(x, y)$ is bounded above at the origin by using the Concus–Finn comparison principle (Theorem A.1). In order to apply Theorem A.1, we need to construct a surface that satisfies (A-1)–(A-4). The most difficult part of this proof is to construct a surface that satisfies both (A-2) and (A-3). Our unique idea is to construct a surface that satisfies (1-4) exactly hence (A-2) and also satisfies (A-3). Such surface can be constructed by a surface with contour lines parallel to the boundary $\partial\Omega_1$. In other words by letting the height of the surface only depends on the distance from the boundary $\partial\Omega_1$, we can easily construct a surface with exact constant contact angle γ on this boundary. We choose a surface so that the height and the mean curvature is bounded so that Inequalities (A-1) and (A-4) can easily be satisfied by shifting this surface upwards.

We now translate the above statement to the precise language of mathematics. Without loss of generality we assume $0 \leq \gamma \leq \pi/2$. First we define a coordinate system such that the one family of the coordinate curves is parallel curves of the boundary $\partial\Omega_1$ and another family of the coordinate curves is lines perpendicular to the boundary $\partial\Omega_1$. Let s and t be new coordinate variables defined implicitly as the following (note that s here has different meaning from s used in Section 2):

$$(3-2) \quad (x, y) = (s, f_1(s)) - t \vec{v}_1(s),$$

where $\vec{v}_1(s)$ is the exterior unit normal vector of the boundary $\partial\Omega_1$ at $(s, f_1(s))$. More explicitly, the coordinate variables of Cartesian coordinate system x and y can be written using the new coordinate variables s and t as follows:

$$(3-3) \quad x = s + t \frac{f_1'(s)}{\sqrt{1 + (f_1'(s))^2}}, \quad y = f_1(s) - t \frac{1}{\sqrt{1 + (f_1'(s))^2}}.$$

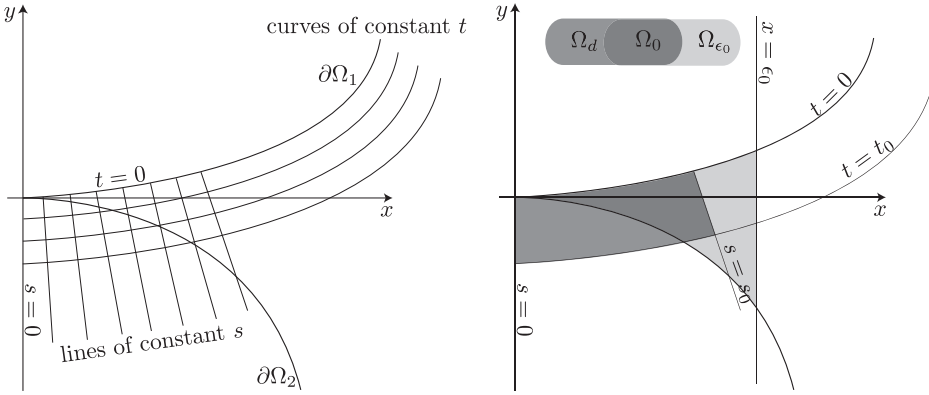


Figure 4. Left: coordinate lines of the s - t coordinate system. Right: the domain Ω_0 .

The variable t can be interpreted as the distance of the point from the boundary $\partial\Omega_1$. The coordinate curves are sketched in Figure 4, left.

The Jacobian of (3-3) is calculated to be

$$\frac{\partial(x, y)}{\partial(s, t)} = \frac{f_1'(s)^2 - 1}{\sqrt{1 + (f_1'(s))^2}} \left(1 + t \frac{f_1''(s)}{(1 + (f_1'(s))^2)^{3/2}} \right).$$

This gives that the point (x, y) in the Cartesian coordinate system can be specified uniquely by the new coordinate variables (s, t) defined by (3-3) if both

$$(3-4) \quad f_1'(s)^2 - 1 \neq 0$$

and

$$(3-5) \quad 1 + t \frac{f_1''(s)}{(1 + (f_1'(s))^2)^{3/2}} \neq 0.$$

Since $f_1(s) \in C^2([0, \epsilon_0])$ and $\lim_{s \rightarrow 0^+} f_1(s) = 0$, there exists $0 < s_0 \leq \epsilon_0$ so that (3-4) is satisfied for all $s \in [0, s_0]$. Also due to the smoothness of $f_1(s)$, we can find $t_0 > 0$ such that (3-5) holds for all $t \in [0, t_0]$ in $s \in [0, s_0]$. That is to say, the coordinate system defined in (3-3) is valid in the domain

$$\Omega_d := \{(s, f_1(s)) - t \vec{v}_1(s) \in \mathbb{R}^2 : 0 \leq s \leq s_0, 0 \leq t \leq t_0\}.$$

Then we choose the subdomain

$$\Omega_0 := \Omega_d \cap \Omega_{\epsilon_0},$$

where $\Omega_{\epsilon_0} := \{(x, y) \in \mathbb{R}^2 : 0 < x < \epsilon_0, f_2(x) < y < f_1(x)\}$, as depicted in Figure 4, right. Since $\bar{\Omega}_0$ contains the cusp at the origin, finding an upper bound for the surface u in domain Ω_0 by using Theorem A.1 would prove that the capillary surface

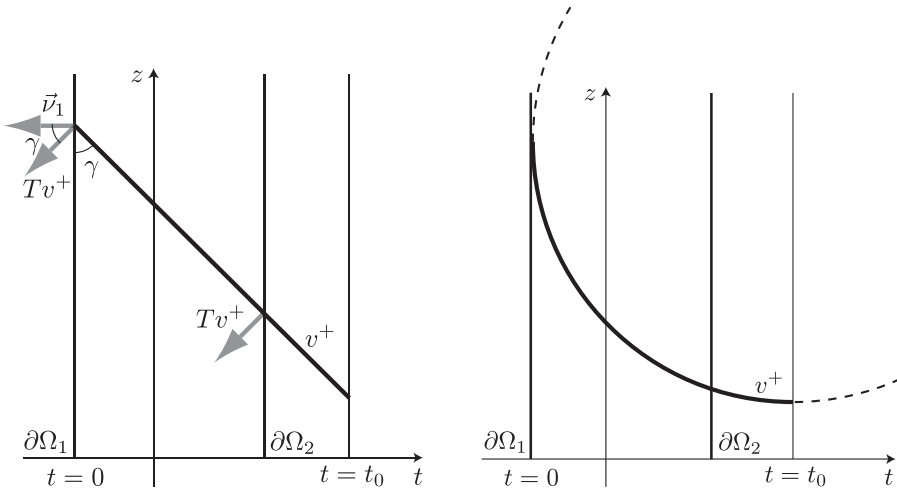


Figure 5. Cross section of a surface $v^+(s, t)$ on the line of constant s : Choice of function $g(t)$ for $\gamma \neq 0$ (left) and for $\gamma = 0$ (right).

is bounded above at the cusp. Using the parameters t and s , we now construct a surface $v^+(s, t)$ in Ω_0 , with components (x, y, z) , as follows:

(3-6)

$$x(s, t) = s + t \frac{f'_1(s)}{\sqrt{1+(f'_1(s))^2}}, \quad y(s, t) = f_1(s) - t \frac{1}{\sqrt{1+(f'_1(s))^2}}, \quad z(s, t) = g(t).$$

The choice of the height function $g(t)$ depends on the contact angle γ . In our opinion, the simplest choice such that the surface v^+ satisfies (1-4) exactly and also satisfies (A-3) is

(3-7)

$$g(t) = \begin{cases} -\cot \gamma t + K & \text{for } \gamma \neq 0, \\ -\sqrt{t_0^2 - (t - t_0)^2} + K & \text{for } \gamma = 0, \end{cases}$$

where K is a constant that we will specify later. The cross section of this surface on a line of constant s is depicted in Figure 5, left.

The surface $v^+(s, t)$ can be sketched as in Figure 6. For example, if the curve $\partial\Omega_1$ is a part of a circle, then the surface $v^+(s, t)$ for the case $\gamma \neq 0$ becomes a part of a cone, and for the case $\gamma = 0$ it becomes a part of a torus.

We now verify that the surface $v^+(s, t)$ satisfies (1-4) exactly and also satisfies (A-3). We first consider the case $\gamma \neq 0$, as the vector Tv^+ can be interpreted as a unit downwards vector of the surface v^+ , it follows immediately from Figure 5 (left) that $Tv^+(s, t)$ can be written as

$$Tv^+ = \cos \gamma \vec{v}_1 - \sin \gamma \hat{z},$$

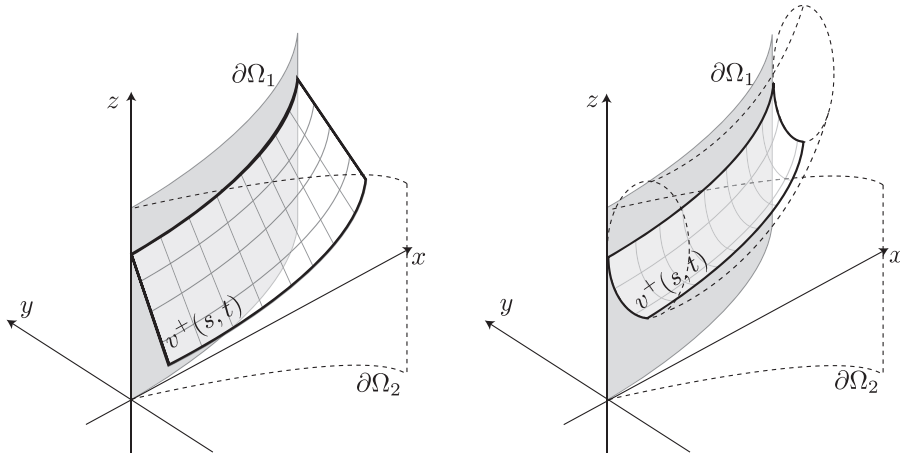


Figure 6. Sketch of the surface $v^+(s, t)$ for $\gamma \neq 0$ (left) and for $\gamma = 0$ (right).

where \hat{z} is a unit vector in z direction. Noting that the vector \vec{v}_1 is orthogonal to \hat{z} , we obtain that (1-4) is satisfied exactly by the surface $v^+(s, t)$, i.e.,

$$\vec{v}_1 \cdot T v^+ = \cos \gamma \quad \text{on } \partial\Omega_1 \cap \partial\Omega_0.$$

We now verify that the surface $v^+(s, t)$ satisfies Inequality (A-3). By noticing \vec{v}_2 and \hat{z} are orthogonal and both \vec{v}_1 and \vec{v}_2 are unit vectors, we obtain the inequality

$$\vec{v}_2 \cdot T v^+ = \cos \gamma \vec{v}_1 \cdot \vec{v}_2, > -\cos \gamma, = \cos(\pi - \gamma).$$

Although the case of $\gamma = 0$ may look complicated, it follows immediately from Figure 5 (right) that the angle between the unit downward normal vector of v^+ and \vec{v}_1 are parallel on the boundary, on $\partial\Omega_1 \cap \partial\Omega_0$,

$$\vec{v}_1 \cdot T v^+ = 1 = \cos 0.$$

Also it follows immediately from the definition of the differential operator T that $|T v^+| \leq 1$; see (1-6). By noting that \vec{v}_2 is a unit vector, i.e., $|\vec{v}_2| = 1$, we have

$$\vec{v}_2 \cdot T v^+ > -1 = \cos(\pi - 0).$$

Hence the surface $v^+(s, t)$ defined by (3-6)–(3-7) satisfies Inequalities (A-2) and (A-3). We now show that the surface $v^+(s, t)$ satisfies (A-1) by choosing large enough constant K .

Since $\nabla \cdot T v^+$ is twice the mean curvature of the surface v^+ , it is given by the well-known formula (see [Moon and Spencer 1970], for example)

$$\nabla \cdot T v^+ = -2H(v^+) = -\frac{EN + GL - 2FM}{EG - F^2},$$

where

$$E = (x_s)^2 + (y_s)^2 + (z_s)^2, \quad F = x_s x_t + y_s y_t + z_s z_t, \quad G = (x_t)^2 + (y_t)^2 + (z_t)^2,$$

and

$$L = \frac{\begin{vmatrix} x_{ss} & y_{ss} & z_{ss} \\ x_s & y_s & z_s \\ x_t & y_t & z_t \end{vmatrix}}{\sqrt{EG - F^2}}, \quad M = \frac{\begin{vmatrix} x_{st} & y_{st} & z_{st} \\ x_s & y_s & z_s \\ x_t & y_t & z_t \end{vmatrix}}{\sqrt{EG - F^2}}, \quad N = \frac{\begin{vmatrix} x_{tt} & y_{tt} & z_{tt} \\ x_s & y_s & z_s \\ x_t & y_t & z_t \end{vmatrix}}{\sqrt{EG - F^2}}.$$

After some calculation we obtain

$$\begin{aligned} \nabla \cdot T v^+ &= \frac{g_1''(t)}{(1 + (g'(t))^2)^{3/2}} \\ &+ \frac{f_1''(s)}{(1 + (f_1'(s))^2)^{3/2}} \left(1 + t \frac{f_1''(s)}{(1 + (f_1'(s))^2)^{3/2}} \right) \frac{g'(t)}{\sqrt{1 + (g'(t))^2}}. \end{aligned}$$

Recalling that we have chosen the domain Ω_0 so that (3-5) holds in Ω_0 and that $f_1''(s) \in C^2([0, \epsilon_0])$, in order to show $\nabla \cdot T v^+$ is bounded, all we need to show is that $g_1''(t)/(1 + (g'(t))^2)^{3/2}$ is bounded, that is to say, the curvature of the curve $g(t)$ is bounded. For the case of $\gamma \neq 0$, we have chosen $g(t)$ to be a linear function, so $g''(t)$ is zero. For the case of $\gamma = 0$, we have chosen $g(t)$ to be the part of a circle with radius t_0 , so $g_1''(t)/(1 + (g'(t))^2)^{3/2} = 1/t_0$. In either case, it follows that $\nabla \cdot T v^+$ is bounded. We now consider the quantity $\nabla \cdot T v^+ - v^+$, which can be written as

$$\nabla \cdot T v^+ - v^+ = \nabla \cdot T v^+ - (g(t) + K).$$

It follows immediately from the choice of $g(t)$ that it is bounded in the domain $\bar{\Omega}_0$ and also we have shown that twice the mean curvature $\nabla \cdot T v^+$ is bounded and does not depend on K . Hence there exists a constant K_0 such that

$$\nabla \cdot T v^+ - v^+ = \nabla \cdot T v^+ - (g(t) + K) \leq 0 \quad \text{for all } K \geq K_0.$$

Thus we have shown that the surface v^+ satisfies the (A-1) when $K > K_0$.

We now put the last piece of the puzzle in place by showing v^+ satisfies (A-4) for an appropriate choice of the constant K . Corollary A.1 implies that the capillary surface u is bounded away from the cusp, hence it is bounded on

$$\partial\Omega_0 \setminus (\partial\Omega_1 \cup \partial\Omega_2 \cup \{(0, 0)\}).$$

Since $g(t)$ is bounded in the domain $\bar{\Omega}_0$, there exists a constant $K_1 \geq K_0$ such that $g(t) + K_1 > u$ on $\partial\Omega_0 \setminus (\partial\Omega_1 \cup \partial\Omega_2 \cup \{(0, 0)\})$. Thus the surface v^+ satisfies (A-4) when $K = K_1$.

We have shown that the surface $v^+(s, t)$ defined in (3-6)–(3-7) satisfies inequalities (A-1)–(A-4), so by the Concus–Finn comparison principle we have

$$v^+(s, t) \geq u(x, y) \quad \text{in } \Omega_0.$$

Therefore the capillary surface at the cusp is bounded above when $\gamma_1 + \gamma_2 = \pi$ and each boundary $(\partial\Omega_1, \partial\Omega_2)$ has finite curvature near the cusp.

We can follow the similar steps for constructing the subsurface to show that this capillary surface is bounded below. We first construct a coordinate system such that one of the families of the coordinate curves is parallel curves of the boundary $\partial\Omega_2$ and another is perpendicular lines of the boundary $\partial\Omega_2$. Then choose a surface v^- so that the height only depends on the distance from $\partial\Omega_2$ which satisfies the contact angle condition exactly on $\partial\Omega_2$ and also it satisfies $\vec{v}_1 \cdot T v^- - \cos \gamma \leq 0$. By choosing v^- to have the bounded height and the finite mean curvature, we can shift this surface downwards enough to satisfy $\nabla \cdot T v^- - v^- \geq 0$ in Ω_0 and $v^- \leq u$ on $\partial\Omega_0 \setminus (\partial\Omega_1 \cup \partial\Omega_2 \cup \{(0, 0)\})$. Then using the Concus–Finn comparison principle, we can prove that $u(x, y)$ is bounded below.

Thus by showing that there exist bounded sub- and supersolutions of the Laplace–Young capillary surface equation, we have proven that the capillary surface is bounded if the contact angles of the boundaries are supplementary angles and boundaries have finite curvatures near the cusp. \square

3B. Proof of the continuity of the capillary surface when $\gamma_1 + \gamma_2 = \pi$.

Theorem 3.2. *If the capillary surface satisfies the conditions in Theorem 3.1, it is continuous at the cusp.*

Proof. Having established the boundedness of the solution, we can use the methods of [Lancaster and Siegel 1996] to establish a parametric description of the surface, with parameter domain at first the unit disk. The above comparison surface is needed in proving Case 5 (page 173) in that reference. Assuming the surface is discontinuous at the corner implies that an arc of the unit circle corresponds to the points on the surface above the corner point. A change of coordinates allows us to use the half-unit disk as the parameter domain, where the boundary line segment corresponds to the points on the surface above the corner point. Following the proof of Step 3 (page 175) of [Lancaster and Siegel 1996], for two different heights, there are level curves going through the corner point, and this leads to a contradiction (last paragraph of page 175 of the same reference). \square

4. Concluding remarks

We have shown that the validity of the statement “[the capillary surface] rises with the same order like the order of contact of the two arcs, which form the cusp”

[Scholz 2004] is not restricted to power-law cusps; it can be extended further. Our proof directly uses the functions $f_1(x)$ and $f_2(x)$ without approximating them by series. This idea has given us an advantage in the sense that our leading order term expression gives clearer intuitive understanding of the relationship between the shape of the domain and the shape of the singular capillary surface. Also as shown in an Example in Subsection 2.4.1, our leading order term gives first three terms of the power series asymptotic expansion, owing to the fact we have avoided approximating the boundary by the power series.

Even though we have extended the results beyond power-series cusps, our results still suffer from certain restrictions, including (2-8)–(2-12). Also a complete asymptotic series solution maybe desirable in order to claim a complete understanding of the asymptotic behavior; however, this will require further assumptions to the boundary functions f_1 and f_2 . The authors suspect that functions f_1 and f_2 of a form similar to the right-hand side of (2-27) can be potential candidates for a type of cusp for which a complete asymptotic series can be determined.

Also we have shown the previously unknown phenomenon of a bounded capillary surface in a cusp domain is possible when the contact angles of the two walls are supplementary (i.e., $\gamma_1 + \gamma_2 = \pi$). Although our proof covers most of the cases when the boundaries are smooth except at the cusp, the behavior of the capillary surface is unknown when the curvature of the boundary is not finite at the cusp. For example, it is unknown whether or not the capillary surface is bounded in a cusp domain bounded by $f_1 = x^{3/2}$ and $f_2 = -x^{3/2}$ when the contact angles of the two walls are supplementary.

The phenomenon that the capillary surface can be bounded or unbounded in a cusp domain depending on the contact angle can be interesting physically, as it indicates that a gradual change in the contact angle (e.g., by changing the temperature of the liquid) can cause a dramatic change in the liquid surface from unbounded to bounded. However, as the bounded capillary surface in a cusp domain only appears when the contact angles are exactly supplementary, it is not unknown to the authors how easily this phenomena can be observed through an experiment.

Thus we end this paper by remarking that the further exploration of singular capillary surfaces through theoretical, experimental and possibly numerical analyses is desired.

Appendix: The Concus–Finn comparison principle

In Sections 2C and 3A we have used the Concus–Finn comparison principle. We present it here for readers unfamiliar with it; see [Finn 1986, pages 110–113; 1989] for detailed discussions and proofs. We use the following formulation of the comparison principle:

Theorem A.1 (supersolution). *Let $u(x, y)$ be a solution of the boundary value problem (1-3)–(1-5) and let Ω_0 be a subdomain of Ω , with boundary $\partial\Omega_0$. Suppose a function $v^+(x, y)$ satisfies the inequalities*

$$(A-1) \quad \nabla \cdot T v^+ - v^+ \leq 0 \quad \text{in } \Omega_0,$$

$$(A-2) \quad \vec{v}_1 \cdot T v^+ - \cos \gamma_1 \geq 0 \quad \text{on } \partial\Omega_1 \cap \partial\Omega_0,$$

$$(A-3) \quad \vec{v}_2 \cdot T v^+ - \cos \gamma_2 \geq 0 \quad \text{on } \partial\Omega_2 \cap \partial\Omega_0,$$

$$(A-4) \quad v^+(x, y) \geq u(x, y) \quad \text{on } \partial\Omega_0 \setminus (\partial\Omega_1 \cup \partial\Omega_2 \cup \{(0, 0)\}).$$

Then $v^+(x, y)$ is a supersolution of the boundary value problem (1-3)–(1-5), i.e.,

$$v^+(x, y) \geq u(x, y) \quad \text{in } \Omega_0.$$

A similar statement holds for subsolutions.

Also we make use of one of the corollaries of the comparison principle to construct an upper bound for the solution; see [Concus and Finn 1974] or pages 113–114 of [Finn 1986].

Corollary A.1 (bound by hemispheres). *Let $u(x, y)$ be a solution of the boundary value problem (1-3)–(1-5) and $B_{r_0}(x_0, y_0)$ a disk of radius $r_0 > 0$ centered at (x_0, y_0) . If $B_{r_0}(x_0, y_0) \subseteq \Omega$, then*

$$(A-5) \quad -\left(\frac{1}{r_0} + r_0\right) \leq u(x, y) \leq \frac{1}{r_0} + r_0 \quad \text{in } B_{r_0}(x_0, y_0).$$

Recalling from (1-2) that the boundary is assumed to be of class C^3 away from the origin, it follows immediately from Corollary A.1 that $u(x, y)$ can only be unbounded at the origin (cusp).

Acknowledgements

The authors thank Professor Robert Finn of Stanford University for useful advice and interesting discussions during his visit to the University of Waterloo in May 2007.

Aoki thanks Professor John Wainwright of the University of Waterloo, whose comments during the author's master's degree thesis defence led to refinements in some of the results in Section 2. He also thanks the National Institute of Informatics, since part of this manuscript was prepared during his internship at the institute in 2010.

Both authors thank the University of Waterloo and the Natural Science and Engineering Council of Canada, which funded this research.

References

- [Aoki 2007] Y. Aoki, “Analysis of asymptotic solutions for cusp problems in capillarity”, Master’s thesis, University of Waterloo, 2007, available at <http://hdl.handle.net/10012/3352>.
- [Concus and Finn 1969] P. Concus and R. Finn, “On the behavior of a capillary surface in a wedge”, *Proc. Nat. Acad. Sci. USA* **63**:2 (1969), 292–299. Zbl 0219.76104
- [Concus and Finn 1974] P. Concus and R. Finn, “On capillary free surfaces in a gravitational field”, *Acta Math.* **132** (1974), 207–223. MR 58 #32327c Zbl 0382.76005
- [Finn 1986] R. Finn, *Equilibrium capillary surfaces*, Grundlehren der Math. Wiss. **284**, Springer, New York, 1986. MR 88f:49001 Zbl 0583.35002
- [Finn and Hwang 1989] R. Finn and J.-F. Hwang, “On the comparison principle for capillary surfaces”, *J. Fac. Sci. Univ. Tokyo Sect. IA Math.* **36**:1 (1989), 131–134. MR 90h:35099 Zbl 0684.35007
- [King et al. 1999] J. R. King, J. R. Ockendon, and H. Ockendon, “The Laplace–Young equation near a corner”, *Quart. J. Mech. Appl. Math.* **52**:1 (1999), 73–97. MR 2000f:76017 Zbl 0932.76012
- [Lancaster and Siegel 1996] K. E. Lancaster and D. Siegel, “Existence and behavior of the radial limits of a bounded capillary surface at a corner”, *Pacific J. Math.* **176**:1 (1996), 165–194. MR 98g:58030a Zbl 0866.76018
- [Miersemann 1993] E. Miersemann, “Asymptotic expansion at a corner for the capillary problem: the singular case”, *Pacific J. Math.* **157**:1 (1993), 95–107. MR 93m:35039 Zbl 0796.76020
- [Moon and Spencer 1970] P. Moon and D. E. Spencer, *Field theory handbook: including coordinate systems, differential equations and their solutions*, 2nd ed., Springer, Berlin, 1970. MR 89i:00026 Zbl 0097.39403
- [Norbury et al. 2005] J. Norbury, G. C. Sander, and C. F. Scott, “Corner solutions of the Laplace–Young equation”, *Quart. J. Mech. Appl. Math.* **58**:1 (2005), 55–71. MR 2006b:76017 Zbl 1064.76020
- [Scholz 2001] M. Scholz, *Über das Verhalten von Kapillarflächen in Spitzen*, Ph.D. thesis, Universität Leipzig, 2001, available at <http://www.people.imise.uni-leipzig.de/markus.scholz/pdf/p1.pdf>.
- [Scholz 2004] M. Scholz, “On the asymptotic behaviour of capillary surfaces in cusps”, *Z. Angew. Math. Phys.* **55**:2 (2004), 216–234. MR 2005e:76023 Zbl 1058.35059

Received June 17, 2011. Revised December 16, 2011.

YASUNORI AOKI
DEPARTMENT OF APPLIED MATHEMATICS
UNIVERSITY OF WATERLOO
200 UNIVERSITY AVENUE WEST
WATERLOO, ON N2L 3G1
CANADA
yaoki@uwaterloo.ca

DAVID SIEGEL
DEPARTMENT OF APPLIED MATHEMATICS
UNIVERSITY OF WATERLOO
200 UNIVERSITY AVENUE WEST
WATERLOO, ON N2L 3G1
CANADA
dsiegel@uwaterloo.ca

ON ORTHOGONAL POLYNOMIALS WITH RESPECT TO CERTAIN DISCRETE SOBOLEV INNER PRODUCT

FRANCISCO MARCELLÁN, RAMADAN ZEJNULLAHU,
BUJAR FEJZULLAHU AND EDMUNDO HUERTAS

In this paper we deal with sequences of polynomials orthogonal with respect to the discrete Sobolev inner product

$$\langle f, g \rangle_S = \int_0^\infty \omega(x) f(x) g(x) dx + Mf(\xi)g(\xi) + Nf'(\xi)g'(\xi),$$

where ω is a weight function, $\xi \leq 0$, and $M, N \geq 0$. The location of the zeros of discrete Sobolev orthogonal polynomials is given in terms of the zeros of standard polynomials orthogonal with respect to the weight function ω . In particular, for $\omega(x) = x^\alpha e^{-x}$ we obtain the asymptotics for discrete Laguerre–Sobolev orthogonal polynomials.

1. Introduction

Polynomials orthogonal with respect to an inner product

$$(1) \quad \langle f, g \rangle = \int_E \omega(x) f(x) g(x) dx + Mf(\xi)g(\xi) + Nf'(\xi)g'(\xi),$$

where ξ is a real number and $d\mu$ is a positive Borel measure supported on an infinite subset E of the real line have been considered by several authors (see, for instance, [Alfaro et al. 1992; López et al. 1995; Marcellán and Ronveaux 1990; Marcellán and Van Assche 1993] and the references therein). They are known in the literature as Sobolev-type or discrete Sobolev orthogonal polynomials. Special attention has been paid to their algebraic and analytic properties of these polynomials, in particular, the distribution of their zeros taking into account the location of the point ξ with respect to the set E .

When E is the interval $[0, +\infty)$ and $\xi = 0$, Meijer [1993a] analyzed some analytic properties of the zeros of the so called discrete Sobolev orthogonal polynomials (1). Some results of [Meijer 1993a] are direct generalizations of the results of [Koekoek and Meijer 1993], where the weight function is the Laguerre

MSC2010: primary 33C47; secondary 42C05.

Keywords: orthogonal polynomials, discrete Sobolev polynomials, Laguerre polynomials, asymptotics.

weight $\omega(x) = x^\alpha e^{-x}$. Koekoek and Meijer established properties of the discrete Laguerre–Sobolev polynomials such as their representation as a hypergeometric series, an holonomic second order linear differential equation associated with them, properties of the zeros, and a higher-order recurrence relation that such polynomials satisfy. The asymptotic properties of these discrete Laguerre–Sobolev polynomials have been studied in [Álvarez-Nodarse and Moreno-Balcázar 2004; Marcellán and Moreno-Balcázar 2006], while the analysis of convergence of the Fourier expansions in terms of such polynomials was done in [Fejzullahu and Marcellán 2009].

In this paper we consider the discrete Sobolev polynomials $\{\hat{S}_n\}_{n \geq 0}$ orthogonal with respect to (1) where $E = [0, +\infty)$ and $\xi < 0$. We show that these polynomials can be expressed as

$$\hat{S}_n(x) = \hat{P}_n(x) + A_{n,1}(x - \xi) \hat{P}_{n-1}^{[2]}(x) + A_{n,2}(x - \xi)^2 \hat{P}_{n-2}^{[4]}(x),$$

where $\{\hat{P}_n\}_{n \geq 0}$ and $\{\hat{P}_n^{[k]}\}_{n \geq 0}$, $k \in \mathbb{N}$, are the sequences of monic polynomials orthogonal with respect to the weight functions $\omega(\cdot)$ and $(\cdot - \xi)^k \omega(\cdot)$, respectively. Moreover, the behavior of the coefficients $A_{n,1}$ and $A_{n,2}$ is studied in more detail. In particular, when ω is the Laguerre weight, we obtain some asymptotic properties for the sequence of discrete Laguerre–Sobolev orthogonal polynomials.

The structure of the manuscript is as follows. In Section 2 we give some basic background concerning polynomial perturbations of a measure as well as interlacing properties for the zeros of the corresponding orthogonal polynomials. We point out that the results presented therein are of independent interest in terms of the core of our contribution. Indeed, they constitute an alternative approach in the subject. In Section 3, a representation of monic polynomials orthogonal with respect to the inner product (1) is given in terms of polynomial orthogonal with respect to polynomial perturbations of the weight function. Some results about their zeros are deduced. In Section 4 we focus our attention on the asymptotics of discrete Laguerre–Sobolev orthogonal polynomials. More precisely, we obtain outer relative asymptotics, a Mehler–Heine formula and the Plancherel–Rotach outer asymptotics for such orthogonal polynomials.

Throughout this paper positive constants are denoted by c, c_1, \dots , and they may vary at every occurrence. The notation $u_n \cong v_n$ means that the sequence $\{u_n/v_n\}_n$ converges to 1. We will denote by $k(\pi_n)$ the leading coefficient of any polynomial π_n and $\hat{\pi}_n(x) = (k(\pi_n))^{-1} \pi_n(x)$.

2. Auxiliary results

Let ω denote a weight function on $(0, \infty)$, i.e., $\omega(x) \geq 0$ and all moments

$$c_n = \int_0^\infty \omega(x) x^n dx, \quad n = 0, 1, \dots$$

exist. Let $\{\hat{P}_n\}_{n \geq 0}$ denote the sequence of monic polynomials orthogonal (SMOP, in short) with respect to the standard inner product

$$\langle f, g \rangle = \int_0^\infty \omega(x) f(x) g(x) dx.$$

In particular, from the moments we get an explicit expression of the SMOP. Indeed, we get

$$\hat{P}_0(x) = 1$$

and

$$(2) \quad \hat{P}_n(x) = \frac{1}{\Delta_{n-1}} \begin{vmatrix} c_0 & c_1 & c_2 & \dots & c_n \\ c_1 & c_2 & c_3 & \dots & c_{n+1} \\ c_2 & c_3 & c_4 & \dots & c_{n+2} \\ \cdot & \cdot & \cdot & \dots & \cdot \\ c_{n-1} & c_n & c_{n+1} & \dots & c_{2n-1} \\ 1 & x & x^2 & \dots & x^n \end{vmatrix}, \quad n \geq 1,$$

where

$$\Delta_{n-1} = \begin{vmatrix} c_0 & c_1 & c_2 & \dots & c_{n-1} \\ c_1 & c_2 & c_3 & \dots & c_n \\ c_2 & c_3 & c_4 & \dots & c_{n+1} \\ \cdot & \cdot & \cdot & \dots & \cdot \\ c_{n-1} & c_n & c_{n+1} & \dots & c_{2n-2} \end{vmatrix}, \quad n \geq 1,$$

are the Gram determinants.

The n -th reproducing kernel for ω is

$$K_n(x, y) = \sum_{k=0}^n \frac{\hat{P}_k(x) \hat{P}_k(y)}{\|\hat{P}_k\|_\omega^2}.$$

Here, $\|\hat{P}_n\|_\omega^2 = \langle \hat{P}_n, \hat{P}_n \rangle$. Because of the Christoffel–Darboux formula, it may also be expressed as

$$K_n(x, y) = \frac{1}{\|\hat{P}_n\|_\omega^2} \frac{\hat{P}_{n+1}(x) \hat{P}_n(y) - \hat{P}_n(x) \hat{P}_{n+1}(y)}{x - y}.$$

The confluent formula reads as

$$(3) \quad K_n(x, x) = \sum_{k=0}^n \frac{(\hat{P}_k(x))^2}{\|\hat{P}_k\|_\omega^2} = \frac{1}{\|\hat{P}_n\|_\omega^2} (\hat{P}'_{n+1}(x) \hat{P}_n(x) - \hat{P}'_n(x) \hat{P}_{n+1}(x)).$$

In the same way we can describe the SMOP $\{\hat{P}_n^{[k]}\}_{n \geq 0}$, orthogonal with respect to the inner product

$$\langle f, g \rangle_k = \int_0^\infty (x - \xi)^k \omega(x) f(x) g(x) dx,$$

where $\xi \leq 0$. For $n \geq 1$, they are given by the determinant (2) where c_i is replaced by d_i^k , $k \in \mathbb{N}$, where

$$(4) \quad d_n^k = \int_0^\infty (x - \xi)^k \omega(x) x^n dx = d_{n+1}^{k-1} - \xi d_n^{k-1}, \quad n = 0, 1, \dots,$$

and $c_n = d_n^0$. In the sequel, we will set

$$\|\hat{P}_n^{[k]}\|_{\omega,k}^2 = \int_0^\infty (x - \xi)^k \omega(x) (\hat{P}_n^{[k]}(x))^2 dx.$$

Proposition 1. Let $D_{n-1}^k = \det[a_{ij}^k]_{0 \leq i,j \leq n-1}$, where $a_{ij}^k = d_{i+j}^k$, $k \in \mathbb{N}$. Then

$$(5) \quad D_{n-1}^k = (-1)^n D_{n-1}^{k-1} \hat{P}_n^{[k-1]}(\xi),$$

with $D_{n-1}^0 = \Delta_{n-1}$.

Proof. For $n \geq 1$ and $k \in \mathbb{N}$,

$$(6) \quad \hat{P}_n^{[k-1]}(x) = \frac{1}{D_{n-1}^{[k-1]}} \begin{vmatrix} d_0^{k-1} & d_1^{k-1} & \dots & d_n^{k-1} \\ d_1^{k-1} & d_2^{k-1} & \dots & d_{n+1}^{k-1} \\ \cdot & \cdot & \dots & \cdot \\ d_{n-1}^{k-1} & d_n^{k-1} & \dots & d_{2n-1}^{k-1} \\ 1 & x & \dots & x^n \end{vmatrix},$$

with $\hat{P}_n = \hat{P}_n^{[0]}$. The determinant in (6) becomes [Szegő 1975, Formula (2.2.9)]

$$\hat{P}_n^{[k-1]}(x) = \frac{(-1)^n}{D_{n-1}^{k-1}} \begin{vmatrix} d_1^{k-1} - d_0^{k-1}x & d_2^{k-1} - d_1^{k-1}x & \dots & d_n^{k-1} - d_{n-1}^{k-1}x \\ d_2^{k-1} - d_1^{k-1}x & d_3^{k-1} - d_2^{k-1}x & \dots & d_{n+1}^{k-1} - d_n^{k-1}x \\ \cdot & \cdot & \dots & \cdot \\ d_n^{k-1} - d_{n-1}^{k-1}x & d_{n+1}^{k-1} - d_n^{k-1}x & \dots & d_{2n-1}^{k-1} - d_{2n-2}^{k-1}x \end{vmatrix}.$$

Now, by using (4), (5) follows. □

Next we will compute some integrals involving the polynomials $\hat{P}_n^{[k]}$.

Proposition 2. (i) The integral $\int_0^\infty (x - \xi)^{k-1} \omega(x) \hat{P}_n^{[k]}(x) dx$ is given by

$$\frac{\|\hat{P}_n^{[k-1]}\|_{\omega,k-1}^2}{\hat{P}_n^{[k-1]}(\xi)} = \begin{cases} \frac{\|\hat{P}_n\|_{\omega}^2}{\hat{P}_n(\xi)} & \text{if } k = 1, \\ \frac{(-1)^{k-1}}{\hat{P}_n^{[k-1]}(\xi)} \prod_{i=1}^{k-1} \frac{\hat{P}_{n+1}^{[i-1]}(\xi)}{\hat{P}_n^{[i-1]}(\xi)} \|\hat{P}_n\|_{\omega}^2 & \text{if } k \geq 2. \end{cases}$$

(ii) The integral $\int_0^\infty (x - \xi)^{k-2} \omega(x) \hat{P}_n^{[k]}(x) dx$ is given by

$$\frac{(\hat{P}_{n+1}^{[k-2]}(x))'_{x=\xi} \|\hat{P}_n^{[k-2]}\|_{\omega, k-2}^2}{\hat{P}_n^{[k-1]}(\xi) \hat{P}_n^{[k-2]}(\xi)} = \begin{cases} \frac{(\hat{P}_{n+1}(x))'_{x=\xi} \|\hat{P}_n\|_{\omega}^2}{\hat{P}_n(\xi) \hat{P}_n^{[1]}(\xi)} & \text{if } k = 2, \\ \frac{(-1)^k (\hat{P}_{n+1}^{[k-2]}(x))'_{x=\xi}}{\hat{P}_n^{[k-1]}(\xi) \hat{P}_n^{[k-2]}(\xi)} \prod_{i=1}^{k-2} \frac{\hat{P}_{n+1}^{[i-1]}(\xi)}{\hat{P}_n^{[i-1]}(\xi)} \|\hat{P}_n\|_{\omega}^2 & \text{if } k \geq 3. \end{cases}$$

Proof. (i) Using (4) recursively as well as properties of determinants, we have

$$\begin{aligned} D_{n-1}^k \int_0^\infty (x - \xi)^{k-1} \omega(x) \hat{P}_n^{[k]}(x) dx &= \begin{vmatrix} d_0^k & d_1^k & d_2^k & \dots & d_n^k \\ d_1^k & d_2^k & d_3^k & \dots & d_{n+1}^k \\ \cdot & \cdot & \cdot & \dots & \cdot \\ d_{n-1}^k & d_n^k & d_{n+1}^k & \dots & d_{2n-1}^k \\ d_0^{k-1} & d_1^{k-1} & d_2^{k-1} & \dots & d_n^{k-1} \end{vmatrix} \\ &= \begin{vmatrix} d_1^{k-1} & d_2^{k-1} & d_3^{k-1} & \dots & d_{n+1}^{k-1} \\ d_2^{k-1} & d_3^{k-1} & d_4^{k-1} & \dots & d_{n+2}^{k-1} \\ \cdot & \cdot & \cdot & \dots & \cdot \\ d_n^{k-1} & d_{n+1}^{k-1} & d_{n+2}^{k-1} & \dots & d_{2n}^{k-1} \\ d_0^{k-1} & d_1^{k-1} & d_2^{k-1} & \dots & d_n^{k-1} \end{vmatrix} \\ &= (-1)^n D_n^{k-1}. \end{aligned}$$

On the other hand,

$$\|\hat{P}_n^{[k-1]}\|_{\omega, k-1}^2 = \int_0^\infty (x - \xi)^{k-1} \omega(x) x^n \hat{P}_n^{[k-1]}(x) dx = \frac{D_n^{k-1}}{D_{n-1}^{k-1}},$$

and by using (5) we get

$$(7) \quad \int_0^\infty (x - \xi)^{k-1} \omega(x) \hat{P}_n^{[k]}(x) dx = \frac{(-1)^n D_{n-1}^{k-1} \|\hat{P}_n^{[k-1]}\|_{\omega, k-1}^2}{D_{n-1}^k} = \frac{\|\hat{P}_n^{[k-1]}\|_{\omega, k-1}^2}{\hat{P}_n^{[k-1]}(\xi)}.$$

On the other hand, we have from [Szegő 1975, Theorem 2.5]

$$(8) \quad (x - \xi) \hat{P}_n^{[k]}(x) = \hat{P}_{n+1}^{[k-1]}(x) - \frac{\hat{P}_{n+1}^{[k-1]}(\xi)}{\hat{P}_n^{[k-1]}(\xi)} \hat{P}_n^{[k-1]}(x).$$

Therefore,

$$\|\hat{P}_n^{[k]}\|_{\omega,k}^2 = -\frac{\hat{P}_{n+1}^{[k-1]}(\xi)}{\hat{P}_n^{[k-1]}(\xi)} \|\hat{P}_n^{[k-1]}\|_{\omega,k-1}^2.$$

Using this relation recursively we obtain

$$(9) \quad \|\hat{P}_n^{[k]}\|_{\omega,k}^2 = (-1)^k \prod_{i=1}^k \frac{\hat{P}_{n+1}^{[i-1]}(\xi)}{\hat{P}_n^{[i-1]}(\xi)} \|\hat{P}_n\|_{\omega}^2, \quad k \geq 2.$$

Combining (7) and (9), our statement follows.

(ii) We have

$$(10) \quad (\hat{P}_{n+1}^{[k-2]}(x))' = \frac{1}{D_n^{k-2}} \begin{vmatrix} d_0^{k-2} & d_1^{k-2} & d_2^{k-2} & \dots & d_{n+1}^{k-2} \\ d_1^{k-2} & d_2^{k-2} & d_3^{k-2} & \dots & d_{n+2}^{k-2} \\ \cdot & \cdot & \cdot & \dots & \cdot \\ d_n^{k-2} & d_{n+1}^{k-2} & d_{n+2}^{k-2} & \dots & d_{2n+1}^{k-2} \\ 0 & 1 & 2x & \dots & nx^{n-1} \end{vmatrix}, \quad n \geq 0.$$

Now, adding to the last column the n -th and $(n-1)$ -th columns multiplied by $-2x$ and x^2 , respectively, and repeating this operation for each of the preceding columns, we obtain

$$(11) \quad (\hat{P}_{n+1}^{[k-2]}(x))' = \frac{1}{D_n^{k-2}} \begin{vmatrix} d_0^{k-2} & d_1^{k-2} & d_2^{k-2} - 2xd_1^{k-2} + x^2d_0^{k-2} & \dots & d_{n+1}^{k-2} - 2xd_n^{k-2} + x^2d_{n-1}^{k-2} \\ d_1^{k-2} & d_2^{k-2} & d_3^{k-2} - 2xd_2^{k-2} + x^2d_1^{k-2} & \dots & d_{n+2}^{k-2} - 2xd_{n+1}^{k-2} + x^2d_n^{k-2} \\ \cdot & \cdot & \cdot & \dots & \cdot \\ d_n^{k-2} & d_{n+1}^{k-2} & d_{n+2}^{k-2} - 2xd_{n+1}^{k-2} + x^2d_n^{k-2} & \dots & d_{2n+1}^{k-2} - 2xd_{2n}^{k-2} + x^2d_{2n-1}^{k-2} \\ 0 & 1 & 0 & \dots & 0 \end{vmatrix} \\ = \frac{1}{D_n^{k-2}} \begin{vmatrix} d_2^{k-2} - 2xd_1^{k-2} + x^2d_0^{k-2} & d_3^{k-2} - 2xd_2^{k-2} + x^2d_1^{k-2} & \dots & d_{n+2}^{k-2} - 2xd_{n+1}^{k-2} + x^2d_n^{k-2} \\ d_3^{k-2} - 2xd_2^{k-2} + x^2d_1^{k-2} & d_4^{k-2} - 2xd_3^{k-2} + x^2d_2^{k-2} & \dots & d_{n+3}^{k-2} - 2xd_{n+2}^{k-2} + x^2d_{n+1}^{k-2} \\ \cdot & \cdot & \dots & \cdot \\ d_{n+1}^{k-2} - 2xd_n^{k-2} + x^2d_{n-1}^{k-2} & d_{n+2}^{k-2} - 2xd_{n+1}^{k-2} + x^2d_n^{k-2} & \dots & d_{2n+1}^{k-2} - 2xd_{2n}^{k-2} + x^2d_{2n-1}^{k-2} \\ \cdot & d_0^{k-2} & \dots & d_1^{k-2} \\ \cdot & \cdot & \dots & \cdot \end{vmatrix}.$$

On the other hand,

$$D_{n-1}^k \int_0^\infty (x - \xi)^{k-2} \omega(x) \hat{P}_n^{[k]}(x) dx = \begin{vmatrix} d_0^k & d_1^k & \dots & d_n^k \\ d_1^k & d_2^k & \dots & d_{n+1}^k \\ \cdot & \cdot & \dots & \cdot \\ d_{n-1}^k & d_n^k & \dots & d_{2n-1}^k \\ d_0^{k-2} & d_1^{k-2} & \dots & d_n^{k-2} \end{vmatrix},$$

and by using (5), (9), and (11) we get

$$\begin{aligned} \int_0^\infty (x - \xi)^{k-2} \omega(x) \hat{P}_n^{[k]}(x) dx &= \frac{D_n^{k-2} (\hat{P}_{n+1}^{[k-2]}(x))'_{x=\xi}}{D_{n-1}^{k-2} \hat{P}_n^{[k-1]}(\xi) \hat{P}_n^{[k-2]}(\xi)} = \frac{(\hat{P}_{n+1}^{[k-2]}(x))'_{x=\xi} \|\hat{P}_n^{[k-2]}\|_{\omega, k-2}^2}{\hat{P}_n^{[k-1]}(\xi) \hat{P}_n^{[k-2]}(\xi)} \\ &= \frac{(-1)^k (\hat{P}_{n+1}^{[k-2]}(x))'_{x=\xi}}{\hat{P}_n^{[k-1]}(\xi) \hat{P}_n^{[k-2]}(\xi)} \prod_{i=1}^{k-2} \frac{\hat{P}_{n+1}^{[i-1]}(\xi)}{\hat{P}_n^{[i-1]}(\xi)} \|\hat{P}_n\|_{\omega}^2. \quad \square \end{aligned}$$

Denote by $x_{r,n}^{[k]}$, $r = 1, 2, \dots, n$, the zeros of $\hat{P}_n^{[k]}(x)$ in increasing order.

Proposition 3. (i) *The zeros of $\hat{P}_n^{[k]}(x)$ interlace with both the zeros of $\hat{P}_{n+1}^{[k-1]}(x)$ and $\hat{P}_n^{[k-1]}(x)$, i.e.,*

$$x_{r,n}^{[k-1]} < x_{r,n}^{[k]} < x_{r+1,n+1}^{[k-1]}, \quad r = 1, 2, \dots, n.$$

(ii) *Between two consecutive zeros of $\hat{P}_{n+1}^{[k-2]}$, $k \geq 2$, there is exactly one zero of $\hat{P}_n^{[k]}$.*

(iii) $\text{sgn } \hat{P}_n^{[k-2]}(x_{r,n-1}^{[k]}) = (-1)^{n-r} = -\text{sgn } \hat{P}_{n-2}^{[k+2]}(x_{r,n-1}^{[k]})$ for $r = 1, 2, \dots, n-1$.

Proof. (i) Here we will use the same argument as in [Chihara 1978, page 65] (see also [Bracciali et al. 2002, Lemma 1]). It is well known that the zeros of $\hat{P}_{n+1}^{[k-1]}$ interlace with the zeros of $\hat{P}_n^{[k-1]}$, i.e.,

$$0 < x_{1,n+1}^{[k-1]} < x_{1,n}^{[k-1]} < x_{2,n+1}^{[k-1]} < \dots < x_{n,n}^{[k-1]} < x_{n+1,n+1}^{[k-1]} < \infty.$$

From (5) $\hat{P}_{n+1}^{[k-1]}(\xi)/\hat{P}_n^{[k-1]}(\xi) < 0$ and taking (8) into account we have

$$\begin{aligned} \text{sgn } \hat{P}_n^{[k]}(x_{r,n+1}^{[k-1]}) &= \text{sgn } \hat{P}_n^{[k-1]}(x_{r,n+1}^{[k-1]}) = (-1)^{n-r+1} \quad \text{for } r = 1, 2, \dots, n+1, \\ \text{sgn } \hat{P}_n^{[k]}(x_{r,n}^{[k-1]}) &= \text{sgn } \hat{P}_{n+1}^{[k-1]}(x_{r,n}^{[k-1]}) = (-1)^{n-r+1} \quad \text{for } r = 1, 2, \dots, n. \end{aligned}$$

Thus, there exist zeros $x_{r,n}^{[k]}$, $r = 2, 3, \dots, n$, of $\hat{P}_n^{[k]}(x)$ satisfying

$$x_{r,n}^{[k-1]} < x_{r,n}^{[k]} < x_{r+1,n+1}^{[k-1]}, \quad r = 1, 2, \dots, n.$$

(ii) By using (8) and the recurrence relation we obtain

$$(x - \xi)^2 \hat{P}_n^{[k]}(x) = (d_{1,n}x + d_{2,n}) \hat{P}_{n+1}^{[k-2]}(x) + d_{3,n} \hat{P}_n^{[k-2]}(x).$$

Since $\hat{P}_{n+1}^{[k-2]}(\xi) \neq 0$ we have $d_{3,n} \neq 0$. Now, the rest of the proof can be done in a similar way as in [Meijer 1993a, Lemma 6.1]; see also [Meijer 1993b, Lemma 4.1].

(iii) From (ii) we have $x_{r,n}^{[k-2]} < x_{r,n-1}^{[k]} < x_{r+1,n}^{[k-2]}$ for $r = 1, 2, \dots, n - 1$. Therefore,

$$\operatorname{sgn} \hat{P}_n^{[k-2]}(x_{r,n-1}^{[k]}) = (-1)^{n-r}.$$

Again, according to (ii), $x_{r-1,n-2}^{[k+2]} < x_{r,n-1}^{[k]} < x_{r,n-2}^{[k+2]}$ for $r = 1, 2, \dots, n - 2$, and $x_{n-2,n-2}^{[k+2]} < x_{n-1,n-1}^{[k]}$. Therefore,

$$\operatorname{sgn} \hat{P}_{n-2}^{[k+2]}(x_{r,n-1}^{[k]}) = (-1)^{n-r-1} \quad \text{and} \quad \operatorname{sgn} \hat{P}_{n-2}^{[k+2]}(x_{n-1,n-1}^{[k]}) = 1.$$

As a conclusion,

$$\operatorname{sgn} \hat{P}_n^{[k-2]}(x_{r,n-1}^{[k]}) = -\operatorname{sgn} \hat{P}_{n-2}^{[k+2]}(x_{r,n-1}^{[k]}), \quad r = 1, 2, \dots, n - 1. \quad \square$$

3. Discrete Sobolev orthogonal polynomials

Connection formula. We consider the inner product

$$(12) \quad \langle f, g \rangle_S = \int_0^\infty \omega(x) f(x) g(x) dx + Mf(\xi)g(\xi) + Nf'(\xi)g'(\xi),$$

where $\xi \leq 0$, and $M, N \geq 0$. Let $\{\hat{S}_n\}_{n \geq 0}$ denote the SMOP with respect to the discrete Sobolev inner product (12)).

Theorem 1. *Let $M \geq 0$ and $N \geq 0$. There are real constants $A_{n,1}$ and $A_{n,2}$ such that*

$$\hat{S}_n(x) = \hat{P}_n(x) + A_{n,1}(x - \xi) \hat{P}_{n-1}^{[2]}(x) + A_{n,2}(x - \xi)^2 \hat{P}_{n-2}^{[4]}(x),$$

where

$$\begin{aligned} A_{n,1} &= \frac{NI_{2,n}(\xi) \hat{P}'_n(\xi) - MI_{3,n}(\xi) \hat{P}_n(\xi)}{I_{1,n}(\xi) I_{3,n}(\xi) - NI_{2,n}(\xi) \hat{P}_{n-1}^{[2]}(\xi)}, \\ A_{n,2} &= \frac{MN \hat{P}_n(\xi) \hat{P}_{n-1}^{[2]}(\xi) - NI_{1,n}(\xi) \hat{P}'_n(\xi)}{I_{1,n}(\xi) I_{3,n}(\xi) - NI_{2,n}(\xi) \hat{P}_{n-1}^{[2]}(\xi)}, \\ I_{1,n}(\xi) &= -\frac{\hat{P}_n(\xi)}{K_{n-1}(\xi, \xi)}, \\ I_{2,n}(\xi) &= \frac{\hat{P}_{n-1}(\xi) \hat{P}_{n-1}^{[1]}(\xi) \hat{P}_{n-1}^{[2]'}(\xi)}{\hat{P}_{n-2}(\xi) \hat{P}_{n-2}^{[1]}(\xi) \hat{P}_{n-2}^{[2]}(\xi) \hat{P}_{n-2}^{[3]}(\xi)} \|\hat{P}_{n-2}\|_\omega^2, \\ I_{3,n}(\xi) &= -\frac{\hat{P}_{n-1}(\xi) \hat{P}_{n-1}^{[1]}(\xi) \hat{P}_{n-1}^{[2]}(\xi)}{\hat{P}_{n-2}(\xi) \hat{P}_{n-2}^{[1]}(\xi) \hat{P}_{n-2}^{[2]}(\xi) \hat{P}_{n-2}^{[3]}(\xi)} \|\hat{P}_{n-2}\|_\omega^2. \end{aligned}$$

Proof. We will prove that

$$\langle \hat{S}_n, (\cdot - \xi)^k \rangle_S = 0 \quad \text{for } k = 0, 1, \dots, n - 1.$$

For $k \geq 2$ and $n > k$,

$$\begin{aligned} &\langle \hat{S}_n, (\cdot - \xi)^k \rangle_S \\ &= \int_0^\infty \omega(x) \hat{S}_n(x) (x - \xi)^k dx \\ &= \int_0^\infty \omega(x) \hat{P}_n(x) (x - \xi)^k dx + A_{n,1} \int_0^\infty (x - \xi)^2 \omega(x) \hat{P}_{n-1}^{[2]}(x) (x - \xi)^{k-1} dx \\ &\quad + A_{n,2} \int_0^\infty (x - \xi)^4 \omega(x) \hat{P}_{n-2}^{[4]}(x) (x - \xi)^{k-2} dx \\ &= 0, \end{aligned}$$

Now consider $k = 0$ and $n \geq 1$. We have

$$\begin{aligned} \langle \hat{S}_n, 1 \rangle_S &= \int_0^\infty \omega(x) \hat{S}_n(x) dx + M \hat{S}_n(\xi) \\ &= A_{n,1} \int_0^\infty (x - \xi) \omega(x) \hat{P}_{n-1}^{[2]}(x) dx + A_{n,2} \int_0^\infty (x - \xi)^2 \omega(x) \hat{P}_{n-2}^{[4]}(x) dx \\ &\quad + M \hat{P}_n(\xi). \end{aligned}$$

On the other hand, by using Proposition 2(i),

$$(13) \quad I_{1,n}(\xi) = \int_0^\infty (x - \xi) \omega(x) \hat{P}_{n-1}^{[2]}(x) dx = -\frac{\hat{P}_n(\xi)}{\hat{P}_{n-1}(\xi) \hat{P}_{n-1}^{[1]}(\xi)} \|\hat{P}_{n-1}\|_\omega^2,$$

and taking derivatives in (8) and then substituting $x = \xi$ we get

$$(14) \quad \hat{P}_{n-1}^{[k]}(\xi) = \left(\hat{P}_n^{[k-1]}(x) \right)'_{x=\xi} - \frac{\hat{P}_n^{[k-1]}(\xi)}{\hat{P}_{n-1}^{[k-1]}(\xi)} \left(\hat{P}_{n-1}^{[k-1]}(x) \right)'_{x=\xi}.$$

Combining (3), (13), and (14), we get

$$I_{1,n}(\xi) = -\frac{\hat{P}_n(\xi)}{K_{n-1}(\xi, \xi)}.$$

Using Proposition 2(ii),

$$\begin{aligned} (15) \quad I_{2,n}(\xi) &= \int_0^\infty (x - \xi)^2 \omega(x) \hat{P}_{n-2}^{[4]}(x) dx = \frac{\left(\hat{P}_{n-1}^{[2]}(x) \right)'_{x=\xi} \|\hat{P}_{n-2}\|_{\omega,2}^2}{\hat{P}_{n-2}^{[2]}(\xi) \hat{P}_{n-2}^{[3]}(\xi)} \\ &= \frac{\hat{P}_{n-1}(\xi) \hat{P}_{n-1}^{[1]}(\xi) \left(\hat{P}_{n-1}^{[2]}(x) \right)'_{x=\xi}}{\hat{P}_{n-2}(\xi) \hat{P}_{n-2}^{[1]}(\xi) \hat{P}_{n-2}^{[2]}(\xi) \hat{P}_{n-2}^{[3]}(\xi)} \|\hat{P}_{n-2}\|_\omega^2. \end{aligned}$$

Therefore,

$$\langle \hat{S}_n, 1 \rangle_S = A_{n,1} I_{1,n}(\xi) + A_{n,2} I_{2,n}(\xi) + M \hat{P}_n(\xi).$$

In the same way, for $k = 1$ and $n \geq 2$, we have

$$\begin{aligned} \langle \hat{S}_n, (\cdot - \xi) \rangle_S &= \int_0^\infty \omega(x) \hat{S}_n(x) (x - \xi) dx + N \hat{S}'_n(\xi) \\ &= A_{n,2} I_{3,n}(\xi) + N A_{n,1} \hat{P}_{n-1}^{[2]}(\xi) + N \hat{P}'_n(\xi), \end{aligned}$$

where

$$\begin{aligned} I_{3,n}(\xi) &= \int_0^\infty (x - \xi)^3 \omega(x) \hat{P}_{n-2}^{[4]}(x) dx = \frac{\|\hat{P}_{n-2}^{[3]}\|_{\omega,3}^2}{\hat{P}_n^{[3]}(\xi)} \\ &= -\frac{\hat{P}_{n-1}(\xi) \hat{P}_{n-1}^{[1]}(\xi) \hat{P}_{n-1}^{[2]}(\xi)}{\hat{P}_{n-2}(\xi) \hat{P}_{n-2}^{[1]}(\xi) \hat{P}_{n-2}^{[2]}(\xi) \hat{P}_{n-2}^{[3]}(\xi)} \|\hat{P}_{n-2}\|_{\omega}^2. \end{aligned}$$

Finally, using the expressions of $A_{n,1}$ and $A_{n,2}$, our statement follows. □

Next, we will study the behavior of the coefficients $A_{n,1}$ and $A_{n,2}$.

Proposition 4.

(i) $I_{1,n}(\xi) I_{3,n}(\xi) - N I_{2,n}(\xi) \hat{P}_{n-1}^{[2]}(\xi) = -I_{2,n}(\xi) \hat{P}_{n-1}^{[2]}(\xi) (N + \alpha_n \beta_n)$, where

$$0 < \alpha_n = \frac{I_{1,n}(\xi)}{\hat{P}_{n-1}^{[2]}(\xi)} < d_0^1 \quad \text{and} \quad \frac{d_0^3}{d_0^2} < -\frac{\hat{P}_{n-1}^{[2]'}(\xi)}{\hat{P}_{n-1}^{[2]}(\xi)} = \frac{I_{2,n}(\xi)}{I_{3,n}(\xi)} = \frac{1}{\beta_n} < -\frac{n}{\xi}.$$

(ii) $N I_{2,n}(\xi) \hat{P}'_n(\xi) - M I_{3,n}(\xi) \hat{P}_n(\xi) = I_{2,n}(\xi) \hat{P}'_n(\xi) (N + M \beta_n \gamma_n)$, where

$$\frac{d_0^1}{c_0} < -\frac{\hat{P}'_n(\xi)}{\hat{P}_n(\xi)} = \frac{1}{\gamma_n} < -\frac{n}{\xi}.$$

(iii) $M N \hat{P}_n(\xi) \hat{P}_{n-1}^{[2]}(\xi) - N I_{1,n}(\xi) \hat{P}'_n(\xi) = N \hat{P}_n(\xi) \hat{P}_{n-1}^{[2]}(\xi) \left(M + \frac{\alpha_n}{\gamma_n} \right)$.

Proof. (i) From the Christoffel–Darboux formula for polynomials $\{\hat{P}_n^{[2]}\}_{n \geq 0}$ we have

$$\begin{aligned} (16) \quad (x - \xi) \sum_{k=0}^n \frac{\hat{P}_k^{[2]}(x) \hat{P}_k^{[2]}(y)}{\|\hat{P}_k^{[2]}\|_{\omega,2}^2} - \sum_{k=0}^n \frac{\hat{P}_k^{[2]}(x)}{\|\hat{P}_k^{[2]}\|_{\omega,2}^2} (y - \xi) \hat{P}_k^{[2]}(y) \\ = \frac{1}{\|\hat{P}_n^{[2]}\|_{\omega,2}^2} (\hat{P}_{n+1}^{[2]}(x) \hat{P}_n^{[2]}(y) - \hat{P}_n^{[2]}(x) \hat{P}_{n+1}^{[2]}(y)). \end{aligned}$$

If we multiply (16) by $(y - \xi)\omega(y)$ and integrate over $(0, \infty)$, evaluation at $x = \xi$ yields

$$\begin{aligned}
 - \sum_{k=0}^n \frac{\hat{P}_k^{[2]}(\xi)}{\|\hat{P}_k^{[2]}\|_{\omega,2}^2} \int_0^\infty (y - \xi)^2 \omega(y) \hat{P}_k^{[2]}(y) dy \\
 = \frac{1}{\|\hat{P}_n^{[2]}\|_{\omega,2}^2} (\hat{P}_{n+1}^{[2]}(\xi) I_{1,n+1}(\xi) - \hat{P}_n^{[2]}(\xi) I_{1,n+2}(\xi)).
 \end{aligned}$$

Since

$$\int_0^\infty (y - \xi)^2 \omega(y) \hat{P}_k^{[2]}(y) dy = 0 \quad \text{for } k = 1, 2, \dots, n$$

and $\hat{P}_0^{[2]} = 1$, the left-hand side is negative. Therefore,

$$\hat{P}_{n+1}^{[2]}(\xi) I_{1,n+1}(\xi) - \hat{P}_n^{[2]}(\xi) I_{1,n+2}(\xi) < 0.$$

From (5) we have

$$\text{sgn } \hat{P}_{n+1}^{[2]}(\xi) = (-1)^{n+1} \quad \text{and} \quad \text{sgn } \hat{P}_n^{[2]}(\xi) = (-1)^n.$$

Thus, $\hat{P}_{n+1}^{[2]}(\xi) \hat{P}_n^{[2]}(\xi)$ is negative and, as a consequence,

$$\frac{I_{1,n+2}(\xi)}{\hat{P}_{n+1}^{[2]}(\xi)} < \frac{I_{1,n+1}(\xi)}{\hat{P}_n^{[2]}(\xi)}.$$

Using this relation recursively, we get

$$\frac{I_{1,n}(\xi)}{\hat{P}_{n-1}^{[2]}(\xi)} < I_{1,1}(\xi) = d_0^1.$$

On the other hand, (5) and (13) imply that $\text{sgn } I_{1,n}(\xi) = (-1)^{n+1}$; therefore,

$$0 < \frac{I_{1,n}(\xi)}{\hat{P}_{n-1}^{[2]}(\xi)} < d_0^1.$$

From (16)

$$0 < \sum_{k=0}^n \frac{(\hat{P}_k^{[2]}(\xi))^2}{\|\hat{P}_k^{[2]}\|_{\omega,2}^2} = \frac{1}{\|\hat{P}_n^{[2]}\|_{\omega,2}^2} (\hat{P}_{n+1}^{[2]'}(\xi) \hat{P}_n^{[2]}(\xi) - \hat{P}_n^{[2]'}(\xi) \hat{P}_{n+1}^{[2]}(\xi)).$$

Since $\hat{P}_{n+1}^{[2]}(\xi) \hat{P}_n^{[2]}(\xi)$ is negative this yields

$$\frac{\hat{P}_{n+1}^{[2]'}(\xi)}{\hat{P}_{n+1}^{[2]}(\xi)} < \frac{\hat{P}_n^{[2]'}(\xi)}{\hat{P}_n^{[2]}(\xi)}.$$

Using this relation recursively, we obtain

$$\frac{\hat{P}_{n+1}^{[2]'}(\xi)}{\hat{P}_{n+1}^{[2]}(\xi)} < \frac{\hat{P}_1^{[2]'}(\xi)}{\hat{P}_1^{[2]}(\xi)} = -\frac{d_0^3}{d_0^2}.$$

Let $0 < x_{1,n}^{[2]} < x_{2,n}^{[2]} < \dots < x_{n,n}^{[2]}$ denote the zeros of $\hat{P}_n^{[2]}$. Then

$$-\frac{\hat{P}_n^{[2]'}(\xi)}{\hat{P}_n^{[2]}(\xi)} = \frac{1}{x_{1,n}^{[2]} - \xi} + \frac{1}{x_{2,n}^{[2]} - \xi} + \dots + \frac{1}{x_{n,n}^{[2]} - \xi} < -\frac{n}{\xi}.$$

Statements (ii) and (iii) can be proved in a similar way as (i). □

Proposition 5. *Let $M, N \geq 0$ and not both zero. Then*

$$\operatorname{sgn} A_{n,1} = -1 \quad \text{and} \quad \operatorname{sgn} A_{n,2} = -\operatorname{sgn} N.$$

Proof. From (5) and Proposition 4

$$\operatorname{sgn} A_{n,1} = -\operatorname{sgn} \frac{\hat{P}_n'(\xi)}{\hat{P}_{n-1}^{[2]}(\xi)} = \operatorname{sgn} \left(-\frac{\hat{P}_n'(\xi)}{\hat{P}_n(\xi)} \right) \operatorname{sgn} \frac{\hat{P}_n(\xi)}{\hat{P}_{n-1}^{[2]}(\xi)} = -1.$$

In a similar way,

$$\begin{aligned} \operatorname{sgn} A_{n,2} &= -\operatorname{sgn} N \operatorname{sgn} \frac{\hat{P}_n(\xi)}{I_{2,n}} \\ &= \operatorname{sgn} N \operatorname{sgn} \left(-\frac{\hat{P}_{n-1}^{[2]}(\xi)}{\hat{P}_{n-1}^{[2]'}(\xi)} \right) \operatorname{sgn} \frac{\hat{P}_n(\xi) \hat{P}_{n-2}(\xi) \hat{P}_{n-2}^{[1]}(\xi) \hat{P}_{n-2}^{[2]}(\xi) \hat{P}_{n-2}^{[3]}(\xi)}{\hat{P}_{n-1}(\xi) \hat{P}_{n-1}^{[1]}(\xi) \hat{P}_{n-1}^{[2]}(\xi)} \\ &= -\operatorname{sgn} N. \end{aligned} \quad \square$$

The zeros. We now analyze the zeros of the polynomials \hat{S}_n . The techniques are the same as those used by Meijer [1993a; 1993b].

Theorem 2. *The discrete Sobolev orthogonal polynomial \hat{S}_n has n real simple zeros and at most one of them is outside of $[\xi, \infty)$.*

Proof. Since for $N = 0$, \hat{S}_n is a standard orthogonal polynomial, in the sequel we will consider the cases when $N > 0$ and $M \geq 0$. Let $v_1 < v_2 < \dots < v_k$ be the zeros of $\hat{S}_n(x)$ on (ξ, ∞) with odd multiplicity. Let us introduce the polynomial

$$\phi(x) = (x - v_1)(x - v_2) \cdots (x - v_k).$$

Notice that $\phi(\xi)$ and $\phi'(\xi)$ have opposite signs and $\phi(x)\hat{S}_n(x)$ does not change sign on $[\xi, \infty)$. If $\deg \phi \leq n - 2$, then

$$0 = \langle \phi, \hat{S}_n \rangle_S = \int_0^\infty \omega(x) \phi(x) \hat{S}_n(x) dx + M\phi(\xi) \hat{S}_n(\xi) + N\phi'(\xi) \hat{S}_n'(\xi)$$

and

$$0 = \langle (\cdot - \xi)\phi, \hat{S}_n \rangle_S = \int_0^\infty \omega(x)(x - \xi)\phi(x)\hat{S}_n(x) dx + N\phi(\xi)\hat{S}'_n(\xi).$$

This means that $\phi'(\xi)\hat{S}'_n(\xi)$ and $\phi(\xi)\hat{S}'_n(\xi)$ have the same sign, and therefore $\phi'(\xi)$ and $\phi(\xi)$ have the same sign. This yields a contradiction.

As a conclusion, $\deg \phi = n - 1$ or $\deg \phi = n$, which proves our statement. \square

Next, we prove that the zeros of $\hat{S}_n(x)$ interlace with the zeros of $\hat{P}_{n-1}^{[2]}(x)$ if $\hat{S}_n(x)$ has a zero outside $[\xi, \infty)$. Notice that, by Theorem 1, $\hat{S}_n(\xi) \neq 0$.

Theorem 3. Denote by $v_{r,n}$, $r = 1, 2, \dots, n$, the zeros of $\hat{S}_n(x)$ in increasing order. Suppose that $v_{1,n} < \xi$. Then $2\xi - x_{1,n-1}^{[2]} < v_{1,n} < \xi$ and

$$\xi < v_{2,n} < x_{1,n-1}^{[2]} < \dots < v_{n,n} < x_{n-1,n-1}^{[2]}.$$

Proof. From Theorem 1 we have

$$\hat{S}_n(x_{r,n-1}^{[2]}) = \hat{P}_n(x_{r,n-1}^{[2]}) + A_{n,2}(x_{r,n-1}^{[2]} - \xi)^2 \hat{P}_{n-2}^{[4]}(x_{r,n-1}^{[2]}), \quad r = 1, 2, \dots, n - 1.$$

Then from Proposition 3(iii) and Proposition 5 we get

$$\text{sgn } \hat{S}_n(x_{r,n-1}^{[2]}) = (-1)^{n-r}, \quad r = 1, 2, \dots, n - 1,$$

On the other hand, from (5) and Theorem 1,

$$\text{sgn } \hat{S}_n(\xi) = \text{sgn } \hat{P}_n(\xi) = (-1)^n.$$

Therefore, every interval $(\xi, x_{1,n-1}^{[2]})$ and $(x_{r,n-1}^{[2]}, x_{r+1,n-1}^{[2]})$, for $r = 1, \dots, n - 2$, contains an odd number of zeros of $\hat{S}_n(x)$. Since \hat{S}_n has n real zeros and at most one of them is outside of (ξ, ∞) , then

$$\xi < v_{2,n} < x_{1,n-1}^{[2]} < \dots < v_{n,n} < x_{n-1,n-1}^{[2]}.$$

Now, we will prove that $2\xi - x_{1,n-1}^{[2]} < v_{1,n} < \xi$. Let

$$\hat{S}_n(x) = (x - v_{1,n})(x - v_{2,n}) \cdots (x - v_{n,n}).$$

By Theorem 1 and Proposition 4,

$$\hat{S}'_n(\xi) = \hat{P}'_n(\xi) + A_{n,1}\hat{P}_{n-2}^{[2]}(\xi) = \frac{\beta_n \hat{P}_n(\xi)(M + \alpha_n/\gamma_n)}{N + \alpha_n \beta_n}.$$

Therefore,

$$\text{sgn } \hat{S}'_n(\xi) = \text{sgn } \hat{P}_n(\xi) = \text{sgn } \hat{S}_n(\xi)$$

and

$$0 < \frac{\hat{S}'_n(\xi)}{\hat{S}_n(\xi)} = \frac{1}{\xi - v_{1,n}} - \frac{1}{v_{2,n} - \xi} - \dots - \frac{1}{v_{n,n} - \xi}.$$

Hence $\frac{1}{\xi - \nu_{1,n}} > \frac{1}{\nu_{2,n} - \xi}$, which implies successively

$$x_{1,n-1}^{[2]} - \xi > \nu_{2,n} - \xi > \xi - \nu_{1,n} \quad \text{and} \quad 2\xi - x_{1,n-1}^{[2]} < \nu_{1,n}.$$

Our statement follows. □

4. Discrete Laguerre–Sobolev orthogonal polynomials: asymptotics

Laguerre polynomials. For $\alpha \in \mathbb{R}$, the Laguerre polynomials are defined by

$$L_n^{(\alpha)}(x) = \sum_{k=0}^n \binom{n+\alpha}{n-k} \frac{(-x)^k}{k!}.$$

For $\alpha > -1$, the $\{L_n^{(\alpha)}(x)\}_{n \geq 0}$ are orthogonal on $[0, +\infty)$ with respect to the weight function $\omega(x) = x^\alpha e^{-x}$ [Szegő 1975, Chapter V]. Let $\{L_n^{(\alpha,k)}\}_{n=0}^\infty$, $k \in \mathbb{N}$, denote the sequence of polynomials orthogonal with respect to the modified Laguerre weight $(x - \xi)^k \omega(x)$, $\xi < 0$, normalized by the condition that $L_n^{(\alpha,k)}$ has the same leading coefficient as the classical Laguerre orthogonal polynomial $L_n^{(\alpha)} = L_n^{(\alpha,0)}$. That is, $k(L_n^{(\alpha,k)}) = (-1)^n/n!$.

We summarize some properties of the $L_n^{(\alpha,k)}(x)$, $k \in \mathbb{N} \cup \{0\}$, to be used later.

Proposition 6 [Fejzullahu 2011]. (i) For $\alpha > -1$,

$$\|L_n^{(\alpha)}\|_\alpha^2 = \int_0^\infty (L_n^{(\alpha)}(x))^2 x^\alpha e^{-x} dx = \frac{\Gamma(n+\alpha+1)}{\Gamma(n+1)}.$$

(ii) For every $n \in \mathbb{N}$,

$$(L_n^{(\alpha)}(x))' = -L_{n-1}^{(\alpha+1)}(x).$$

(iii) (Perron’s formula) Let $\alpha \in \mathbb{R}$. Then

$$L_n^{(\alpha)}(x) = 2^{-1} \pi^{-1/2} e^{x/2} (-x)^{-\alpha/2-1/4} n^{\alpha/2-1/4} e^{2\sqrt{-nx}} (1 + O(n^{-1/2})).$$

This relation holds for x in the complex plane cut along the positive real semiaxis; both $(-x)^{-\alpha/2-1/4}$ and $\sqrt{-x}$ must be taken real and positive if $x < 0$. The bound of the remainder holds uniformly in every closed domain which does not overlap the positive real semiaxis.

Moreover, we get the outer ratio asymptotics

$$\lim_{n \rightarrow \infty} n^{(l-j)/2} \frac{L_{n+k}^{(\alpha+j)}(x)}{L_{n+h}^{(\alpha+l)}(x)} = (-x)^{(l-j)/2}, \quad j, l \in \mathbb{R}, \quad h, k \in \mathbb{Z},$$

$$\lim_{n \rightarrow \infty} \frac{L_n^{(\alpha,k)}(x)}{n^{k/2} L_n^{(\alpha)}(x)} = \frac{1}{(\sqrt{-x} + \sqrt{-\xi})^k},$$

uniformly on compact subsets of $\mathbb{C} \setminus [0, \infty)$.

(iv) (Mehler–Heine formula) *Uniformly on compact subsets of \mathbb{C} , we have*

$$\lim_{n \rightarrow \infty} \frac{L_n^{(\alpha)}(x/(n+j))}{n^\alpha} = x^{-\alpha/2} J_\alpha(2\sqrt{x})$$

and

$$\lim_{n \rightarrow \infty} \frac{L_n^{(\alpha,k)}(x/(n+j))}{n^{\alpha+k/2}} = \frac{1}{(\sqrt{-\xi})^k} x^{-\alpha/2} J_\alpha(2\sqrt{x})$$

where $j \in \mathbb{N} \cup 0$ and J_α is the Bessel function of the first kind.

(v) (Plancherel–Rotach type outer asymptotics for $L_n^{(\alpha,N)}$) *Uniformly on compact subsets of $\mathbb{C} \setminus [0, 4]$ and uniformly on $j \in \mathbb{N} \cup \{0\}$, we have*

$$\lim_{n \rightarrow \infty} \frac{L_{n-1}^{(\alpha)}((n+j)x)}{L_n^{(\alpha)}((n+j)x)} = -\frac{1}{\phi((x-2)/2)}$$

and

$$\lim_{n \rightarrow \infty} \frac{L_n^{(\alpha,N)}((n+j)x)}{L_n^{(\alpha)}((n+j)x)} = \left(\frac{\phi((x-2)/2) + 1}{x} \right)^N,$$

where ϕ is the conformal mapping of $\mathbb{C} \setminus [-1, 1]$ onto the exterior of the unit circle given by

$$\phi(x) = x + \sqrt{x^2 - 1}, \quad x \in \mathbb{C} \setminus [-1, 1],$$

with $\sqrt{x^2 - 1} > 0$ when $x > 1$.

Proposition 7. $L_n^{(\alpha,2)'}(\xi) \cong \frac{n}{4\xi} L_n^{(\alpha+1)}(\xi).$

Proof. Using integration by parts we have

$$\int_0^\infty (L_n^{(\alpha,2)}(x))' L_k^{(\alpha+1,3)}(x) (x-\xi)^3 x^{\alpha+1} e^{-x} dx = \begin{cases} 0 & \text{if } k \leq n-3, \\ n(n-1) \|\hat{L}_n^{(\alpha,2)}\|_{\alpha,2}^2 & \text{if } k = n-2. \end{cases}$$

Therefore,

$$(L_n^{(\alpha,2)}(x))' = -L_{n-1}^{(\alpha+1,3)}(x) + H_n L_{n-2}^{(\alpha+1,3)}(x),$$

where

$$H_n = \frac{n(n-1) \|\hat{L}_n^{(\alpha,2)}\|_{\alpha,2}^2}{\|\hat{L}_{n-2}^{(\alpha+1,3)}\|_{\alpha+1,3}^2}.$$

Using (8) and Proposition 6(iii),

$$\begin{aligned} H_n &= \frac{(n+1)^2(n+\alpha)}{(n-1)^3} \frac{L_{n-2}^{(\alpha+1,2)}(\xi)}{L_{n-1}^{(\alpha+1,2)}(\xi)} \prod_{i=1}^2 \frac{L_{n-2}^{(\alpha+1,i-1)}(\xi)}{L_{n-1}^{(\alpha+1,i-1)}(\xi)} \frac{L_{n+1}^{(\alpha,i-1)}(\xi)}{L_n^{(\alpha,i-1)}(\xi)} \\ &= \frac{L_{n-2}^{(\alpha+1,2)}(\xi)}{L_{n-1}^{(\alpha+1,2)}(\xi)} \prod_{i=1}^2 \frac{L_{n-2}^{(\alpha+1,i-1)}(\xi)}{L_{n-1}^{(\alpha+1,i-1)}(\xi)} \frac{L_{n+1}^{(\alpha,i-1)}(\xi)}{L_n^{(\alpha,i-1)}(\xi)} + O\left(\frac{1}{n}\right). \end{aligned}$$

On the other hand, [Fejzullahu 2011, Proposition 2.2] gives

$$(17) \quad (L_n^{(\alpha,2)}(x))' = -L_{n-1}^{(\alpha,3)}(x) + G_n L_{n-2}^{(\alpha+1,3)}(x),$$

where

$$\begin{aligned} G_n &= H_n - \frac{n^3}{(n-1)^3} \prod_{i=1}^3 \frac{L_{n-2}^{(\alpha+1,i-1)}(\xi)}{L_{n-1}^{(\alpha+1,i-1)}(\xi)} \frac{L_n^{(\alpha,i-1)}(\xi)}{L_{n-1}^{(\alpha,i-1)}(\xi)} \\ &= \prod_{i=1}^3 \frac{L_{n-2}^{(\alpha+1,i-1)}(\xi)}{L_{n-1}^{(\alpha+1,i-1)}(\xi)} \left(\frac{L_{n+1}^{(\alpha)}(\xi) L_{n+1}^{(\alpha,1)}(\xi)}{L_n^{(\alpha)}(\xi) L_n^{(\alpha,1)}(\xi)} - \frac{L_n^{(\alpha)}(\xi) L_n^{(\alpha,1)}(\xi) L_n^{(\alpha,2)}(\xi)}{L_{n-1}^{(\alpha)}(\xi) L_{n-1}^{(\alpha,1)}(\xi) L_{n-1}^{(\alpha,2)}(\xi)} \right) + O\left(\frac{1}{n}\right). \end{aligned}$$

Again, from [Fejzullahu 2011, Proposition 2.2],

$$\begin{aligned} \frac{L_{n+1}^{(\alpha)}(\xi) L_{n+1}^{(\alpha,1)}(\xi)}{L_n^{(\alpha)}(\xi) L_n^{(\alpha,1)}(\xi)} &= \frac{L_{n+1}^{(\alpha)}(\xi) L_{n+1}^{(\alpha-1,1)}(\xi)}{L_n^{(\alpha)}(\xi) L_n^{(\alpha,1)}(\xi)} + \frac{L_{n+2}^{(\alpha-1)}(\xi)}{L_{n+1}^{(\alpha-1)}(\xi)} + O\left(\frac{1}{n}\right), \\ \frac{L_n^{(\alpha)}(\xi) L_n^{(\alpha,1)}(\xi) L_n^{(\alpha,2)}(\xi)}{L_{n-1}^{(\alpha)}(\xi) L_{n-1}^{(\alpha,1)}(\xi) L_{n-1}^{(\alpha,2)}(\xi)} &= \frac{L_n^{(\alpha)}(\xi) L_n^{(\alpha,1)}(\xi) L_n^{(\alpha-1,2)}(\xi)}{L_{n-1}^{(\alpha)}(\xi) L_{n-1}^{(\alpha,1)}(\xi) L_{n-1}^{(\alpha,2)}(\xi)} \\ &\quad + \frac{L_{n+1}^{(\alpha-1)}(\xi) L_{n+1}^{(\alpha-1,1)}(\xi)}{L_n^{(\alpha-1)}(\xi) L_n^{(\alpha-1,1)}(\xi)} + O\left(\frac{1}{n}\right), \end{aligned}$$

and

$$\begin{aligned} \frac{L_{n+2}^{(\alpha-1)}(\xi)}{L_{n+1}^{(\alpha-1)}(\xi)} - \frac{L_{n+1}^{(\alpha-1)}(\xi) L_{n+1}^{(\alpha-1,1)}(\xi)}{L_n^{(\alpha-1)}(\xi) L_n^{(\alpha-1,1)}(\xi)} &= \frac{L_{n+2}^{(\alpha-2)}(\xi)}{L_{n+1}^{(\alpha-1)}(\xi)} + 1 \\ &\quad - \frac{L_{n+1}^{(\alpha-1)}(\xi) L_{n+1}^{(\alpha-2,1)}(\xi)}{L_n^{(\alpha-1)}(\xi) L_n^{(\alpha-1,1)}(\xi)} - \frac{L_{n+2}^{(\alpha-2)}(\xi)}{L_{n+1}^{(\alpha-1)}(\xi)} + O\left(\frac{1}{n}\right) \\ &= \frac{L_{n+2}^{(\alpha-2)}(\xi)}{L_{n+1}^{(\alpha-1)}(\xi)} - \frac{L_{n+1}^{(\alpha-1)}(\xi) L_{n+1}^{(\alpha-2,1)}(\xi)}{L_n^{(\alpha-1)}(\xi) L_n^{(\alpha-1,1)}(\xi)} - \frac{L_{n+2}^{(\alpha-3)}(\xi)}{L_{n+1}^{(\alpha-2)}(\xi)} + O\left(\frac{1}{n}\right). \end{aligned}$$

Therefore, by using Proposition 6(iii),

$$\sqrt{n}G_n \cong -\sqrt{-\xi}.$$

and taking into account (17) the result follows. □

Discrete Laguerre–Sobolev orthogonal polynomials. Let $\{S_n\}_{n \geq 0}$ be the sequence of polynomials orthogonal with respect to the discrete Sobolev inner product (12), where $\omega(x) = x^\alpha e^{-x}$ and $\xi < 0$, normalized by the condition that S_n has the same leading coefficient as the classical Laguerre orthogonal polynomial $L_n^{(\alpha)}$, i.e., $k(S_n) = (-1)^n/n!$.

Theorem 4. *Let $M \geq 0$ and $N \geq 0$. There are real constants $B_{n,0}$, $B_{n,1}$, and $B_{n,2}$ such that*

$$(18) \quad S_n(x) = B_{n,0}L_n^{(\alpha)}(x) + B_{n,1}(x - \xi)L_{n-1}^{(\alpha,2)}(x) + B_{n,2}(x - \xi)^2L_{n-2}^{(\alpha,4)}(x),$$

where $B_{n,0} = \frac{1}{1 + A_{n,1} + A_{n,2}}$, $B_{n,1} = -\frac{A_{n,1}}{n(1 + A_{n,1} + A_{n,2})}$, and

$$B_{n,2} = \frac{A_{n,2}}{n(n-1)(1 + A_{n,1} + A_{n,2})}.$$

Moreover:

(i) *If $M > 0$ and $N > 0$, then*

$$(19) \quad B_{n,0} \cong \frac{8\xi n^\alpha}{M(L_n^{(\alpha)}(\xi))^2}, \quad B_{n,1} \cong -\frac{32\xi\sqrt{-\xi}n^{\alpha-1/2}}{M(L_n^{(\alpha)}(\xi))^2}, \quad B_{n,2} \cong \frac{1}{n^2}.$$

(ii) *If $M = 0$ and $N > 0$, then*

$$B_{n,0} \cong \frac{1}{4\sqrt{-\xi}n}, \quad B_{n,1} \cong -\frac{1}{n}, \quad B_{n,2} \cong \frac{1}{4n^2\sqrt{-\xi}n}.$$

(iii) *If $M > 0$ and $N = 0$, then*

$$B_{n,0} \cong \frac{\sqrt{-\xi}}{Mn^{1/2-\alpha}(L_{n-1}^{(\alpha)}(\xi))^2}, \quad B_{n,1} \cong -\frac{1}{n}, \quad B_{n,2} = 0.$$

Proof. From Theorem 1,

$$S_n(x) = \frac{(-1)^n \hat{S}_n(x)}{n!(1 + A_{n,1} + A_{n,2})}$$

and, as a consequence,

$$S_n(x) = B_{n,0}L_n^{(\alpha)}(x) + B_{n,1}(x - \xi)L_{n-1}^{(\alpha,2)}(x) + B_{n,2}(x - \xi)^2L_{n-2}^{(\alpha,4)}(x),$$

where $B_{n,0}$, $B_{n,1}$, and $B_{n,2}$ are as in the statement of the theorem.

Now, from Proposition 4 we can obtain the behavior of the coefficients $B_{n,0}$, $B_{n,1}$ and $B_{n,2}$ for n large enough. In order to estimate $A_{n,1}$ and $A_{n,2}$, first we

compute $\alpha_n \beta_n$, α_n / γ_n , $\beta_n \gamma_n$ and $I_{2,n}(\xi)$. From (13) and Proposition 6, we can write

$$\begin{aligned} \alpha_n \beta_n &= -\frac{I_{1,n}(\xi)}{\hat{L}_{n-1}^{(\alpha,2)'}(\xi)} = \frac{\hat{L}_n^{(\alpha)}(\xi)}{\hat{L}_{n-1}^{(\alpha)}(\xi) \hat{L}_{n-1}^{(\alpha,1)}(\xi) \hat{L}_{n-1}^{(\alpha,2)'}(\xi)} \|\hat{L}_{n-1}^{(\alpha)}\|_\alpha^2 \\ &= -\frac{\Gamma(n+\alpha)}{\Gamma(n)} \frac{nL_n^{(\alpha)}(\xi)}{L_{n-1}^{(\alpha)}(\xi) L_{n-1}^{(\alpha,1)}(\xi) L_{n-1}^{(\alpha,2)'}(\xi)} \cong \frac{8(-\xi)^{3/2} n^{\alpha-1/2}}{L_n^{(\alpha)}(\xi) L_n^{(\alpha+1)}(\xi)}, \\ \frac{\alpha_n}{\gamma_n} &= -\frac{I_{1,n}(\xi) \hat{L}_n^{(\alpha)'}(\xi)}{\hat{L}_n^{(\alpha)}(\xi) \hat{L}_{n-1}^{(\alpha,2)}(\xi)} = \frac{\hat{L}_n^{(\alpha)'}(\xi)}{\hat{L}_{n-1}^{(\alpha)}(\xi) \hat{L}_{n-1}^{(\alpha,1)}(\xi) \hat{L}_{n-1}^{(\alpha,2)}(\xi)} \|\hat{L}_{n-1}^{(\alpha)}\|_\alpha^2 \\ &= \frac{\Gamma(n+\alpha)}{\Gamma(n)} \frac{nL_{n-1}^{(\alpha+1)}(\xi)}{L_{n-1}^{(\alpha)}(\xi) L_{n-1}^{(\alpha,1)}(\xi) L_{n-1}^{(\alpha,2)}(\xi)} \cong \frac{8(-\xi)^{3/2} n^{\alpha-1/2} L_n^{(\alpha+1)}(\xi)}{\left(L_n^{(\alpha)}(\xi)\right)^3}, \\ \beta_n \gamma_n &= \alpha_n \beta_n \frac{\gamma_n}{\alpha_n} \cong \left(\frac{L_n^{(\alpha)}(\xi)}{L_n^{(\alpha+1)}(\xi)}\right)^2 \cong -\frac{\xi}{n}, \\ I_{2,n}(\xi) &\cong (-1)^{n-1} (n-2)! n^{\alpha+3} \frac{L_{n-1}^{(\alpha)}(\xi) L_{n-1}^{(\alpha,1)}(\xi) L_{n-1}^{(\alpha,2)'}(\xi)}{L_{n-2}^{(\alpha)}(\xi) L_{n-2}^{(\alpha,1)}(\xi) L_{n-2}^{(\alpha,2)}(\xi) L_{n-3}^{(\alpha,3)}(\xi)} \\ &\cong \frac{8\xi (-1)^{n-1} (n-2)! n^{\alpha+2}}{L_n^{(\alpha)}(\xi)}. \end{aligned}$$

Next, we will analyze the following three situations.

(i) Let $M > 0$ and $N > 0$. Then,

$$A_{n,1} \cong -\frac{\hat{L}_n^{(\alpha)'}(\xi)}{\hat{L}_{n-1}^{(\alpha,2)}(\xi)} = \frac{nL_n^{(\alpha)'}(\xi)}{L_{n-1}^{(\alpha,2)}(\xi)} = -\frac{nL_{n-1}^{(\alpha+1)}(\xi)}{L_{n-1}^{(\alpha,2)}(\xi)} \cong -4\sqrt{-\xi n}$$

and

$$A_{n,2} \cong -\frac{M \hat{L}_n^{(\alpha)}(\xi)}{I_{n,2}(\xi)} \cong \frac{M(L_n^{(\alpha)}(\xi))^2}{8\xi n^\alpha}.$$

Therefore,

$$B_{n,0} \cong \frac{8\xi n^\alpha}{M(L_n^{(\alpha)}(\xi))^2}, \quad B_{n,1} \cong \frac{32\xi \sqrt{-\xi} n^{\alpha-1/2}}{M(L_n^{(\alpha)}(\xi))^2}, \quad B_{n,2} \cong \frac{1}{n^2}.$$

(ii) Let $M = 0$ and $N > 0$. Then,

$$A_{n,1} \cong -4\sqrt{-\xi n} \quad \text{and} \quad A_{n,2} = -\frac{\hat{L}_n^{(\alpha)}(\xi)}{I_{n,2}(\xi)} \frac{\alpha_n}{\gamma_n} \cong -1.$$

Therefore,

$$B_{n,0} \cong -\frac{1}{4\sqrt{-\xi n}}, \quad B_{n,1} \cong -\frac{1}{n}, \quad B_{n,2} \cong \frac{1}{4n^2\sqrt{-\xi n}}.$$

(iii) Let $M > 0$ and $N = 0$. Then,

$$A_{n,1} = \frac{M\hat{L}_n^{(\alpha)}(\xi)}{I_{n,1}(\xi)} = -\frac{M\hat{L}_{n-1}^{(\alpha)}(\xi)\hat{L}_{n-1}^{(\alpha,1)}(\xi)}{\|L_{n-1}^{(\alpha)}\|_\alpha^2} \cong -\frac{Mn^{1/2-\alpha}}{\sqrt{-\xi}}(L_{n-1}^{(\alpha)}(\xi))^2, \quad A_{n,2} = 0.$$

Therefore,

$$B_{n,0} \cong -\frac{\sqrt{-\xi}}{Mn^{1/2-\alpha}(L_{n-1}^{(\alpha)}(\xi))^2}, \quad B_{n,1} \cong -\frac{1}{n}, \quad B_{n,2} = 0. \quad \square$$

Next we deduce several asymptotic properties for discrete Laguerre–Sobolev polynomials when $M, N \geq 0$. (For $M > 0$ and $N = 0$, the same asymptotic results for corresponding Laguerre-type polynomials has been deduced in [Dueñas et al. 2011] and [Fejzullahu and Zejnullahu 2010].)

Theorem 5. (i) (Outer relative asymptotics) *Uniformly on compact subsets of $\mathbb{C} \setminus [0, \infty)$ we have:*

- If $M > 0$ and $N > 0$, then

$$\lim_{n \rightarrow \infty} \frac{S_n(x)}{L_n^{(\alpha)}(x)} = \left(\frac{\sqrt{-x} - \sqrt{-\xi}}{\sqrt{-x} + \sqrt{-\xi}} \right)^2.$$

Notice that, according to the Hurwitz’s Theorem, the point ξ attracts two negative zeros of $S_n(x)$ for n large enough.

- If $M = 0$ and $N > 0$ or $M > 0$ and $N = 0$, then

$$\lim_{n \rightarrow \infty} \frac{S_n(x)}{L_n^{(\alpha)}(x)} = \frac{\sqrt{-x} - \sqrt{-\xi}}{\sqrt{-x} + \sqrt{-\xi}}.$$

Notice that, according to the Hurwitz’s Theorem, the point ξ attracts one negative zero of $S_n(x)$ for n large enough.

(ii) (Mehler–Heine formula)

- If $M > 0$ and $N > 0$

$$\lim_{n \rightarrow \infty} \frac{S_n(x/n)}{n^\alpha} = x^{-\alpha/2} J_\alpha(2\sqrt{x}),$$

- If $M = 0$ and $N > 0$ or $M > 0$ and $N = 0$

$$\lim_{n \rightarrow \infty} \frac{S_n(x/n)}{n^\alpha} = -x^{-\alpha/2} J_\alpha(2\sqrt{x}),$$

uniformly on compact subsets of \mathbb{C} .

(iii) (Plancherel–Rotach type outer asymptotics for S_n)

- If $M \geq 0$ and $N \geq 0$, then

$$\lim_{n \rightarrow \infty} \frac{S_n(nx)}{L_n^{(\alpha)}(nx)} = 1,$$

uniformly on compact subsets of $\mathbb{C} \setminus [0, 4]$.

Proof. We will prove the theorem when $M > 0$ and $N > 0$. The proofs of the other cases can be done in a similar way.

(i) From (18)

$$\frac{S_n(x)}{L_n^{(\alpha)}(x)} = B_{n,0} + nB_{n,1}(x - \xi) \frac{L_{n-1}^{(\alpha,2)}(x)}{nL_n^{(\alpha)}(x)} + n^2B_{n,2}(x - \xi)^2 \frac{L_{n-2}^{(\alpha,4)}(x)}{n^2L_n^{(\alpha)}(x)}.$$

Now, Proposition 6(iii) and (19) yield

$$\lim_{n \rightarrow \infty} \frac{S_n(x)}{L_n^{(\alpha)}(x)} = (x - \xi)^2 \lim_{n \rightarrow \infty} \frac{L_{n-2}^{(\alpha,4)}(x)}{n^2L_n^{(\alpha)}(x)} = \left(\frac{\sqrt{-x} - \sqrt{-\xi}}{\sqrt{-x} + \sqrt{-\xi}} \right)^2.$$

(ii) Scaling the variable as $x \rightarrow x/n$ in (18) then dividing by n^α we get

$$\begin{aligned} & \frac{S_n(x/n)}{n^\alpha} \\ &= B_{n,0} \frac{L_n^{(\alpha)}(x/n)}{n^\alpha} + nB_{n,1}(x/n - \xi) \frac{L_{n-1}^{(\alpha,2)}(x/n)}{n^{\alpha+1}} + n^2B_{n,2}(x/n - \xi)^2 \frac{L_{n-2}^{(\alpha,4)}(x/n)}{n^{\alpha+2}}. \end{aligned}$$

Now, Proposition 6(iv) and (19) yield

$$\lim_{n \rightarrow \infty} \frac{S_n(x/n)}{n^\alpha} = (-\xi)^2 \lim_{n \rightarrow \infty} \frac{L_{n-2}^{(\alpha,4)}(x)}{n^{\alpha+2}} = x^{-\alpha/2} J_\alpha(2\sqrt{x}).$$

(iii) Dividing (18) by $L_n^\alpha(x)$ then scaling the variable as $x \rightarrow nx$ we get

$$\begin{aligned} \frac{S_n(nx)}{L_n^{(\alpha)}(nx)} &= B_{n,0} + nB_{n,1} \frac{nx - \xi}{n} \frac{L_{n-1}^{(\alpha,2)}(nx)}{L_{n-1}^{(\alpha)}(nx)} \frac{L_{n-1}^{(\alpha)}(nx)}{L_n^{(\alpha)}(nx)} \\ &\quad + n^2B_{n,2} \frac{(nx - \xi)^2}{n^2} \frac{L_{n-2}^{(\alpha,4)}(nx)}{L_{n-2}^{(\alpha)}(nx)} \frac{L_{n-2}^{(\alpha)}(nx)}{L_n^{(\alpha)}(nx)}. \end{aligned}$$

From Proposition 6(v) and (19)

$$\lim_{n \rightarrow \infty} \frac{S_n(nx)}{L_n^{(\alpha)}(nx)} = x^2 \left(\frac{\phi((x-2)/2) + 1}{x} \right)^4 \frac{1}{(\phi((x-2)/2))^2}.$$

Now, using the fact that $(\phi(z) + 1)^2 = 2(z+1)\phi(z)$ if $|z| > 1$, we get our result. \square

Acknowledgements

The third author (BF) acknowledges the kind reception of the Departamento de Matemáticas of Universidad Carlos III de Madrid during his visit from December 2010 to March 2011. In this period the content of manuscript was discussed in several seminars. The work of the first (FM) and fourth (EH) authors has been supported by Dirección General de Investigación, Ministerio de Ciencia e Innovación of Spain, grant MTM2009-12740-C03-01. The authors thank the valuable comments by the referees which have contributed to improve the presentation of the manuscript.

References

- [Alfaro et al. 1992] M. Alfaro, F. Marcellán, M. L. Rezola, and A. Ronveaux, “On orthogonal polynomials of Sobolev type: algebraic properties and zeros”, *SIAM J. Math. Anal.* **23**:3 (1992), 737–757. MR 93g:42015 Zbl 0764.33003
- [Álvarez-Nodarse and Moreno-Balcázar 2004] R. Álvarez-Nodarse and J. J. Moreno-Balcázar, “Asymptotic properties of generalized Laguerre orthogonal polynomials”, *Indag. Math. (N.S.)* **15**:2 (2004), 151–165. MR 2005e:33003 Zbl 1064.41022
- [Bracciali et al. 2002] C. F. Bracciali, D. K. Dimitrov, and A. Sri Ranga, “Chain sequences and symmetric generalized orthogonal polynomials”, *J. Comput. Appl. Math.* **143**:1 (2002), 95–106. MR 2003f:33007 Zbl 1006.33008
- [Chihara 1978] T. S. Chihara, *An introduction to orthogonal polynomials*, Mathematics and its Applications **13**, Gordon and Breach, New York, 1978. MR 58 #1979 Zbl 0389.33008
- [Dueñas et al. 2011] H. Dueñas, E. J. Huertas, and F. Marcellán, “Analytic properties of Laguerre-type orthogonal polynomials”, *Integral Transforms Spec. Funct.* **22**:2 (2011), 107–122. MR 2011k:33029 Zbl 1213.33017
- [Fejzullahu 2011] B. X. Fejzullahu, “Asymptotics for orthogonal polynomials with respect to the Laguerre measure modified by a rational factor”, *Acta Sci. Math. (Szeged)* **77**:1-2 (2011), 73–85. MR 2841145 Zbl 05990930
- [Fejzullahu and Marcellán 2009] B. X. Fejzullahu and F. Marcellán, “A Cohen type inequality for Laguerre–Sobolev expansions”, *J. Math. Anal. Appl.* **352**:2 (2009), 880–889. MR 2010a:42093 Zbl 1160.42312
- [Fejzullahu and Zejnullahu 2010] B. X. Fejzullahu and R. X. Zejnullahu, “Orthogonal polynomials with respect to the Laguerre measure perturbed by the canonical transformations”, *Integral Transforms Spec. Funct.* **21**:8 (2010), 569–580. MR 2011i:33022 Zbl 1215.33007
- [Koekoek and Meijer 1993] R. Koekoek and H. G. Meijer, “A generalization of Laguerre polynomials”, *SIAM J. Math. Anal.* **24**:3 (1993), 768–782. MR 94b:33007 Zbl 0780.33007
- [López et al. 1995] G. López, F. Marcellán, and W. Van Assche, “Relative asymptotics for polynomials orthogonal with respect to a discrete Sobolev inner product”, *Constr. Approx.* **11**:1 (1995), 107–137. MR 96c:42051 Zbl 0840.42017
- [Marcellán and Moreno-Balcázar 2006] F. Marcellán and J. J. Moreno Balcázar, “Asymptotics and zeros of Sobolev orthogonal polynomials on unbounded supports”, *Acta Appl. Math.* **94**:2 (2006), 163–192. MR 2007i:42002 Zbl 1137.42312

- [Marcellán and Ronveaux 1990] F. Marcellán and A. Ronveaux, “On a class of polynomials orthogonal with respect to a discrete Sobolev inner product”, *Indag. Math. (N.S.)* **1**:4 (1990), 451–464. MR 92f:42029 Zbl 0732.42016
- [Marcellán and Van Assche 1993] F. Marcellán and W. Van Assche, “Relative asymptotics for orthogonal polynomials with a Sobolev inner product”, *J. Approx. Theory* **72**:2 (1993), 193–209. MR 94h:42037 Zbl 0771.42014
- [Meijer 1993a] H. G. Meijer, “Laguerre polynomials generalized to a certain discrete Sobolev inner product space”, *J. Approx. Theory* **73**:1 (1993), 1–16. MR 94d:42027 Zbl 0771.42015
- [Meijer 1993b] H. G. Meijer, “Zero distribution of orthogonal polynomials in a certain discrete Sobolev space”, *J. Math. Anal. Appl.* **172**:2 (1993), 520–532. MR 94a:42027 Zbl 0780.42016
- [Szegő 1975] G. Szegő, *Orthogonal polynomials*, 4th ed., Amer. Math. Soc. Colloq. Publ. **23**, Amer. Math. Soc., Providence, RI, 1975. MR 51 #8724 Zbl 0305.42011

Received June 20, 2011. Revised January 23, 2012.

FRANCISCO MARCELLÁN
DEPARTAMENTO DE MATEMÁTICAS
ESCUELA POLITÉCNICA SUPERIOR
UNIVERSIDAD CARLOS III DE MADRID
AVENIDA DE LA UNIVERSIDAD 30
28911 LEGANÉS
SPAIN
pacomarc@ing.uc3m.es

RAMADAN ZEJNULLAHU
FACULTY OF MATHEMATICS AND SCIENCES
UNIVERSITY OF PRISHTINA
MOTHER TERESA 5
10000 PRISHTINA
KOSOVO
zejnullahu@yahoo.com

BUJAR FEJZULLAHU
FACULTY OF MATHEMATICS AND SCIENCES
UNIVERSITY OF PRISHTINA
MOTHER TERESA 5
1000 PRISHTINA
KOSOVO
bujar.feizullahu@uni-pr.edu

EDMUNDO HUERTAS
DEPARTAMENTO DE MATEMÁTICAS
ESCUELA POLITÉCNICA SUPERIOR
UNIVERSIDAD CARLOS III DE MADRID
AVENIDA DE LA UNIVERSIDAD 30
28911 LEGANÉS
SPAIN
ehuertas@math.uc3m.es

GREEN VERSUS LEMPERT FUNCTIONS: A MINIMAL EXAMPLE

PASCAL THOMAS

The Lempert function for a set of poles in a domain of \mathbb{C}^n at a point z is obtained by taking a certain infimum over all analytic disks going through the poles and the point z ; it majorizes the corresponding multipole pluricomplex Green function. Coman proved that both coincide in the case of sets of two poles in the unit ball. We give an example of a set of three poles in the unit ball where this equality fails.

1. Introduction

Let Ω be a domain in \mathbb{C}^n , and $a_j \in \Omega$, $j = 1, \dots, N$. The pluricomplex Green function with logarithmic singularities at $S := \{a_1, \dots, a_N\}$ is defined by

$$G_S(z) := \sup\{u \in \text{PSH}(\Omega, \mathbb{R}_-) : u(z) \leq \log|z - a_j| + C_j, j = 1, \dots, N\},$$

where $\text{PSH}(\Omega, \mathbb{R}_-)$ stands for the set of all negative plurisubharmonic functions in Ω . When Ω is hyperconvex, this solves the Monge–Ampère equation with right hand side equal to $\sum_{i=1}^N \delta_{a_i}$.

Pluricomplex Green functions have been studied by many authors at different levels of generality. See [Demailly 1987; Zahariuta 1984; Lempert 1981; Lelong 1989; Lárusson and Sigurdsson 1998].

A deep result due to Poletsky [1993], and see also [Lárusson and Sigurdsson 1998; Edigarian 1997], is that the Green function may be computed from analytic disks:

$$(1-1) \quad G_S(z) = \inf \left\{ \sum_{\alpha: \varphi(\alpha) \in S} \log|\alpha| : \text{such that there exists } \varphi \in \mathcal{O}(\mathbb{D}, \Omega) \text{ with } \varphi(0) = z \right\}.$$

MSC2010: 32U35, 32F45.

Keywords: pluricomplex Green function, Lempert function, analytic disks, Schwarz Lemma.

However, it is tempting to pick only one $\alpha_j \in \varphi^{-1}(a_j)$ in the range $1 \leq j \leq N$, which motivated the Coman's definition of the Lempert function [2000]:

$$(1-2) \quad \ell_S(z) := \inf \left\{ \sum_{j=1}^N \log|\zeta_j| : \varphi(0) = z, \varphi(\zeta_j) = a_j, j = 1, \dots, N \right. \\ \left. \text{for some } \varphi \in \mathbb{C}(\mathbb{D}, \Omega) \right\},$$

where \mathbb{D} is the unit disc in \mathbb{C} .

One easily sees that $\ell_S(z) \geq G_S(z)$ without recourse to (1-1); the fact that equality holds when $N = 1$ and Ω is convex is part of Lempert's celebrated theorem [1981], which was, in fact, the starting point for many of the notions defined above; see also [Edigarian 1995]. Coman [2000] proved that equality holds when $N = 2$ and $\Omega = \mathbb{B}^2$, the unit ball of \mathbb{C}^2 . The goal of this note is to present an example that shows that this is as far as it can go.

Theorem 1.1. *There exists a set of 3 points $S \subset \mathbb{B}^2$ such that $\ell_S(z) > G_S(z)$ for some $z \in \mathbb{B}^2$.*

Other examples in the same vein have been found in [Carlehed and Wiegerinck 2003; Thomas and Trao 2003; Nikolov and Zwonek 2005]. The interesting features of this one are that it involves no multiplicities and is minimal in the ball. Examples with an arbitrary number of points can be deduced from it. Let $z_0 \in \mathbb{B}^2$ satisfy $\ell_S(z_0) - G_S(z_0) =: \varepsilon_0 > 0$. Consider $S' := S \cup \{a_4, \dots, a_N\}$ with all the a_j close enough to the boundary so that $\ell_{S'}(z_0) \geq \ell_S(z_0) - \varepsilon_0/2$ (the Schwarz lemma shows that $|\zeta_j| \rightarrow 1$ when $\varphi(\zeta_j) = a_j$ and $|a_j| \rightarrow 1$). Then $\ell_{S'}(z_0) > G_S(z_0) \geq G_{S'}(z_0)$, as was to be shown. (I thank Nikolai Nikolov for sharing this observation with me).

Moreover, the corresponding Green function can be recovered, up to a bounded error, by using an analytic disk with just one more preimage than the number of points: One of the points has exactly two preimages and each of the other two points, only one; see [Magnússon et al. 2012, §6.8.2, Lemma 6.16].

More specifically, the theorem will follow from a precise calculation in the bidisk \mathbb{D}^2 . Let $S_\varepsilon = \{(0, 0), (\rho(\varepsilon), 0), (0, \varepsilon)\} \subset \mathbb{D}^2$, where $\lim_{\varepsilon \rightarrow 0} \rho(\varepsilon)/\varepsilon = 0$.

Proposition 1.2. *There exists $C_1 > 0$ such that for any $\delta \in (0, 1/4)$ there exists $\varepsilon_0 = \varepsilon_0(z, \delta) > 0$ and $r_0 = r_0(\delta) > 0$ such that*

$$(1-3) \quad G_{S_\varepsilon}(z) \leq 2 \log|z_2| + C_1,$$

$$(1-4) \quad \ell_{S_\varepsilon}(z) \geq (2 - \delta) \log|z_2|.$$

for any ε with $|\varepsilon| < \varepsilon_0$ and any $z = (z_1, z_2) \in \mathbb{D}^2$ such that

$$(1-5) \quad \frac{1}{2}|z_2|^{3/2} \leq |z_1| \leq |z_2|^{3/2} \quad \text{and} \quad \|z\| < r_0.$$

Proof of Theorem 1.1. If U and V are domains, and $S \subset U \subset V$, then the definitions of the Green and Lempert functions imply that $G_S^U(z) \geq G_S^V(z)$ and $\ell_S^U(z) \geq \ell_S^V(z)$.

For $|\varepsilon|$ small enough, $S_\varepsilon \subset \mathbb{B}^2$. When $|z_1| = |z_2|^{3/2}$, so that z verifies (1-5), the inclusion $\mathbb{B}^2 \subset \mathbb{D}^2$ implies

$$\ell_{S_\varepsilon}^{\mathbb{B}^2}(z) \geq \ell_{S_\varepsilon}^{\mathbb{D}^2}(z) \geq (2 - \delta) \log |z_2|.$$

Using the fact that $\mathbb{D}^2/\sqrt{2} \subset \mathbb{B}^2$ and the invariance of the Green function under biholomorphic mappings, we have

$$G_{S_\varepsilon}^{\mathbb{B}^2}(z) \leq G_{S_\varepsilon}^{\mathbb{D}^2/\sqrt{2}}(z) = G_{\sqrt{2}S_\varepsilon}^{\mathbb{D}^2}(\sqrt{2}z) \leq 2 \log |z_2| + \log 2 + C_1.$$

The last inequality follows from the fact that $\sqrt{2}z$ still verifies (1-5), and $\sqrt{2}S_\varepsilon$ has the same form as S_ε , so we can apply (1-3).

Comparing the last two estimates, we see that $G_{S_\varepsilon}^{\mathbb{B}^2}(z) < \ell_{S_\varepsilon}^{\mathbb{B}^2}(z)$ for $|z_2|$ small enough and $|\varepsilon| < \varepsilon_0$. \square

Open questions

This example is minimal in the ball, in terms of number of poles; what is the situation for the bidisk? Are the Green and Lempert functions equal when one takes two poles, not lying on a line parallel to the coordinate axes? Do they at least have the same order of singularity as one pole tends to the other?

What is the precise order of the singularity of the limit as $\varepsilon \rightarrow 0$ of the Lempert function in this case? Looking at the available analytic disks that give the correct order of the singularity of the limit of the Green function, one finds $\frac{3}{2} \log |z_2|$, so one would hope that the proposition can still be proved at least for $\delta < 1/2$.

Do the analytic disks from [Magnússon et al. 2012] yield the Green function itself, without any bounded error term?

More generally, when one is given a finite number of points in a given bounded (hyperconvex) domain, is there a bound on the number of preimages required to attain the Green function in the Poletsky formula? For instance, is 4 the largest possible number of preimages required when looking at 3 points in the ball?

2. Upper estimate for the Green function

Proof of (1-3) of Proposition 1.2. The upper bound (1-3) follows from [Magnússon et al. 2012, §6.8.2, Lemma 6.16]. For the reader's convenience, and since that paper is not generally available, we repeat the proof here in the case that concerns us.

We now construct an analytic disk passing twice through one of the poles. Our disk will be a perturbation of the Neil parabola $\zeta \mapsto (\zeta^3, \zeta^2)$.

We write $s(\varepsilon) = \rho(\varepsilon)/\varepsilon = o(1)$.

Choose complex numbers λ and μ such that

$$\lambda^2 := \frac{z_1}{z_2(z_2 - \varepsilon)} \left(\frac{z_1}{z_2 - \varepsilon} + s(\varepsilon) \right) \quad \text{and} \quad \mu^2 := \varepsilon + \left(\frac{s(\varepsilon)}{2\lambda} \right)^2.$$

Let

$$\Psi_{\lambda, \mu}(\zeta) := \left((\lambda\zeta - \frac{1}{2}s(\varepsilon))(\zeta^2 - \mu^2), \zeta^2 - \left(\frac{s(\varepsilon)}{2\lambda} \right)^2 \right).$$

Then by construction $\Psi_{\lambda, \mu}(\mu) = \Psi_{\lambda, \mu}(-\mu) = (0, \varepsilon)$,

$$\Psi_{\lambda, \mu} \left(\frac{s(\varepsilon)}{2\lambda} \right) = (0, 0) \quad \text{and} \quad \Psi_{\lambda, \mu} \left(-\frac{s(\varepsilon)}{2\lambda} \right) = (\varepsilon s(\varepsilon), 0),$$

so we have a disk passing through all three poles of G_ε . Furthermore, choosing

$$\zeta_z := \frac{1}{\lambda} \left(\frac{z_1}{z_2 - \varepsilon} + \frac{s(\varepsilon)}{2} \right),$$

we have $\Psi_{\lambda, \mu}(\zeta_z) = z$. Notice that

$$\zeta_z^2 = \frac{z_2(z_2 - \varepsilon)}{z_1} \left(\frac{z_1}{z_2 - \varepsilon} + \frac{s(\varepsilon)}{2} \right)^2 \left(\frac{z_1}{z_2 - \varepsilon} + s(\varepsilon) \right)^{-1},$$

so for any $\eta > 0$ there exists $\varepsilon_0(\delta, \eta) > 0$ such that for $|\varepsilon| < \varepsilon_0(\delta, \eta)$

$$(2-1) \quad \left| |\zeta_z| - |z_2|^{1/2} \right| \leq \eta$$

for any z such that $\delta \leq \frac{1}{2}|z_2|^{3/2} \leq |z_1| \leq |z_2|^{3/2} \leq 1$. In particular, by choosing η small enough we ensure that $\zeta_z \in \mathbb{D}$. We need a more general fact.

Claim. *Let $\eta > 0$, and $\delta > 0$. Then there exists $\varepsilon_1 = \varepsilon_1(\delta, \eta) > 0$ such that for any ε with $|\varepsilon| \leq \varepsilon_1$, we have $\Psi_{\lambda, \mu}(D(0, 1 - \eta)) \subset \mathbb{D}^2$ for any z such that $\delta \leq \frac{1}{2}|z_2|^{3/2} \leq |z_1| \leq |z_2|^{3/2} \leq 1$.*

Proof. For $|\varepsilon| \leq \delta^{2/3}/2$, we have $|z_2|/2 \leq |z_2 - \varepsilon| \leq 2|z_2|$, so

$$|\lambda|^2 \geq \left| \frac{z_1}{2z_2^2} \right| \left(\left| \frac{z_1}{2z_2} \right| - |s(\varepsilon)| \right) \geq \left| \frac{z_1^2}{8z_2^3} \right| \geq \frac{1}{32}$$

for ε small enough. So when $|\zeta| \leq 1 - \eta$,

$$|\Psi_{\lambda, \mu, 2}(\zeta)| \leq (1 - \eta)^2 + 256|s(\varepsilon)|^2 < 1$$

for ε small enough.

In a similar way, given η' , for ε small enough depending on δ and η' , we have $|z_2| \leq (1 + \eta')|z_2 - \varepsilon|$, so

$$|\lambda|^2 \leq (1 + \eta')^2 \left| \frac{z_1}{z_2^2} \right| \left(\left| \frac{z_1}{z_2} \right| + \frac{|s(\varepsilon)|}{(1 + \eta')} \right) \leq (1 + \eta')^3 \left| \frac{z_1^2}{z_2^3} \right| \leq (1 + \eta')^3$$

for ε small enough. Choose η' so that $(1 + \eta')^3 = (1 + \eta)$. When $|\zeta| \leq 1 - \eta$,

$$|\Psi_{\lambda, \mu, 1}(\zeta)| \leq \left((1 + \eta)(1 - \eta) + \frac{1}{2}|s(\varepsilon)| \right) \left((1 - \eta)^2 + |\varepsilon| + 64^2 |s(\varepsilon)|^2 \right) < 1$$

for ε small enough. □

So now the function $v(\zeta) := G_\varepsilon(\Psi_{\lambda, \mu}((1 - \eta)\zeta))$ is negative and subharmonic on \mathbb{D} . Furthermore, it has logarithmic poles at the points

$$\pm \frac{\mu}{1 - \eta} \quad \text{and} \quad \pm \frac{s(\varepsilon)}{2\lambda(1 - \eta)};$$

in the cases when $\mu = 0$ or $s(\varepsilon) = 0$, we get a double logarithmic pole at the corresponding point.

Denote by $d_G(\zeta, \xi) := |(\zeta - \xi)/(1 - \zeta\bar{\xi})|$ the invariant (pseudohyperbolic) distance between points of the unit disk. Then

$$\begin{aligned} G_\varepsilon(z) = v(\zeta_z) &\leq \log d_G\left(\zeta_z, \frac{\mu}{1 - \eta}\right) + \log d_G\left(\zeta_z, -\frac{\mu}{1 - \eta}\right) \\ &\quad + \log d_G\left(\zeta_z, \frac{s(\varepsilon)}{2\lambda(1 - \eta)}\right) + \log d_G\left(\zeta_z, -\frac{s(\varepsilon)}{2\lambda(1 - \eta)}\right). \end{aligned}$$

By (2-1), choosing $m(\delta, \eta)$ accordingly, we have $G_\varepsilon(z) \leq 4 \log|z_2|^{1/2} + O(\eta)$ for $|\varepsilon| \leq m$. □

3. Lower estimate for the Lempert function

Proof of (1-4) of Proposition 1.2. The proof will follow the methods and notations of [Thomas 2007]. We will make repeated use of the involutive automorphisms of the unit disk given by $\phi_a(\zeta) := (a - \zeta)/(1 - \bar{a}\zeta)$ for $a \in \mathbb{D}$, which exchange 0 and a . Notice that the invariant (pseudohyperbolic) distance verifies

$$d_G(a, b) := |\phi_a(b)| = |\phi_b(a)|.$$

Write $\rho(\varepsilon) = \varepsilon s(\varepsilon)$ with $\lim_{\varepsilon \rightarrow 0} s(\varepsilon) = 0$.

We will assume that the conclusion fails. That is, for any $\delta \in (0, 1/4)$, there exist arbitrarily small values of $|z_2| = \max(|z_1|, |z_2|)$, and $|\varepsilon|$ such that

$$(3-1) \quad \ell_{S_\varepsilon}(z) < (2 - \delta) \log|z_2|.$$

After applying, for each analytic disk, an automorphism of the disk that exchanges the preimage of $(0, 0)$ and 0, the assumption implies that there exists a holomorphic map φ from \mathbb{D} to \mathbb{D}^2 and points $\zeta_j \in \mathbb{D}$, depending on z and ε , satisfying the conditions

$$(3-2) \quad \begin{aligned} \varphi(0) &= (0, 0), & \varphi(\zeta_1) &= (\varepsilon s(\varepsilon), 0), \\ \varphi(\zeta_0) &= (z_1, z_2), & \varphi(\zeta_2) &= (0, \varepsilon), \end{aligned}$$

with

$$(3-3) \quad \log|\zeta_0| + \log|\phi_{\zeta_0}(\zeta_1)| + \log|\phi_{\zeta_0}(\zeta_2)| \leq (2 - \delta) \log|z_2|.$$

The interpolation conditions in (3-2) are equivalent to the existence of holomorphic functions h_1 and h_2 from \mathbb{D} to itself such that

$$\varphi(\zeta) = (\zeta \phi_{\zeta_2}(\zeta) h_1(\zeta), \zeta \phi_{\zeta_1}(\zeta) h_2(\zeta)),$$

such that furthermore

$$(3-4) \quad h_1(\zeta_1) = \frac{\varepsilon s(\varepsilon)}{\zeta_1 \phi_{\zeta_2}(\zeta_1)} =: w_1,$$

$$(3-5) \quad h_1(\zeta_0) = \frac{z_1}{\zeta_0 \phi_{\zeta_2}(\zeta_0)} =: w_2,$$

$$(3-6) \quad h_2(\zeta_2) = \frac{\varepsilon}{\zeta_2 \phi_{\zeta_1}(\zeta_2)} =: w_4,$$

$$(3-7) \quad h_2(\zeta_0) = \frac{z_2}{\zeta_0 \phi_{\zeta_1}(\zeta_0)} =: w_3.$$

By the invariant Schwarz lemma, the existence of a holomorphic function h_1 mapping \mathbb{D} to itself and satisfying (3-4) and (3-5) is equivalent to

$$(3-8) \quad |w_1| < 1, \quad |w_2| < 1 \quad \text{and} \quad d_G(w_1, w_2) < d_G(\zeta_1, \zeta_0) = |\phi_{\zeta_1}(\zeta_0)|.$$

In the same way, the existence of h_2 is equivalent to

$$(3-9) \quad |w_3| < 1, \quad |w_4| < 1 \quad \text{and} \quad d_G(w_3, w_4) < d_G(\zeta_2, \zeta_0) = |\phi_{\zeta_2}(\zeta_0)|.$$

As in [Thomas 2007], we start by remarking that (3-3) can be rewritten as

$$(3-10) \quad -\log|w_2| - \log|w_3| = \log \left| \frac{\zeta_0 \phi_{\zeta_1}(\zeta_0)}{z_2} \right| + \log \left| \frac{\zeta_0 \phi_{\zeta_0}(\zeta_2)}{z_1} \right| \\ \leq \log|\zeta_0| + (2 - \delta) \log|z_2| - \log|z_1| - \log|z_2| \\ \leq \log|\zeta_0| - \left(\frac{1}{2} + \delta\right) \log|z_2| + \log 2,$$

by (1-5). We can rewrite this in a more symmetric fashion:

$$(3-11) \quad \log \frac{1}{|w_2|} + \log \frac{1}{|w_3|} + \log \frac{1}{|\zeta_0|} \leq \left(\frac{1}{2} + \delta\right) \log \frac{1}{|z_2|} + \log 2.$$

Since all terms are positive by (3-8) and (3-9), each of the terms on the left hand side is bounded by the right hand side.

We will proceed as follows: We have used the contradiction hypothesis (3-3) to prove that $|\zeta_0|$ and $|w_3|$ are relatively big. We will prove that $|\phi_{\zeta_2}(\zeta_0)|$ has to be relatively small, which by (3-9) forces $|w_4|$ to be roughly as large as $|w_3|$. This then allows us to bound $|\phi_{\zeta_1}(\zeta_2)|$ by a quantity that becomes as small as desired

when ε can be made small, and hence allows us to bound $|\phi_{\zeta_1}(\zeta_0)|$ by the triangle inequality.

The final contradiction will concern $w_2 = z_1/(\zeta_0\phi_{\zeta_2}(\zeta_0))$. On the one hand, (3-11) guarantees that it is not too small; but an explicit computation of the quotient w_1/w_4 shows that w_1 must be small, and by (3-8) and the estimate on $|\phi_{\zeta_1}(\zeta_0)|$, $|w_2|$ is small as well.

We provide the details. From (3-11),

$$(3-12) \quad \log|w_3| \geq \left(\frac{1}{2} + \delta\right) \log|z_2| - \log 2.$$

From (3-5) and (3-10),

$$(3-13) \quad \begin{aligned} \log|\phi_{\zeta_2}(\zeta_0)| &= \log|z_1/\zeta_0| - \log|w_2| \\ &\leq \log|z_1/\zeta_0| + \log|\zeta_0| - \left(\frac{1}{2} + \delta\right) \log|z_2| + \log 2 \\ &\leq (1 - \delta) \log|z_2| + \log 2. \end{aligned}$$

Since $\delta < 1/4$, (3-13) and (3-12) imply that $|\phi_{\zeta_2}(\zeta_0)| < \frac{1}{2}|w_3|$ for $|z_2| \leq r_1(\delta)$, so by (3-9) and the triangle inequality for d_G ,

$$(3-14) \quad |w_4| \geq \frac{1}{2}|w_3|.$$

We now prove that both ζ_1 and ζ_2 must be close to ζ_0 and even closer to each other. First, since (3-11) implies that $\log|\zeta_0| \geq \left(\frac{1}{2} + \delta\right) \log|z_2| - \log 2$, by (3-13), $|\phi_{\zeta_2}(\zeta_0)| \leq \frac{1}{2}|\zeta_0|$ for $|z_2| \leq r_2(\delta)$. By the triangle inequality for d_G ,

$$(3-15) \quad \frac{1}{2}|\zeta_0| \leq |\zeta_2| \leq \frac{3}{2}|\zeta_0|.$$

On the other hand, from (3-11),

$$\log|w_3| + \log|\zeta_0| \geq \left(\frac{1}{2} + \delta\right) \log|z_2| - \log 2, \quad \text{that is, } |w_3\zeta_0| \geq \frac{1}{2}|z_2|^{\delta+1/2}.$$

Therefore, applying (3-14) and (3-15),

$$(3-16) \quad |\phi_{\zeta_1}(\zeta_2)| = \left| \frac{\varepsilon}{\zeta_2 w_4} \right| \leq 4 \left| \frac{\varepsilon}{\zeta_0 w_3} \right| \leq 8|\varepsilon||z_2|^{-\delta-1/2}.$$

In particular, for

$$(3-17) \quad |\varepsilon| < \frac{1}{8}|z_2|^{3/2},$$

this implies $|\phi_{\zeta_1}(\zeta_2)| < |z_2|^{1-\delta}$, and by the triangle inequality,

$$(3-18) \quad |\phi_{\zeta_1}(\zeta_0)| < |\phi_{\zeta_2}(\zeta_0)| + |\phi_{\zeta_1}(\zeta_2)| < 3|z_2|^{1-\delta}.$$

We now establish the two (contradictory) estimates for w_2 . On the one hand, (3-11) implies that

$$(3-19) \quad \log|w_2| \geq \left(\frac{1}{2} + \delta\right) \log|z_2| - \log 2, \quad \text{that is, } |w_2| \geq \frac{1}{2}|z_2|^{\delta+1/2}.$$

On the other hand,

$$\left| \frac{w_1}{w_4} \right| = \left| \frac{\varepsilon s(\varepsilon)}{\zeta_1 \phi_{\zeta_2}(\zeta_1)} \frac{\zeta_2 \phi_{\zeta_1}(\zeta_2)}{\varepsilon} \right| = \left| s(\varepsilon) \frac{\zeta_2}{\zeta_1} \right|.$$

By the triangle inequality for d_G , when (3-17) holds, the lower bound in (3-15) and the corollary to (3-16) imply

$$|\zeta_1| \geq |\zeta_2| - |\phi_{\zeta_1}(\zeta_2)| \geq \frac{1}{2}|\zeta_0| - |z_2|^{1-\delta} \geq \frac{1}{4}|\zeta_0|$$

for $|z_2|$ small enough, because of (3-11) again. So finally, using the upper bound in (3-15), $|w_1/w_4| \leq 6|s(\varepsilon)|$. We choose $\varepsilon_0 < \frac{1}{8}|z_2|^{3/2}$ so that for any ε with $|\varepsilon| \leq \varepsilon_0$,

$$(3-20) \quad |s(\varepsilon)| < |z_2|^{1-\delta}.$$

The triangle inequality for d_G and (3-18) imply that when $|\varepsilon| \leq \varepsilon_0$,

$$|w_2| \leq |w_1| + |\phi_{\zeta_1}(\zeta_0)| \leq 6|s(\varepsilon)| + 3|z_2|^{1-\delta} \leq 9|z_2|^{1-\delta}.$$

Finally, if we choose $|z_2| \leq r_0(\delta)$, with

$$r_0(\delta) \leq \min(r_1(\delta), r_2(\delta)) \quad \text{and} \quad 9r_0(\delta)^{1-\delta} < \frac{1}{2}r_0(\delta)^{1/2+\delta},$$

we see that for any ε with $|\varepsilon| \leq \varepsilon_0$, this last bound contradicts (3-19). \square

Acknowledgments

Part of this work was done when I was a guest of the Hanoi University of Education in March 2011. I wish to thank my colleagues there for their hospitality, in particular Do Duc Thai and Nguyen Van Trao. A detailed exposition of the results in a workshop organized in Bedlewo by colleagues from the Jagiellonian University in Cracow resulted in some streamlining of the proof. Finally, the Semester in Complex Analysis and Spectral Theory of the Centre de Recerca Matemàtica at the Universitat Autònoma de Barcelona provided the occasion to mention the result in a talk, and to put the finishing touches on the submitted manuscript.

I thank Nguyen Van Trao for useful discussions on this topic, and the referee for pointing out and correcting a mistake in the original exposition.

References

- [Carlehed and Wiegerinck 2003] M. Carlehed and J. Wiegerinck, “Le cône des fonctions plurisous-harmoniques négatives et une conjecture de Coman”, *Ann. Polon. Math.* **80** (2003), 93–108. MR 2004e:32034 Zbl 1026.32066
- [Coman 2000] D. Coman, “The pluricomplex Green function with two poles of the unit ball of \mathbb{C}^n ”, *Pacific J. Math.* **194**:2 (2000), 257–283. MR 2001g:32081 Zbl 1015.32029

- [Demailly 1987] J.-P. Demailly, “Mesures de Monge–Ampère et mesures pluriharmoniques”, *Math. Z.* **194**:4 (1987), 519–564. MR 88g:32034 Zbl 0595.32006
- [Edigarian 1995] A. Edigarian, “A remark on the Lempert theorem”, *Univ. Iagel. Acta Math.* **32** (1995), 83–88. MR 96h:32034 Zbl 0833.32005
- [Edigarian 1997] A. Edigarian, “On definitions of the pluricomplex Green function”, *Ann. Polon. Math.* **67**:3 (1997), 233–246. MR 99f:32039 Zbl 0909.31007
- [Lárusson and Sigurdsson 1998] F. Lárusson and R. Sigurdsson, “Plurisubharmonic functions and analytic discs on manifolds”, *J. Reine Angew. Math.* **501** (1998), 1–39. MR 99e:32020 Zbl 0901.31004
- [Lelong 1989] P. Lelong, “Fonction de Green pluricomplexe et lemmes de Schwarz dans les espaces de Banach”, *J. Math. Pures Appl.* (9) **68**:3 (1989), 319–347. MR 91c:46065 Zbl 0633.32019
- [Lempert 1981] L. Lempert, “La métrique de Kobayashi et la représentation des domaines sur la boule”, *Bull. Soc. Math. France* **109** (1981), 427–474. MR 84d:32036 Zbl 0492.32025
- [Magnússon et al. 2012] J. I. Magnússon, A. Rashkovskii, R. Sigurdsson, and P. J. Thomas, “Limits of multipole pluricomplex Green functions”, *Internat. J. Math* **23**:6 (2012), #1250065. arXiv 1103.2296
- [Nikolov and Zwonek 2005] N. Nikolov and W. Zwonek, “On the product property for the Lempert function”, *Complex Var. Theory Appl.* **50**:12 (2005), 939–952. MR 2006g:32017 Zbl 1085.32006
- [Poletsky 1993] E. A. Poletsky, “Holomorphic currents”, *Indiana Univ. Math. J.* **42**:1 (1993), 85–144. MR 94c:32007 Zbl 0811.32010
- [Thomas 2007] P. J. Thomas, “An example of limit of Lempert functions”, *Vietnam J. Math.* **35**:3 (2007), 317–330. MR 2009h:32017 Zbl 1158.32015 arXiv math/0601642
- [Thomas and Trao 2003] P. J. Thomas and N. V. Trao, “Pluricomplex Green and Lempert functions for equally weighted poles”, *Ark. Mat.* **41**:2 (2003), 381–400. MR 2004m:32067 Zbl 1038.32028
- [Zahariuta 1984] V. P. Zahariuta, *Spaces of analytic functions and maximal plurisubharmonic functions*, thesis, Rostov State University, Rostov-on-Don, 1984.

Received June 20, 2011. Revised October 18, 2011.

PASCAL THOMAS
INSTITUT DE MATHÉMATIQUES DE TOULOUSE
UNIVERSITÉ DE TOULOUSE
UPS, INSA, UT1, UTM
31062 TOULOUSE
FRANCE
pthomas@math.univ-toulouse.fr
<http://www.math.univ-toulouse.fr/~thomas>

DIFFERENTIAL HARNACK INEQUALITIES FOR NONLINEAR HEAT EQUATIONS WITH POTENTIALS UNDER THE RICCI FLOW

JIA-YONG WU

We prove several differential Harnack inequalities for positive solutions to nonlinear backward heat equations with different potentials coupled with the Ricci flow. We also derive an interpolated Harnack inequality for the nonlinear heat equation under the ε -Ricci flow on a closed surface. These new Harnack inequalities extend the previous differential Harnack inequalities for linear heat equations with potentials under the Ricci flow.

1. Introduction and main results

Background. The study of differential Harnack estimates for parabolic equations originated with the work of P. Li and S.-T. Yau [1986], who first proved a gradient estimate for the heat equation via the maximum principle (though a precursory form of their estimate appeared in [Aronson and B enilan 1979]). Using their gradient estimate, the same authors derived a classical Harnack inequality by integrating the gradient estimate along space-time paths. This result was generalized to Harnack inequalities for some nonlinear heat-type equations in [Yau 1994] and for some non-self-adjoint evolution equations in [Yau 1995]. Recently, J. Li and X. Xu [2011] gave sharper local estimates than previous results for the heat equation on Riemannian manifolds with Ricci curvature bounded below. Surprisingly, R. Hamilton employed similar techniques to obtain Harnack inequalities for the Ricci flow [Hamilton 1993a], and the mean curvature flow [Hamilton 1995]. In dimension two, a differential Harnack estimate for the positive scalar curvature was proved in [Hamilton 1988], and then extended by B. Chow [1991a] when the scalar curvature changes sign. Similar techniques were used to obtain the Harnack inequalities for the Gauss curvature flow [Chow 1991b] and the Yamabe flow [Chow 1992]. H.-D. Cao [1992] proved a Harnack inequality for the K ahler–Ricci

This work is partially supported by the NSFC (No.11101267) and the Science and Technology Program of Shanghai Maritime University (No. 20120061).

MSC2010: 53C44.

Keywords: Harnack inequality, interpolated Harnack inequality, nonlinear heat equation, nonlinear backward heat equation, Ricci flow.

flow. B. Andrews [1994] derived several Harnack inequalities for general curvature flows of hypersurfaces. Chow and Hamilton [1997] gave extensions of the Li–Yau Harnack inequality, which they called constrained and linear Harnack inequalities. For more detailed discussion, we refer the interested reader to [Chow et al. 2006, Chapter 10].

Hamilton [1993b] also generalized the Li–Yau Harnack inequality to a matrix Harnack form on a class of Riemannian manifolds with nonnegative sectional curvature. This result was extended to the constrained matrix Harnack inequalities in [Chow and Hamilton 1997]. H.-D. Cao and L. Ni [2005] proved a matrix Harnack estimate for the heat equation on Kähler manifolds. Chow and Ni [2007] proved a matrix Harnack estimate for Kähler–Ricci flow using interpolation techniques from [Chow 1998].

In another direction, differential Harnack inequalities for (backward) heat-type equations coupled with the Ricci flow have become an important object, which can be traced back to [Hamilton 1988]. This subject was further explored by Chow [1998], Chow and Hamilton [1997], Chow and D. Knopf [2002], and H.-B. Cheng [2006], among others. Perhaps the most spectacular result is G. Perelman’s [2002] differential Harnack inequality for the fundamental solution to the backward heat equation coupled with the Ricci flow without any curvature assumption. Perelman’s Harnack inequality has many important applications (it is essential in proving pseudolocality theorems), and it has been extended by X. Cao [2008] and independently by S.-L. Kuang and Qi S. Zhang [2008]. Those authors proved a differential Harnack inequality for all positive solutions to the backward heat equation under the Ricci flow on closed manifolds with nonnegative scalar curvature. X. Cao and Qi S. Zhang [2011a] have established Gaussian upper and lower bounds for the fundamental solution to the backward heat equation under the Ricci flow.

On the subject of differential Harnack inequalities for the linear heat equation coupled with the Ricci flow, there have been many important contributions; see, for example, [Bailesteanu et al. 2010; Cao and Hamilton 2009; Chau et al. 2011; Chow et al. 2010; Guenther 2002; Liu 2009; Wu and Zheng 2010; Zhang 2006].

In recent years there has been increasing interest in the study of the nonlinear heat-type equations coupled with the Ricci flow. A nice example of a nonlinear heat equation, introduced by L. Ma [2006], is

$$(1-1) \quad \frac{\partial}{\partial t} f = \Delta f - af \ln f - bf,$$

where a and b are real constants. Ma first proved a local gradient estimate for positive solutions to the corresponding elliptic equation

$$(1-2) \quad \Delta f - af \ln f - bf = 0$$

on a complete manifold with a fixed metric. Indeed, F. R. K. Chung and S.-T. Yau [1996] observed that equation (1-2) is linked with the gross logarithmic Sobolev inequality. They also established a logarithmic Harnack inequality for this equation when $a < 0$. Y. Yang [2008] derived local gradient estimates for positive solutions to (1-1) on a complete manifold with a fixed metric; see also [Chen and Chen 2009; Huang and Ma 2010; Wu 2010a; 2010b]. Yang’s result has been generalized by L. Ma [2010a; 2010b], who obtained Hamilton and new Li–Yau type gradient estimates for the nonlinear heat equation (1-1), and also by S.-Y. Hsu [2011], who proved local gradient estimates for the nonlinear heat equation (1-1) under the Ricci flow, similar to the gradient estimates of [Yang 2008] for the fixed metric case.

We remind the reader that equations (1-1) and (1-2) often appear in geometric evolution equations, and are also closely related to the gradient Ricci solitons. See, for example, [Cao and Zhang 2011b; Ma 2006] for nice explanations on this subject.

Very recently, X. Cao and Z. Zhang [2011b] used the argument from [Cao and Hamilton 2009] to prove an interesting differential Harnack inequality for positive solutions to the forward nonlinear heat equation

$$(1-3) \quad \frac{\partial}{\partial t} f = \Delta f - f \ln f + Rf$$

coupled with the Ricci flow equation

$$(1-4) \quad \frac{\partial}{\partial t} g_{ij} = -2R_{ij}$$

on a closed manifold. Here Δ , R and R_{ij} are the Laplacian, scalar curvature and Ricci curvature of the metric $g(t)$ moving under the Ricci flow.

Main results. In this paper, we will be concerned with general time-dependent nonlinear backward heat equations of the type (1-1) with different potentials on closed manifolds under the Ricci flow.

Before studying nonlinear backward heat equations, we first study the nonlinear forward heat equation (1-3) with the metric evolving under the Ricci flow. Suppose $(M, g(t))$, $t \in [0, T)$, is a solution to the ε -Ricci flow ($\varepsilon \geq 0$)

$$(1-5) \quad \frac{\partial}{\partial t} g_{ij} = -\varepsilon Rg_{ij}$$

on a closed surface. Let f be a positive solution to the nonlinear forward heat equation with potential εR , that is,

$$(1-6) \quad \frac{\partial}{\partial t} f = \Delta f - f \ln f + \varepsilon Rf.$$

In this case, we can derive a new differential interpolated Harnack inequality, which is originated with B. Chow [1998].

Theorem 1.1. *Let $(M, g(t))$, $t \in [0, T)$, be a solution to the ε -Ricci flow (1-5) on a closed surface with $R > 0$. Let f be a positive solution to the nonlinear heat equation (1-6), $u = -\ln f$ and $H_\varepsilon = \Delta u - \varepsilon R$. Then, for all time $t \in (0, T)$,*

$$H_\varepsilon \leq \frac{1}{t},$$

that is,

$$\frac{\partial}{\partial t} \ln f - |\nabla \ln f|^2 + \ln f + \frac{1}{t} = \Delta \ln f + \varepsilon R + \frac{1}{t} \geq 0.$$

In Theorem 1.1, if we take $\varepsilon = 0$, we can get the following differential Harnack inequality for the nonlinear heat equation on closed surfaces with a fixed metric:

Corollary 1.2. *If $f : M \times [0, T) \rightarrow \mathbb{R}$, is a positive solution to the nonlinear heat equation*

$$\frac{\partial}{\partial t} f = \Delta f - f \ln f$$

on a closed surface (M, g) with $R > 0$, then, for all time $t \in (0, T)$,

$$\frac{\partial}{\partial t} \ln f - |\nabla \ln f|^2 + \ln f + \frac{1}{t} = \Delta \ln f + \frac{1}{t} \geq 0.$$

If we take $\varepsilon = 1$ in Theorem 1.1, we get:

Corollary 1.3. *Let $(M, g(t))$, $t \in [0, T)$, be a solution to the Ricci flow on a closed surface with $R > 0$. If f is a positive solution to the nonlinear heat equation (1-3), then for all time $t \in (0, T)$,*

$$\frac{\partial}{\partial t} \ln f - |\nabla \ln f|^2 + \ln f + \frac{1}{t} = \Delta \ln f + R + \frac{1}{t} \geq 0.$$

Remark 1.4. X. Cao and Z. Zhang [2011b] have proved a differential Harnack inequality for Equation (1-3) under the Ricci flow on manifolds of any dimension. However, on a closed surface, the result of Corollary 1.3 is better than theirs.

Remark 1.5. Interestingly, Theorem 1.1 is a nonlinear interpolated Harnack inequality which links Corollary 1.2 to Corollary 1.3.

Secondly, we now consider differential Harnack inequalities for positive solutions to the nonlinear backward heat equation with potential $2R$, that is,

$$(1-7) \quad \frac{\partial}{\partial t} f = -\Delta f + f \ln f + 2Rf$$

under the Ricci flow. X. Cao and Z. Zhang [2011b] made nice explanations that the nonlinear forward heat equation (1-3) is closely related to expanding gradient

Ricci solitons. Analogously to the argument of Cao and Zhang, our consideration of the Equation (1-7) is motivated by *shrinking* gradient Ricci solitons proposed in [Hamilton 1993a]. Recall that a shrinking gradient Ricci soliton (M, g) is defined by the form (see [Chow et al. 2006])

$$(1-8) \quad R_{ij} + \nabla_i \nabla_j w = c g_{ij},$$

where w is some Ricci soliton potential and c is a positive constant. Taking the trace of both sides of (1-8) yields

$$(1-9) \quad R + \Delta w = \text{const.}$$

Using the contracted Bianchi identity, we can easily deduce that

$$(1-10) \quad R - 2cw + |\nabla w|^2 = -\text{const.}$$

From (1-9) and (1-10), we get

$$(1-11) \quad 2|\nabla w|^2 = -\Delta w + |\nabla w|^2 + 2cw - 2R.$$

Recall that the Ricci flow solution for a complete gradient Ricci soliton [Chow et al. 2006, Theorem 4.1] is the pullback of g under $\varphi(t)$, up to a scale factor $c(t)$:

$$g(t) = c(t) \cdot \varphi(t)^* g,$$

where $c(t) := -2ct + 1 > 0$ and $\varphi(t)$ is the 1-parameter family of diffeomorphisms generated by

$$\frac{1}{c(t)} \nabla_g w.$$

Then the corresponding Ricci soliton potential $\varphi(t)^* w$ satisfies

$$\frac{\partial}{\partial t} \varphi(t)^* w = |\nabla \varphi(t)^* w|^2.$$

Note that along the Ricci flow, (1-11) becomes

$$2|\nabla \varphi(t)^* w|^2 = -\Delta \varphi(t)^* w + |\nabla \varphi(t)^* w|^2 + \frac{2c}{c(t)} \cdot \varphi(t)^* w - 2R.$$

Hence the evolution equation for the Ricci soliton potential $\varphi(t)^* w$ is

$$2 \frac{\partial \varphi(t)^* w}{\partial t} = -\Delta \varphi(t)^* w + |\nabla \varphi(t)^* w|^2 + \frac{2c}{c(t)} \cdot \varphi(t)^* w - 2R.$$

If we let $\varphi(t)^* w = -\ln \tilde{f}$, this equation becomes

$$(1-12) \quad 2 \frac{\partial \tilde{f}}{\partial t} = -\Delta \tilde{f} + 2R \tilde{f} + \frac{2c}{c(t)} \cdot \tilde{f} \ln \tilde{f}.$$

Notice that (1-7) and (1-12) are closely related and only differ by the time scaling and their last terms.

For the nonlinear backward heat equation (1-7) under the Ricci flow, we have:

Theorem 1.6. *Let $(M, g(t))$, $t \in [0, T]$, be a solution to the Ricci flow on a closed manifold of dimension n . Let f be a positive solution to the nonlinear backward heat equation (1-7), $u = -\ln f$, $\tau = T - t$ and*

$$(1-13) \quad H = 2\Delta u - |\nabla u|^2 + 2R - 2\frac{n}{\tau}.$$

Then, for all time $t \in [0, T]$,

$$H \leq \frac{n}{2}.$$

Remark 1.7. We can easily see that $H \leq n/2$ is equivalent to

$$\frac{|\nabla f|^2}{f^2} - 2 \left(\frac{f_\tau}{f} + \ln f + R \right) \leq 2\frac{n}{\tau} + \frac{n}{2}.$$

In [Yang 2008] (see also [Wu 2010b]), the classical Li–Yau gradient estimate for positive solutions to the nonlinear heat equation (1-1) is

$$\frac{|\nabla f|^2}{f^2} - 2 \left(\frac{f_t}{f} + a \ln f + b \right) \leq 2\frac{n}{t} + na$$

on manifolds with a fixed metric satisfying nonnegative Ricci curvature. Hence our Harnack inequality is similar to the classical Li–Yau gradient estimate for the nonlinear heat equation (1-1).

If we assume instead that our solution to the Ricci flow is defined for $t \in [0, T)$ (where $T < \infty$ is the blow-up time) and is of type I, meaning that

$$(1-14) \quad |\text{Rm}| \leq \frac{d_0}{T-t}$$

for some constant d_0 , then we can show this:

Theorem 1.8. *Let $(M, g(t))$, $t \in [0, T)$ (where $T < \infty$ is the blow-up time) be a solution to the Ricci flow on a closed manifold of dimension n , and assume that g is of type I, that is, it satisfies (1-14), for some constant d_0 . Let f be a positive solution to the nonlinear backward heat equation (1-7), $u = -\ln f$, $\tau = T - t$ and*

$$H = 2\Delta u - |\nabla u|^2 + 2R - d\frac{n}{\tau},$$

where $d = d(d_0, n) \geq 2$ is some constant such that $H(\tau) < 0$ for small τ . Then, for all time $t \in [0, T)$,

$$H \leq \frac{n}{2}.$$

Thirdly, we consider the nonlinear backward heat equation

$$(1-15) \quad \frac{\partial}{\partial t} f = -\Delta f + f \ln f + Rf$$

under the Ricci flow. This equation is very similar to (1-7) and only differs by the last potential. We also find that (1-15) can be regarded as the extension of the linear backward heat equation considered in [Cao 2008, Theorem 1.3] and [Kuang and Zhang 2008, Theorem 2.1]. In fact, we only have the additional term $f \ln f$ in the linear backward heat equation. For this system, we prove:

Theorem 1.9. *Let $(M, g(t)), t \in [0, T]$, be a solution to the Ricci flow on a closed manifold of dimension n with nonnegative scalar curvature. Let f be a positive solution to the nonlinear backward heat equation (1-15), $u = -\ln f$, $\tau = T - t$ and*

$$(1-16) \quad H = 2\Delta u - |\nabla u|^2 + R - 2\frac{n}{\tau}.$$

Then, for all time $t \in [0, T]$,

$$H \leq \frac{n}{4}.$$

By modifying the Harnack quantity of Theorem 1.9, we can deduce the following differential Harnack inequality *without* assuming the nonnegativity of R :

Theorem 1.10. *Let $(M, g(t)), t \in [0, T]$, be a solution to the Ricci flow on a closed manifold of dimension n . Let f be a positive solution to the nonlinear backward heat equation (1-15), $v = -\ln f - \frac{1}{2}n \ln(4\pi\tau)$, $\tau = T - t$, and*

$$P = 2\Delta v - |\nabla v|^2 + R - 3\frac{n}{\tau}.$$

Then, for all time $t \in [T/2, T]$,

$$P \leq \frac{n}{4}.$$

Remark 1.11. Theorems 1.6–1.10 extend to the nonlinear case Theorems 1.1–1.3 and 3.6 of [Cao 2008] and Theorem 2.1 of [Kuang and Zhang 2008].

The proof of all our theorems nearly follows from the arguments of X. Cao [2008], X. Cao and R. Hamilton [2009], X. Cao and Z. Zhang [2011b], and S.-L. Kuang and Qi S. Zhang [Kuang and Zhang 2008], where computations of evolution equations and the maximum principle for parabolic equations are employed. The major differences are that one of our results gives an interpolation Harnack inequality for a nonlinear forward heat equation along the ε -Ricci flow on a closed surface, and the others provide differential Harnack estimates for various *nonlinear backward* heat equations under the Ricci flow.

One interesting feature of this paper is that our differential Harnack inequalities are not only like the Perelman's Harnack inequalities, but also similar to the classical Li–Yau Harnack inequalities for the corresponding nonlinear heat equation (see Remark 1.7 above). Another feature is that our Harnack quantities of nonlinear backward heat equations are nearly the same as those of linear backward heat equations considered by X. Cao [2008], and S.-L. Kuang and Qi S. Zhang [2008]. Due to the fact that Ricci soliton potentials are linked with some nonlinear backward heat equations, we expect that our differential Harnack inequalities will be useful in understanding the Ricci solitons.

The rest of this paper is organized as follows: In Section 2, we will prove a new differential interpolated Harnack inequality on a surface, that is, Theorem 1.1. In Section 3, we firstly derive differential Harnack inequalities for positive solutions to the nonlinear backward heat equation with potential $2R$ under the Ricci flow (Theorems 1.6 and 1.8). Then a classical integral version of the Harnack inequality will be proved (Theorem 3.2). In the latter part of this section, we will establish Harnack inequalities for another nonlinear backward heat equation with potential R under the Ricci flow (Theorem 1.9) as well as its classical Harnack version (Theorem 3.4). By modifying the Harnack quantity of Theorem 1.9, we can prove another differential Harnack inequalities without the nonnegative assumption of scalar curvature (Theorem 1.10). Finally, in Section 4, we will prove gradient estimates for positive and bounded solutions to the nonlinear (including backward) heat equation without potentials under the Ricci flow, that is, Theorems 4.1 and 4.3.

2. Nonlinear heat equation with potentials

In this section, we will prove a differential interpolated Harnack inequality for positive solutions to nonlinear forward heat equations with potentials coupled with the ε -Ricci flow on a closed surface.

Let f be a positive solution to the nonlinear forward heat equation (1-6). By the maximum principle, we conclude that the solution will remain positive along the Ricci flow when scalar curvature is positive. If we let

$$u = -\ln f,$$

then u satisfies the equation

$$\frac{\partial}{\partial t} u = \Delta u - |\nabla u|^2 - \varepsilon R - u.$$

Proof of Theorem 1.1. The proof involves a direct computation and the parabolic maximum principle. Let f and u be defined as above. Under the ε -Ricci flow (1-5)

on a closed surface, we have that

$$\frac{\partial R}{\partial t} = \varepsilon(\Delta R + R^2) \quad \text{and} \quad \frac{\partial}{\partial t}(\Delta) = \varepsilon R \Delta,$$

where the Laplacian Δ is acting on functions. Define the Harnack quantity

$$(2-1) \quad H_\varepsilon = \Delta u - \varepsilon R.$$

Using the evolution equations above, we first compute that

$$\begin{aligned} \frac{\partial}{\partial t} H_\varepsilon &= \Delta \left(\frac{\partial}{\partial t} u \right) + \left(\frac{\partial}{\partial t} \Delta \right) u - \varepsilon \frac{\partial R}{\partial t} \\ &= \Delta(\Delta u - |\nabla u|^2 - \varepsilon R - u) + \varepsilon R \Delta u - \varepsilon \frac{\partial R}{\partial t} \\ &= \Delta H_\varepsilon - \Delta |\nabla u|^2 - \Delta u + \varepsilon R H_\varepsilon + \varepsilon^2 R^2 - \varepsilon \frac{\partial R}{\partial t} \end{aligned}$$

Since

$$\Delta |\nabla u|^2 = 2|\nabla \nabla u|^2 + 2\nabla \Delta u \cdot \nabla u + R|\nabla u|^2$$

on a two-dimensional surface, we then have

$$\begin{aligned} \frac{\partial}{\partial t} H_\varepsilon &= \Delta H_\varepsilon - 2|\nabla \nabla u|^2 - 2\nabla \Delta u \cdot \nabla u - R|\nabla u|^2 + \varepsilon R H_\varepsilon + \varepsilon^2 R^2 - \varepsilon \frac{\partial R}{\partial t} - \Delta u \\ &= \Delta H_\varepsilon - 2|\nabla \nabla u|^2 - 2\nabla H_\varepsilon \cdot \nabla u \\ &\quad - 2\varepsilon \nabla R \cdot \nabla u - R|\nabla u|^2 + \varepsilon R H_\varepsilon + \varepsilon^2 R^2 - \varepsilon \frac{\partial R}{\partial t} - \Delta u \\ &= \Delta H_\varepsilon - 2 \left| \nabla_i \nabla_j u - \frac{\varepsilon}{2} R g_{ij} \right|^2 - 2\varepsilon R \Delta u - 2\nabla H_\varepsilon \cdot \nabla u \\ &\quad - 2\varepsilon \nabla R \cdot \nabla u - R|\nabla u|^2 + \varepsilon R H_\varepsilon + 2\varepsilon^2 R^2 - \varepsilon \frac{\partial R}{\partial t} - \Delta u. \end{aligned}$$

Since $\Delta u = H_\varepsilon + \varepsilon R$ by (2-1), these equalities become

$$\begin{aligned} \frac{\partial}{\partial t} H_\varepsilon &= \Delta H_\varepsilon - 2 \left| \nabla_i \nabla_j u - \frac{\varepsilon}{2} R g_{ij} \right|^2 - \varepsilon R H_\varepsilon - 2\nabla H_\varepsilon \cdot \nabla u \\ &\quad - 2\varepsilon \nabla R \cdot \nabla u - R|\nabla u|^2 - \varepsilon \frac{\partial R}{\partial t} - \Delta u. \end{aligned}$$

Rearranging terms yields

$$\begin{aligned} (2-2) \quad \frac{\partial}{\partial t} H_\varepsilon &= \Delta H_\varepsilon - 2 \left| \nabla_i \nabla_j u - \frac{\varepsilon}{2} R g_{ij} \right|^2 - 2\nabla H_\varepsilon \cdot \nabla u - \varepsilon R H_\varepsilon \\ &\quad - R |\nabla u + \varepsilon \nabla \ln R|^2 - \varepsilon R \left(\frac{\partial \ln R}{\partial t} - \varepsilon |\nabla \ln R|^2 \right) - \Delta u \\ &\leq \Delta H_\varepsilon - H_\varepsilon^2 - 2\nabla H_\varepsilon \cdot \nabla u - (\varepsilon R + 1) H_\varepsilon + \frac{\varepsilon}{t} R - \varepsilon R. \end{aligned}$$

The reason for this last inequality is that the trace Harnack inequality for the ε -Ricci flow on a closed surface proved in [Chow 1998] (see also [Wu and Zheng

2010, Lemma 2.1]) states that

$$\frac{\partial \ln R}{\partial t} - \varepsilon |\nabla \ln R|^2 = \varepsilon(\Delta \ln R + R) \geq -\frac{1}{t},$$

since $g(t)$ has positive scalar curvature. Besides this, we also used (2-1) and the elementary inequality

$$\left| \nabla_i \nabla_j u - \frac{\varepsilon}{2} R g_{ij} \right|^2 \geq \frac{1}{2} (\Delta u - \varepsilon R)^2 = \frac{1}{2} H_\varepsilon^2.$$

Adding $-1/t$ to H_ε in (2-2) yields

$$(2-3) \quad \begin{aligned} \frac{\partial}{\partial t} \left(H_\varepsilon - \frac{1}{t} \right) &\leq \Delta \left(H_\varepsilon - \frac{1}{t} \right) - 2 \nabla \left(H_\varepsilon - \frac{1}{t} \right) \cdot \nabla u \\ &\quad - \left(H_\varepsilon + \frac{1}{t} \right) \left(H_\varepsilon - \frac{1}{t} \right) - (\varepsilon R + 1) \left(H_\varepsilon - \frac{1}{t} \right) - \frac{1}{t} - \varepsilon R. \end{aligned}$$

Clearly, for t small enough we have $H_\varepsilon - 1/t < 0$. Since $R > 0$, applying the maximum principle to the evolution formula (2-3) we conclude that $H_\varepsilon - 1/t \leq 0$ for all time t , and the proof of this theorem is completed. \square

We remark that Theorem 1.1 can be regarded as a nonlinear version of an interpolated Harnack inequality proved by B. Chow:

Theorem 2.1 [Chow 1998]. *Let $(M, g(t))$ be a solution to the ε -Ricci flow (1-5) on a closed surface with $R > 0$. If f is a positive solution to*

$$\frac{\partial}{\partial t} f = \Delta f + \varepsilon R f,$$

then

$$\frac{\partial}{\partial t} \ln f - |\nabla \ln f|^2 + \frac{1}{t} = \Delta \ln f + \varepsilon R + \frac{1}{t} \geq 0.$$

3. Nonlinear backward heat equation with potentials

We next study several differential Harnack inequalities for positive solutions to the nonlinear backward heat equation under the Ricci flow, proving Theorems 1.6, 1.8, 1.9, and 1.10 from the Introduction. The first two of these theorems deal with the case where the potential equals $2R$, and the last two with the potential R . The proofs are largely based on the maximum principle.

Potential 2R. Theorems 1.6 and 1.8 deal with differential Harnack inequalities for positive solutions to the equation

$$\frac{\partial}{\partial t} f = -\Delta f + f \ln f + 2Rf$$

under the Ricci flow. We follow the trick used to prove Theorem 1.1 in [Cao and Zhang 2011b] to simplify a tedious calculation of the evolution equations. Also,

the evolution equation of u in this case is very similar to what is considered in [Cao 2008]. So we can borrow Cao's computation for the very general setting there to simplify our calculation. The only difference is that we have extra terms coming from the time derivative $\partial u/\partial \tau$.

Proof of Theorem 1.6. As before, it is easy to compute that u satisfies

$$(3-1) \quad \frac{\partial}{\partial \tau} u = \Delta u - |\nabla u|^2 + 2R - u.$$

Recall from (1-13) that $H = 2\Delta u - |\nabla u|^2 + 2R - 2n/\tau$. Adapting [Cao 2008, (2.4)] and using (3-1) as well as the elementary inequality

$$\left| \nabla_i \nabla_j u - R_{ij} - \frac{1}{\tau} g_{ij} \right|^2 \geq \frac{1}{n} \left(\Delta u - R - \frac{n}{\tau} \right)^2,$$

we can write

$$\begin{aligned} \frac{\partial}{\partial \tau} H &= \Delta H - 2\nabla H \cdot \nabla u - \frac{2}{\tau} H - \frac{2}{\tau} |\nabla u|^2 - 2|\text{Rc}|^2 - 2 \left| \nabla_i \nabla_j u + R_{ij} - \frac{1}{\tau} g_{ij} \right|^2 \\ &\quad - 2(\Delta u - |\nabla u|^2) \\ &\leq \Delta H - 2\nabla H \cdot \nabla u - \frac{2}{\tau} H - \frac{2}{\tau} |\nabla u|^2 - \frac{2}{n} R^2 - \frac{2}{n} \left(\Delta u + R - \frac{n}{\tau} \right)^2 \\ &\quad - 2(\Delta u - |\nabla u|^2), \end{aligned}$$

By the definition of H , we have

$$-2(\Delta u - |\nabla u|^2) = -2H + 2 \left(\Delta u + R - \frac{n}{\tau} \right) + 2R - \frac{2n}{\tau}.$$

Plugging this into the preceding inequality yields

$$\begin{aligned} \frac{\partial}{\partial \tau} H &\leq \Delta H - 2\nabla H \cdot \nabla u - \left(\frac{2}{\tau} + 2 \right) H - \frac{2}{\tau} |\nabla u|^2 - \frac{2}{n} R^2 \\ &\quad - \frac{2}{n} \left(\Delta u + R - \frac{n}{\tau} - \frac{n}{2} \right)^2 + \frac{n}{2} + 2R - \frac{2n}{\tau} \\ &= \Delta H - 2\nabla H \cdot \nabla u - \left(\frac{2}{\tau} + 2 \right) H - \frac{2}{\tau} |\nabla u|^2 \\ &\quad - \frac{2}{n} \left(\Delta u + R - \frac{n}{\tau} - \frac{n}{2} \right)^2 - \frac{2}{n} \left(R - \frac{n}{2} \right)^2 - \frac{2n}{\tau} + n. \end{aligned}$$

Adding $-n/2$ to H , we then get

$$(3-2) \quad \frac{\partial}{\partial \tau} \left(H - \frac{n}{2} \right) \leq \Delta \left(H - \frac{n}{2} \right) - 2\nabla \left(H - \frac{n}{2} \right) \cdot \nabla u - \left(\frac{2}{\tau} + 2 \right) \left(H - \frac{n}{2} \right) \\ - \frac{2}{\tau} |\nabla u|^2 - \frac{2}{n} \left(\Delta u + R - \frac{n}{\tau} - \frac{n}{2} \right)^2 - \frac{2}{n} \left(R - \frac{n}{2} \right)^2 - \frac{3n}{\tau}.$$

If τ is small enough, $H - n/2 < 0$. Then applying the maximum principle to the evolution equation (3-2) yields $H - n/2 \leq 0$ for all τ , hence for all $t \in [0, T)$. \square

An easy modification of the preceding proof, using (1-14) to ensure that we can apply the maximum principle as $\tau \rightarrow 0$, verifies Theorem 1.8. We omit the details.

Remark 3.1. Theorem 1.6 is also true on a complete noncompact Riemannian manifolds, as long as we can apply the maximum principle.

From Theorem 1.6, we can derive a classical Harnack inequality by integrating along a space-time path.

Theorem 3.2. *Let $(M, g(t))$, $t \in [0, T]$, be a solution to the Ricci flow on a closed manifold of dimension n . Let f be a positive solution to the nonlinear backward heat equation (1-7). Assume that (x_1, t_1) and (x_2, t_2) , $0 \leq t_1 < t_2 < T$, are two points in $M \times [0, T)$. Then we have*

$$e^{t_2} \ln f(x_2, t_2) - e^{t_1} \ln f(x_1, t_1) \leq \frac{1}{2} \int_{t_1}^{t_2} e^{T-t} \left(|\dot{\gamma}|^2 + 2R + \frac{n}{2} + \frac{2n}{T-t} \right) dt,$$

where γ is any space-time path joining (x_1, t_1) and (x_2, t_2) .

Proof. This is similar to Theorem 2.3 in [Cao 2008]; we include the proof for completeness. Consider the solutions to

$$\frac{\partial}{\partial \tau} u = \Delta u - |\nabla u|^2 + 2R - u.$$

Combining this with

$$H - \frac{n}{2} = 2\Delta u - |\nabla u|^2 + 2R - 2\frac{n}{\tau} - \frac{n}{2} \leq 0,$$

we have

$$2\frac{\partial}{\partial \tau} u + |\nabla u|^2 - 2R - 2\frac{n}{\tau} + 2u - \frac{n}{2} \leq 0.$$

If $\gamma(x, t)$ is a space-time path joining (x_2, τ_2) and (x_1, τ_1) , with $\tau_1 > \tau_2 > 0$, we have along γ

$$\begin{aligned} \frac{du}{d\tau} &= \frac{\partial u}{\partial \tau} + \nabla u \cdot \gamma \leq -\frac{1}{2}|\nabla u|^2 + R + \frac{n}{\tau} - u + \frac{n}{4} + \nabla u \cdot \gamma \\ &\leq \frac{1}{2} \left(|\dot{\gamma}|^2 + 2R + \frac{n}{2} \right) + \frac{n}{\tau} - u, \end{aligned}$$

where in the last step we used the inequality $-\frac{1}{2}|\nabla u|^2 + \nabla u \cdot \gamma - \frac{1}{2}|\dot{\gamma}|^2 \leq 0$. Rearranging terms yields

$$\frac{d}{d\tau} (e^\tau \cdot u) \leq \frac{e^\tau}{2} \left(|\dot{\gamma}|^2 + 2R + \frac{n}{2} + \frac{2n}{\tau} \right).$$

Integrating this inequality we obtain

$$e^{\tau_1} \cdot u(x_1, \tau_1) - e^{\tau_2} \cdot u(x_2, \tau_2) \leq \frac{1}{2} \int_{\tau_2}^{\tau_1} e^{\tau} \left(|\dot{\gamma}|^2 + 2R + \frac{n}{2} + \frac{2n}{\tau} \right) d\tau,$$

which can be rewritten as

$$e^{t_1} \cdot u(x_1, t_1) - e^{t_2} \cdot u(x_2, t_2) \leq \frac{1}{2} \int_{t_1}^{t_2} e^{T-t} \left(|\dot{\gamma}|^2 + 2R + \frac{n}{2} + \frac{2n}{T-t} \right) dt.$$

Note that $u = -\ln f$. Hence the desired classical Harnack inequality follows. \square

Potential R . We now turn to the equation with potential R :

$$\frac{\partial}{\partial t} f = -\Delta f + f \ln f + Rf.$$

Here we need to assume that the initial metric $g(0)$ has nonnegative scalar curvature. It is well known that this property is preserved by the Ricci flow.

Proof of Theorem 1.9. This time u satisfies

$$\frac{\partial}{\partial \tau} u = \Delta u - |\nabla u|^2 + R - u.$$

Adapting [Cao 2008, (3.2)], we can write

$$(3-3) \quad \begin{aligned} \frac{\partial}{\partial \tau} H &= \Delta H - 2\nabla H \cdot \nabla u - \frac{2}{\tau} H - \frac{2}{\tau} |\nabla u|^2 - 2\frac{R}{\tau} \\ &\quad - 2 \left| \nabla_i \nabla_j u + R_{ij} - \frac{1}{\tau} g_{ij} \right|^2 - 2(\Delta u - |\nabla u|^2). \end{aligned}$$

Since H is now given by (1-16), we have

$$-2(\Delta u - |\nabla u|^2) = -2H + 2 \left(\Delta u + R - \frac{n}{\tau} \right) - \frac{2n}{\tau}.$$

Plugging this into (3-3), we obtain

$$\begin{aligned} \frac{\partial}{\partial \tau} H &\leq \Delta H - 2\nabla H \cdot \nabla u - \left(\frac{2}{\tau} + 2 \right) H - \frac{2}{\tau} |\nabla u|^2 - 2\frac{R}{\tau} \\ &\quad - \frac{2}{n} \left(\Delta u + R - \frac{n}{\tau} \right)^2 + 2 \left(\Delta u + R - \frac{n}{\tau} \right) - \frac{2n}{\tau} \\ &= \Delta H - 2\nabla H \cdot \nabla u - \left(\frac{2}{\tau} + 2 \right) H - \frac{2}{\tau} |\nabla u|^2 - 2\frac{R}{\tau} \\ &\quad - \frac{2}{n} \left(\Delta u + R - \frac{n}{\tau} - \frac{n}{2} \right)^2 - \frac{2n}{\tau} + \frac{n}{2}. \end{aligned}$$

Adding $-n/4$ to H yields

$$(3-4) \quad \begin{aligned} \frac{\partial}{\partial \tau} \left(H - \frac{n}{4} \right) &\leq \Delta \left(H - \frac{n}{4} \right) - 2\nabla \left(H - \frac{n}{4} \right) \cdot \nabla u - \left(\frac{2}{\tau} + 2 \right) \left(H - \frac{n}{4} \right) \\ &\quad - \frac{2}{\tau} |\nabla u|^2 - 2\frac{R}{\tau} - \frac{2}{n} \left(\Delta u + R - \frac{n}{\tau} - \frac{n}{2} \right)^2 - \frac{5n}{2\tau}. \end{aligned}$$

Since $R \geq 0$, it is easy to see that $H - n/4 < 0$ for τ small enough. Applying the maximum principle to the evolution formula (3-4), we have $H - n/4 \leq 0$ for all τ , hence for all t . This finishes the proof of Theorem 1.9. \square

We easily derive counterparts to Theorem 1.8 and Theorem 3.2:

Theorem 3.3. *Let $(M, g(t))$, $t \in [0, T)$ (where $T < \infty$ is the blow-up time) be a solution to the Ricci flow on a closed manifold of dimension n with nonnegative scalar curvature, and assume that g is of type I, that is, it satisfies (1-14), for some constant d_0 . Let f be a positive solution to the nonlinear backward heat equation (1-15), $u = -\ln f$, $\tau = T - t$ and*

$$H = 2\Delta u - |\nabla u|^2 + R - d\frac{n}{\tau},$$

where $d = d(d_0, n) \geq 1$ is some constant such that $H(\tau) < 0$ for small τ . Then, for all time $t \in [0, T)$,

$$H \leq \frac{n}{4}.$$

Theorem 3.4. *Let $(M, g(t))$, $t \in [0, T]$, be a solution to the Ricci flow on a closed manifold of dimension n with nonnegative scalar curvature. Let f be a positive solution to the nonlinear backward heat equation (1-15). Assume that (x_1, t_1) and (x_2, t_2) , with $0 \leq t_1 < t_2 < T$, are two points in $M \times [0, T)$. Then*

$$e^{t_2} \ln f(x_2, t_2) - e^{t_1} \ln f(x_1, t_1) \leq \frac{1}{2} \int_{t_1}^{t_2} e^{T-t} \left(|\dot{\gamma}|^2 + R + \frac{n}{4} + \frac{2n}{T-t} \right) dt,$$

where γ is any space-time path joining (x_1, t_1) and (x_2, t_2) .

In the rest of this section, we will finish the proof of Theorem 1.10. The interesting feature of Theorem 1.10 is that the differential Harnack inequalities hold without any assumption on the scalar curvature R .

Proof of Theorem 1.10. We first compute that v satisfies

$$(3-5) \quad \frac{\partial}{\partial \tau} v = \Delta v - |\nabla v|^2 + R - \frac{n}{2\tau} - \left(v + \frac{n}{2} \ln(4\pi\tau) \right).$$

If we let

$$\tilde{P} := 2\Delta v - |\nabla v|^2 + R - 2\frac{n}{\tau},$$

then by adapting [Cao 2008, (3.7)], we have

$$\begin{aligned} \frac{\partial}{\partial \tau} \tilde{P} = & \Delta \tilde{P} - 2\nabla \tilde{P} \cdot \nabla v - \frac{2}{\tau} \tilde{P} - \frac{2}{\tau} |\nabla v|^2 - 2\frac{R}{\tau} \\ & - 2 \left| \nabla_i \nabla_j v + R_{ij} - \frac{1}{\tau} g_{ij} \right|^2 - 2(\Delta v - |\nabla v|^2). \end{aligned}$$

Since $P = \tilde{P} - n/\tau$, we have

$$(3-6) \quad \frac{\partial}{\partial \tau} P = \Delta P - 2\nabla P \cdot \nabla v - \frac{2}{\tau} P - \frac{2}{\tau} |\nabla v|^2 - 2\frac{R}{\tau} - \frac{n}{\tau^2} - 2 \left| \nabla_i \nabla_j v + R_{ij} - \frac{1}{\tau} g_{ij} \right|^2 - 2(\Delta v - |\nabla v|^2).$$

According to the definition of P , we have

$$-2a(\Delta v - |\nabla v|^2) = -2P + 2 \left(\Delta v + R - \frac{n}{\tau} \right) - \frac{4n}{\tau}.$$

Substituting this into (3-6), we get

$$(3-7) \quad \begin{aligned} \frac{\partial}{\partial \tau} P &\leq \Delta P - 2\nabla P \cdot \nabla v - \left(\frac{2}{\tau} + 2 \right) P - \frac{2}{\tau} |\nabla v|^2 - 2\frac{R}{\tau} - \frac{n}{\tau^2} \\ &\quad - \frac{2}{n} \left(\Delta v + R - \frac{n}{\tau} \right)^2 + 2 \left(\Delta v + R - \frac{n}{\tau} \right) - \frac{4n}{\tau} \\ &= \Delta P - 2\nabla P \cdot \nabla v - \left(\frac{2}{\tau} + 2 \right) P - \frac{2}{\tau} |\nabla v|^2 - \frac{2}{\tau} \left(R + \frac{n}{2\tau} \right) \\ &\quad - \frac{2}{n} \left(\Delta v + R - \frac{n}{\tau} - \frac{n}{2} \right)^2 - \frac{4n}{\tau} + \frac{n}{2}. \end{aligned}$$

Note that the evolution of scalar curvature under the Ricci flow is

$$\frac{\partial R}{\partial t} = \Delta R + 2|\text{Rc}|^2 \geq \Delta R + \frac{2}{n} R^2.$$

Applying the maximum principle to this inequality yields $R \geq -n/(2t)$. Since $t \geq T/2$, we have $1/t \leq 1/\tau$. Hence

$$R \geq -\frac{n}{2t} \geq -\frac{n}{2\tau},$$

that is,

$$R + \frac{n}{2\tau} \geq 0.$$

Combining this with (3-7), we have

$$\frac{\partial}{\partial \tau} P \leq \Delta P - 2\nabla P \cdot \nabla v - \left(\frac{2}{\tau} + 2 \right) P - \frac{4n}{\tau} + \frac{n}{2}.$$

Adding $-n/4$ to P , we get

$$(3-8) \quad \frac{\partial}{\partial \tau} \left(P - \frac{n}{4} \right) \leq \Delta \left(P - \frac{n}{4} \right) - 2\nabla \left(P - \frac{n}{4} \right) \cdot \nabla v - \left(\frac{2}{\tau} + 2 \right) \left(P - \frac{n}{4} \right) - \frac{9n}{2\tau}.$$

It is easy to see that $P - n/4 < 0$ for τ small enough. Applying the maximum principle to the evolution formula (3-8) yields

$$P - \frac{n}{4} \leq 0$$

for all time $t \geq T/2$. Hence the theorem is proved. \square

Remark 3.5. Motivated by Theorems 3.3 and 3.4, we can prove similar theorems by the standard argument from Theorem 1.10. We omit them in the interests of brevity.

4. Gradient estimates for nonlinear (backward) heat equations

In this section, on one hand we consider the positive solution $f(x, t) < 1$ to the nonlinear heat equation without any potential

$$(4-1) \quad \frac{\partial}{\partial t} f = \Delta f - f \ln f,$$

with the metric evolved by the Ricci flow (1-4) on a closed manifold M . This equation has been considered by S.-Y. Hsu [2011] and L. Ma [2010a]. If we let $u = -\ln f$, then

$$(4-2) \quad \frac{\partial}{\partial t} u = \Delta u - |\nabla u|^2 - u$$

and $u > 0$. Note that $0 < f < 1$ is preserved as time t evolves. In fact the initial assumption says that

$$-\ln \sup_M f(x, 0) \leq u(x, 0) \leq -\ln \inf_M f(x, 0).$$

Applying the maximum principle to (4-2), we have

$$-e^{-t} \ln \sup_M f(x, 0) \leq u(x, t) \leq -e^{-t} \ln \inf_M f(x, 0)$$

and hence

$$0 < u(x, t) \leq -\ln \inf_M f(x, 0)$$

for all $x \in M$ and $t \in [0, T)$. Since $u = -\ln f$, this implies

$$0 < \inf_M f(x, 0) \leq f(x, t) < 1$$

for all $x \in M$ and $t \in [0, T)$.

Following the arguments of [Cao and Hamilton 2009], we let

$$H = |\nabla u|^2 - \frac{u}{t}.$$

Comparing with the equation (5.3) in the same reference, we have

$$(4-3) \quad \begin{aligned} \frac{\partial}{\partial t} H &= \Delta H - 2\nabla H \cdot \nabla u - \frac{1}{t} H - 2|\nabla \nabla u|^2 - 2|\nabla u|^2 + \frac{u}{t} \\ &= \Delta H - 2\nabla H \cdot \nabla u - \left(\frac{1}{t} + 1\right) H - 2|\nabla \nabla u|^2 - |\nabla u|^2. \end{aligned}$$

Notice that if t small enough, then $H < 0$. Then applying the maximum principle to (4-3), we obtain:

Theorem 4.1. *Let $(M, g(t)), t \in [0, T)$, be a solution to the Ricci flow on a closed manifold. Let $f < 1$ be a positive solution to the nonlinear heat equation (4-1), $u = -\ln f$ and*

$$H = |\nabla u|^2 - \frac{u}{t}.$$

Then, for all time $t \in (0, T)$,

$$H \leq 0.$$

Remark 4.2. Theorem 4.1 can be regarded as a nonlinear version of [Cao and Hamilton 2009, Theorem 5.1]. Recently, L. Ma [2010a, Theorem 3] has proved the same estimate as in Theorem 4.1 on a closed manifold with nonnegative Ricci curvature under a static metric. However, in our case, we do not need any curvature assumption.

On the other hand, we can also consider the positive solution $f(x, t) < 1$ to the nonlinear backward heat equation without any potential

$$(4-4) \quad \frac{\partial}{\partial t} f = -\Delta f + f \ln f,$$

with the metric evolved by the Ricci flow (1-4). Let $u = -\ln f$. Then we have

$$\frac{\partial}{\partial \tau} u = \Delta u - |\nabla u|^2 - u$$

and $u > 0$. Using the maximum principle, one can see that $0 < f < 1$ is also preserved under the Ricci flow. In fact from the initial assumption

$$0 < \inf_M f(x, T) \leq f(x, T) \leq \sup_M f(x, T) < 1,$$

one can also show that

$$0 < \inf_M f(x, T) \leq f(x, \tau) < 1$$

for all $x \in M$ and $\tau \in (0, T]$ in the same way as the above arguments.

Following the arguments of [Cao 2008], let

$$H = |\nabla u|^2 - \frac{u}{\tau}.$$

Comparing with the equation (5.3) in [Cao 2008], we have

$$(4-5) \quad \begin{aligned} \frac{\partial}{\partial \tau} H &= \Delta H - 2\nabla H \cdot \nabla u - \frac{1}{\tau} H - 2|\nabla \nabla u|^2 - 4R_{ij}u_i u_j - 2|\nabla u|^2 + \frac{u}{\tau} \\ &= \Delta H - 2\nabla H \cdot \nabla u - \left(\frac{1}{\tau} + 1\right) H - 2|\nabla \nabla u|^2 - 4R_{ij}u_i u_j - |\nabla u|^2. \end{aligned}$$

If we assume $R_{ij}(g(t)) \geq -K$, where $0 \leq K \leq \frac{1}{4}$, then

$$-4R_{ij}u_i u_j - |\nabla u|^2 \leq (4K - 1)|\nabla u|^2 \leq 0.$$

Hence if τ small enough, then $H < 0$. Then applying the maximum principle to (4-5), we have a nonlinear version of [Cao 2008, Theorem 5.1].

Theorem 4.3. *Let $(M, g(t))$, $t \in [0, T]$, be a solution to the Ricci flow on a closed manifold with the Ricci curvature satisfying $R_{ij}(g(t)) \geq -K$, where $0 \leq K \leq \frac{1}{4}$. Let $f < 1$ be a positive solution to the nonlinear backward heat equation (4-4), $u = -\ln f$, $\tau = T - t$ and*

$$H = |\nabla u|^2 - \frac{u}{\tau}.$$

Then, for all time $t \in [0, T)$,

$$H \leq 0.$$

Acknowledgments

The author would like to express his gratitude to the referee for careful readings and many valuable suggestions.

References

- [Andrews 1994] B. Andrews, “Harnack inequalities for evolving hypersurfaces”, *Math. Z.* **217**:2 (1994), 179–197. MR 95j:58178 Zbl 0807.53044
- [Aronson and B enilan 1979] D. G. Aronson and P. B enilan, “R egularit e des solutions de l’ equation des milieux poreux dans \mathbb{R}^N ”, *C. R. Acad. Sci. Paris S er. A-B* **288**:2 (1979), A103–A105. MR 82i:35090
- [Bailesteanu et al. 2010] M. Bailesteanu, X. Cao, and A. Pulemotov, “Gradient estimates for the heat equation under the Ricci flow”, *J. Funct. Anal.* **258**:10 (2010), 3517–3542. MR 2011b:53153 Zbl 1193.53139
- [Cao 1992] H. D. Cao, “On Harnack’s inequalities for the K ahler–Ricci flow”, *Invent. Math.* **109**:2 (1992), 247–263. MR 93f:58227 Zbl 0779.53043
- [Cao 2008] X. Cao, “Differential Harnack estimates for backward heat equations with potentials under the Ricci flow”, *J. Funct. Anal.* **255**:4 (2008), 1024–1038. MR 2009e:35121 Zbl 1146.58014
- [Cao and Hamilton 2009] X. Cao and R. S. Hamilton, “Differential Harnack estimates for time-dependent heat equations with potentials”, *Geom. Funct. Anal.* **19**:4 (2009), 989–1000. MR 2010j:53124 Zbl 1183.53059
- [Cao and Ni 2005] H.-D. Cao and L. Ni, “Matrix Li–Yau–Hamilton estimates for the heat equation on K ahler manifolds”, *Math. Ann.* **331**:4 (2005), 795–807. MR 2006k:53113 Zbl 1083.58024
- [Cao and Zhang 2011a] X. Cao and Q. S. Zhang, “The conjugate heat equation and ancient solutions of the Ricci flow”, *Adv. Math.* **228**:5 (2011), 2891–2919. MR 2838064 Zbl 05969510
- [Cao and Zhang 2011b] X. Cao and Z. Zhang, “Differential Harnack estimates for parabolic equations”, pp. 87–98 in *Complex and differential geometry* (Hannover, 2009), edited by W. Ebeling et al., Springer Proceedings in Mathematics **8**, Springer, Berlin, 2011. Zbl 1228.53078

- [Chau et al. 2011] A. Chau, L.-F. Tam, and C. Yu, “Pseudolocality for the Ricci flow and applications”, *Canad. J. Math.* **63**:1 (2011), 55–85. MR 2012g:53133 Zbl 1214.53053
- [Chen and Chen 2009] L. Chen and W. Chen, “Gradient estimates for a nonlinear parabolic equation on complete non-compact Riemannian manifolds”, *Ann. Global Anal. Geom.* **35**:4 (2009), 397–404. MR 2010k:35501 Zbl 1177.35040
- [Cheng 2006] H.-B. Cheng, “A new Li–Yau–Hamilton estimate for the Ricci flow”, *Comm. Anal. Geom.* **14**:3 (2006), 551–564. MR 2008h:53120 Zbl 1116.53039
- [Chow 1991a] B. Chow, “The Ricci flow on the 2-sphere”, *J. Differential Geom.* **33**:2 (1991), 325–334. MR 92d:53036 Zbl 0734.53033
- [Chow 1991b] B. Chow, “On Harnack’s inequality and entropy for the Gaussian curvature flow”, *Comm. Pure Appl. Math.* **44**:4 (1991), 469–483. MR 93e:58032 Zbl 0734.53035
- [Chow 1992] B. Chow, “The Yamabe flow on locally conformally flat manifolds with positive Ricci curvature”, *Comm. Pure Appl. Math.* **45**:8 (1992), 1003–1014. MR 93d:53045 Zbl 0785.53027
- [Chow 1998] B. Chow, “Interpolating between Li–Yau’s and Hamilton’s Harnack inequalities on a surface”, *J. Partial Differ. Equ.* **11**:2 (1998), 137–140. MR 99h:58182 Zbl 0943.58017
- [Chow and Hamilton 1997] B. Chow and R. S. Hamilton, “Constrained and linear Harnack inequalities for parabolic equations”, *Invent. Math.* **129**:2 (1997), 213–238. MR 98i:53051 Zbl 0903.58054
- [Chow and Knopf 2002] B. Chow and D. Knopf, “New Li–Yau–Hamilton inequalities for the Ricci flow via the space-time approach”, *J. Differential Geom.* **60**:1 (2002), 1–54. MR 2003g:53116 Zbl 1048.53026
- [Chow et al. 2006] B. Chow, P. Lu, and L. Ni, *Hamilton’s Ricci flow*, Graduate Studies in Mathematics **77**, Amer. Math. Soc., Providence, RI, 2006. MR 2008a:53068 Zbl 1118.53001
- [Chow et al. 2010] B. Chow, S.-C. Chu, D. Glickenstein, C. Guenther, J. Isenberg, T. Ivey, D. Knopf, P. Lu, F. Luo, and L. Ni, *The Ricci flow: techniques and applications, III: Geometric-analytic aspects*, Mathematical Surveys and Monographs **163**, Amer. Math. Soc., Providence, RI, 2010. MR 2011g:53142 Zbl 1216.53057
- [Chung and Yau 1996] F. R. K. Chung and S.-T. Yau, “Logarithmic Harnack inequalities”, *Math. Res. Lett.* **3**:6 (1996), 793–812. MR 97k:58182 Zbl 0880.58026
- [Guenther 2002] C. M. Guenther, “The fundamental solution on manifolds with time-dependent metrics”, *J. Geom. Anal.* **12**:3 (2002), 425–436. MR 2003a:58034 Zbl 1029.58018
- [Hamilton 1988] R. S. Hamilton, “The Ricci flow on surfaces”, pp. 237–262 in *Mathematics and general relativity* (Santa Cruz, CA, 1986), edited by J. A. Isenberg, Contemp. Math. **71**, Amer. Math. Soc., Providence, RI, 1988. MR 89i:53029 Zbl 0663.53031
- [Hamilton 1993a] R. S. Hamilton, “The Harnack estimate for the Ricci flow”, *J. Differential Geom.* **37**:1 (1993), 225–243. MR 93k:58052 Zbl 0804.53023
- [Hamilton 1993b] R. S. Hamilton, “A matrix Harnack estimate for the heat equation”, *Comm. Anal. Geom.* **1**:1 (1993), 113–126. MR 94g:58215 Zbl 0799.53048
- [Hamilton 1995] R. S. Hamilton, “Harnack estimate for the mean curvature flow”, *J. Differential Geom.* **41**:1 (1995), 215–226. MR 95m:53055 Zbl 0827.53006
- [Hsu 2011] S.-Y. Hsu, “Gradient estimates for a nonlinear parabolic equation under Ricci flow”, *Differential Integral Equations* **24**:7-8 (2011), 645–652. MR 2830313 Zbl 06033866
- [Huang and Ma 2010] G. Huang and B. Ma, “Gradient estimates for a nonlinear parabolic equation on Riemannian manifolds”, *Arch. Math. (Basel)* **94**:3 (2010), 265–275. MR 2011b:58054 Zbl 1194.58020

- [Kuang and Zhang 2008] S. Kuang and Q. S. Zhang, “A gradient estimate for all positive solutions of the conjugate heat equation under Ricci flow”, *J. Funct. Anal.* **255**:4 (2008), 1008–1023. MR 2009h:53150 Zbl 1146.58017
- [Li and Xu 2011] J. Li and X. Xu, “Differential Harnack inequalities on Riemannian manifolds, I: Linear heat equation”, *Adv. Math.* **226**:5 (2011), 4456–4491. MR 2770456 Zbl 1226.58009
- [Li and Yau 1986] P. Li and S.-T. Yau, “On the parabolic kernel of the Schrödinger operator”, *Acta Math.* **156**:3-4 (1986), 153–201. MR 87f:58156 Zbl 0611.58045
- [Liu 2009] S. Liu, “Gradient estimates for solutions of the heat equation under Ricci flow”, *Pacific J. Math.* **243**:1 (2009), 165–180. MR 2010g:53122 Zbl 1180.58017
- [Ma 2006] L. Ma, “Gradient estimates for a simple elliptic equation on complete non-compact Riemannian manifolds”, *J. Funct. Anal.* **241**:1 (2006), 374–382. MR 2007e:53034 Zbl 1112.58023
- [Ma 2010a] L. Ma, “Hamilton type estimates for heat equations on manifolds”, preprint, 2010. arXiv 1009.0603
- [Ma 2010b] L. Ma, “Gradient estimates for a simple nonlinear heat equation on manifolds”, preprint, 2010. arXiv 1009.0604
- [Ni 2007] L. Ni, “A matrix Li–Yau–Hamilton estimate for Kähler–Ricci flow”, *J. Differential Geom.* **75**:2 (2007), 303–358. MR 2008d:53093 Zbl 1120.53023
- [Perelman 2002] G. Perelman, “The entropy formula for the Ricci flow and its geometric applications”, preprint, 2002. Zbl 1130.53001 arXiv math.DG/0211159v1
- [Wu 2010a] J.-Y. Wu, “Gradient estimates for a nonlinear diffusion equation on complete manifolds”, *J. Partial Differ. Equ.* **23**:1 (2010), 68–79. MR 2011b:58056 Zbl 1224.58022
- [Wu 2010b] J.-Y. Wu, “Li–Yau type estimates for a nonlinear parabolic equation on complete manifolds”, *J. Math. Anal. Appl.* **369**:1 (2010), 400–407. MR 2011b:35432 Zbl 1211.58017
- [Wu and Zheng 2010] J.-Y. Wu and Y. Zheng, “Interpolating between constrained Li–Yau and Chow–Hamilton Harnack inequalities on a surface”, *Arch. Math. (Basel)* **94**:6 (2010), 591–600. MR 2011j:53127 Zbl 1198.53078
- [Yang 2008] Y. Yang, “Gradient estimates for a nonlinear parabolic equation on Riemannian manifolds”, *Proc. Amer. Math. Soc.* **136**:11 (2008), 4095–4102. MR 2009d:58048 Zbl 1151.58013
- [Yau 1994] S.-T. Yau, “On the Harnack inequalities of partial differential equations”, *Comm. Anal. Geom.* **2**:3 (1994), 431–450. MR 96f:58186 Zbl 0841.58059
- [Yau 1995] S.-T. Yau, “Harnack inequality for non-self-adjoint evolution equations”, *Math. Res. Lett.* **2**:4 (1995), 387–399. MR 96k:58211 Zbl 0884.58091
- [Zhang 2006] Q. S. Zhang, “Some gradient estimates for the heat equation on domains and for an equation by Perelman”, *Int. Math. Res. Not.* **2006** (2006), Art. ID 92314. MR 2007f:35116

Received June 25, 2011. Revised February 22, 2012.

JIA-YONG WU
 DEPARTMENT OF MATHEMATICS
 SHANGHAI MARITIME UNIVERSITY
 HAIGANG AVENUE 1550
 SHANGHAI 201306
 CHINA
 jywu81@yahoo.com

ON OVERTWISTED, RIGHT-VEERING OPEN BOOKS

PAOLO LISCA

We exhibit infinitely many overtwisted, right-veering, non-destabilizable open books, thus providing infinitely many counterexamples to a conjecture of Honda, Kazez and Matic. The page of all our open books is a four-holed sphere and the underlying 3-manifolds are lens spaces.

1. Introduction

The purpose of this note is to construct infinitely many counterexamples to a conjecture of Honda, Kazez and Matic from [Honda et al. 2009]. For the basic notions of contact topology not recalled below we refer the reader to [Etnyre 2003; Geiges 2008].

Let S be a compact, oriented surface with boundary and $\text{Map}(S, \partial S)$ the group of orientation-preserving diffeomorphisms of S that restrict to ∂S as the identity, up to isotopies fixing ∂S pointwise. An *open book* (also known as an *abstract open book*) is a pair (S, Φ) where S is a surface as above and $\Phi \in \text{Map}(S, \partial S)$. Giroux [2002] introduced a fundamental operation of *stabilization* $(S, \Phi) \rightarrow (S', \Phi')$ on open books, and proved the existence of a 1-1 correspondence between the set of open books modulo stabilization and the set of contact 3-manifolds modulo isomorphism (see, for example, [Etnyre 2006] for details). Honda, Kazez and Matic [Honda et al. 2007] showed that a contact 3-manifold is tight if and only if it corresponds to an equivalence class of open books (S, Φ) all of whose monodromies Φ are *right-veering* (in the sense of [Honda et al. 2007, Section 2]). In [Goodman 2005; Honda et al. 2007] it is also showed that every open book can be made right-veering after a sequence of stabilizations. Honda, Kazez and Matic [Honda et al. 2009] proved that when S is a holed torus, the contact structure corresponding to (S, Φ) is tight if and only if Φ is right-veering, and conjectured that a non-destabilizable right-veering open book corresponds to a tight contact 3-manifold. The Honda–Kazez–Matic conjecture was recently disproved by Lekili [2011], who produced a counterexample (S, Φ) with S equal to a four-holed sphere and whose underlying 3-manifold is the Poincaré homology sphere.

MSC2010: primary 57R17; secondary 53D10.

Keywords: contact surgery, destabilizable diffeomorphisms, Giroux's correspondence, open books, overtwisted contact structures, right-veering diffeomorphisms.

We shall now describe our examples. Denote by $\delta_\gamma \in \text{Map}(S, \partial S)$ the class of a positive Dehn twist along a simple closed curve $\gamma \subset S$.

Theorem 1.1. *Let S be an oriented four-holed sphere, and a, b, c, d, e the simple closed curves on S shown in Figure 1.*

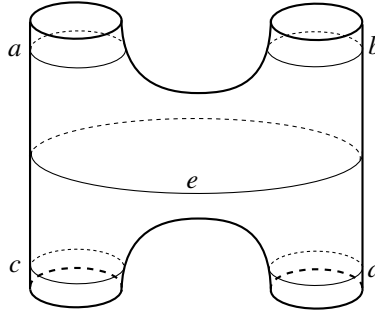


Figure 1. The four-holed sphere S .

Let $h, k \geq 1$ be integers. Define $\Phi_{h,k} := \delta_a^h \delta_b \delta_c \delta_d \delta_e^{-k-1} \in \text{Map}(S, \partial S)$. Then

- the underlying 3-manifold $Y_{(S, \Phi_{h,k})}$ is the lens space

$$L((h+1)(2k-1)+2, (h+1)k+1);$$

- the associated contact structure $\xi_{(S, \Phi_{h,k})}$ is overtwisted;
- $\Phi_{h,k}$ is right-veering;
- $(S, \Phi_{h,k})$ is not destabilizable.

Warning: in the above statement we adopt the convention that the lens space $L(p, q)$ is the oriented 3-manifold obtained by performing a rational surgery along an unknot in S^3 with coefficient $-p/q$.

We prove Theorem 1.1 in Section 2. The proof can be outlined as follows. In Proposition 2.1 we use elementary arguments to determine a contact surgery presentation for the contact 3-manifold $(Y_{(S, \Phi_{h,k})}, \xi_{(S, \Phi_{h,k})})$, and in Corollary 2.2 we apply Proposition 2.1 and a few Kirby calculus moves to identify the underlying 3-manifold $Y_{(S, \Phi_{h,k})}$. In Proposition 2.3 we appeal to calculations from [Lekili 2011] to deduce that the contact Ozsváth–Szabó invariant of $\xi_{(S, \Phi_{h,k})}$ vanishes, and we conclude from the fact that $Y_{(S, \Phi_{h,k})}$ is a lens space that $\xi_{(S, \Phi_{h,k})}$ must be overtwisted. That $\Phi_{h,k}$ is right-veering in Lemma 2.4 follows directly from [Arıkan and Dursöy 2012, Theorem 4.3], but it can also be deduced by imitating the proof of [Lekili 2011, Theorem 1.2], that is, by applying [Honda et al. 2007, Corollary 3.4]. Finally, we use results from [Arıkan 2008; Lekili 2011] to conclude that $(S, \Phi_{h,k})$ is not destabilizable.

2. Proof of Theorem 1.1

Recall that every contact structure has a *contact surgery presentation*. We refer the reader to [Ding and Geiges 2004] for this fact and the basic properties of contact surgeries, and to [Lisca and Stipsicz 2004] for the use of the “front notation” in contact surgery presentations, in particular for the meaning of Figure 2 below.

Proposition 2.1. *For $h, k \geq 1$, the contact structure $\xi_{(\mathcal{S}, \Phi_{h,k})}$ has the contact surgery presentation given by Figure 2.*

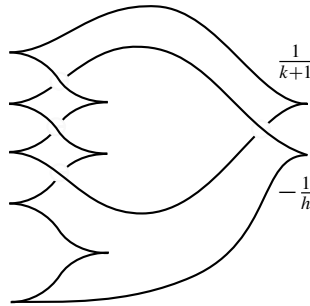


Figure 2. Contact surgery presentation for $\xi_{(\mathcal{S}, \Phi_{h,k})}$, $h, k \geq 1$.

Proof. Figure 3 (a) represents an open book (A, f) , where A is an annulus and f is a positive Dehn twist along the core of A . The associated contact 3-manifold is the standard contact 3-sphere (S^3, ξ_{st}) , the annulus A can be viewed as the page of an open book decomposition of S^3 , and the curve κ in the picture can be made Legendrian via an isotopy of the contact structure, in such a way that the contact framing on κ coincides with the framing induced on it by the page (see [Etnyre 2006, Figure 11]). The knot κ is the unique Legendrian unknot in (S^3, ξ_{st}) having Thurston–Bennequin invariant $tb(\kappa) = -1$ and rotation number $rot(\kappa) = 0$. A suitable choice of orientation for κ uniquely specifies its *negative* oriented Legendrian stabilization κ_- , which satisfies $tb(\kappa_-) = -2$ and $rot(\kappa_-) = -1$. As shown in [Etnyre 2006], κ_- can be realized as sitting on the page of a Giroux stabilization (A', f') of (A, f) . This is illustrated in Figure 3 (b), assuming the orientation on κ was taken to be “counterclockwise” in Figure 3 (a). Finally, Figure 3 (c) shows an open book (\mathcal{S}, f'') obtained by Giroux stabilizing (A', f') and containing both κ_- and $(\kappa_-)_-$ in \mathcal{S} (κ_- was also given the “counterclockwise” orientation in Figure 3 (b)). Clearly (\mathcal{S}, f'') still corresponds to (S^3, ξ_{st}) , and it is well-known that $\kappa_-, (\kappa_-)_-$ are the two Legendrian knots illustrated in Figure 2 (when oriented “clockwise” in that picture). By definition, $\Phi_{h,k}$ is obtained by precomposing f'' with $k + 1$ negative Dehn twists along parallel copies of κ_- and h positive Dehn twists along parallel copies of $(\kappa_-)_-$. Moreover, if $m \neq 0$ is an integer, $\frac{1}{m}$ -contact

surgery along any Legendrian knot λ is equivalent to $\frac{m}{|m|}$ -contact surgeries along $|m|$ Legendrian push-offs of λ [Ding and Geiges 2004]. Since page and contact framings coincide, and by [Etnyre 2006, Theorem 5.7] positive (negative, respectively) Dehn twists correspond to -1 -contact surgeries ($+1$ -contact surgeries, respectively), it is easy to check that the resulting contact structure is given by the contact surgery presentation of Figure 2. \square

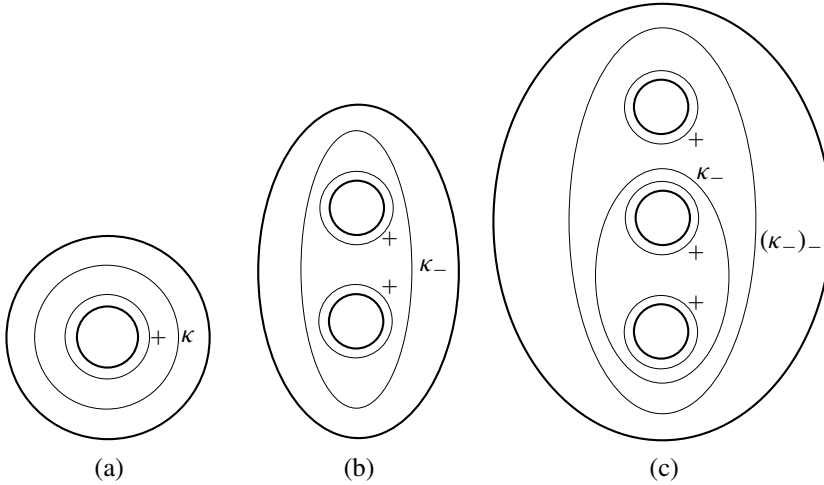


Figure 3. Determination of the contact surgery presentation.

Corollary 2.2. *For $h, k \geq 1$, the oriented 3-manifold underlying the open book $(S, \Phi_{h,k})$ is the lens space $L((h+1)(2k-1)+2, (h+1)k+1)$.*

Proof. Using the fact that the two Legendrian unknots illustrated in Figure 2 have Thurston–Bennequin invariants -2 and -3 , it is easy to check that the topological surgery underlying Figure 2 is given by the first (upper left) picture of Figure 4. Two $+1$ -blowups and two inverse slam-dunks give the second picture, while the

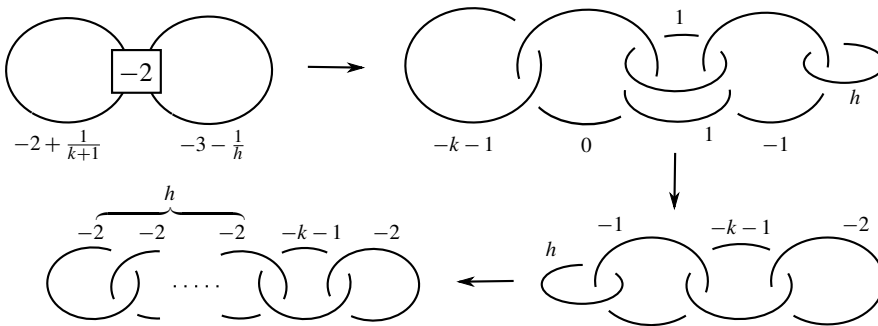


Figure 4. Determination of the underlying 3-manifold.

third picture is obtained from the second one by sliding the -1 -framed knot over the 0 -framed knot and then applying two $+1$ -blow-downs. The last picture is obtained simply converting the h -framed unknot in the third picture into the string of -2 -framed unknots via a sequence of -1 -blowups and a final $+1$ -blowdown. The last picture shows that the underlying 3 -manifold $Y_{(S, \Phi_{h,k})}$ is obtained by performing a rational surgery on an unknot in S^3 with coefficient $-p/q$, where

$$\frac{p}{q} = 2 - \frac{1}{k+1 - \frac{1}{2 - \frac{1}{\ddots - \frac{1}{2}}}} = \frac{(h+1)(2k-1)+2}{(h+1)k+1}.$$

Therefore, according to our conventions $Y_{(S, \Phi_{h,k})}$ can be identified with the lens space $L((h+1)(2k-1)+2, (h+1)k+1)$. □

Proposition 2.3. *For $h, k \geq 1$, the contact structure $\xi_{(S, \Phi_{h,k})}$ is overtwisted.*

Proof. By [Giroux 2000; Honda 2000] a contact structure on a lens space is either overtwisted or Stein fillable. Moreover, Stein fillable contact structures have nonzero contact Ozsváth–Szabó invariant [Ozsváth and Szabó 2005]. Finally, [Lekili 2011, Theorem 1.3] immediately implies that the contact invariant of $(S, \Phi_{h,k})$ vanishes, therefore $\xi_{(S, \Phi_{h,k})}$ must be overtwisted. □

Lemma 2.4. *For $h, k \geq 1$, the diffeomorphism class*

$$\Phi_{h,k} = \delta_a^h \delta_b \delta_c \delta_d \delta_e^{-k-1} \in \text{Map}(S, \partial S)$$

is right-veering.

Proof. The lemma follows immediately from the statement of Theorem 4.3 in [Arikan and Durusoy 2012]. Alternatively, one can imitate the proof of Theorem 1.2 of [Lekili 2011]. Indeed, applying Corollary 3.4 from [Honda et al. 2007] to the monodromy $\Phi_1 = \delta_e^{-k-1}$ and a properly embedded arc $\gamma_{cd} \subset S$ disjoint from the curve e and connecting the components ∂_c and ∂_d of ∂S parallel to the curves c and d shows that $\Phi_2 = \delta_d \delta_e^{-k-1}$ is right-veering with respect to ∂_d . Another application of the corollary to Φ_2 and γ_{cd} shows that $\Phi_3 = \delta_c \delta_d \delta_e^{-k-1}$ is right-veering with respect to ∂_c . Moreover, since δ_c is right-veering with respect to ∂_c and the composition of right-veering diffeomorphisms is still right-veering [Honda et al. 2007], Φ_3 is right-veering with respect to ∂_d as well. Applying the corollary in the same way to Φ_3 and an arc connecting the components of ∂S parallel to the curves a and b yields the statement of the lemma. □

Proof of Theorem 1.1. Corollary 2.2, Proposition 2.3 and Lemma 2.4 establish the first three portions of the statement. Thus we only need to show that $(S, \Phi_{h,k})$ is not destabilizable for every $h, k \geq 1$. If $(S, \Phi_{h,k})$ were destabilizable, it would be a stabilization of an open book (S', Φ') , where S' is a three-holed sphere and $\Phi' = \tau_1^{a_1} \tau_2^{a_2} \tau_3^{a_3}$, where $a_i \in \mathbb{Z}$ and τ_i is a positive Dehn twist along a simple closed curve parallel to the i -th boundary components of S' , $i = 1, 2, 3$. By [Arıkan 2008, Theorem 1.2], $\xi_{(S, \Phi_{h,k})}$ is tight if and only if $a_i \geq 0$, $i = 1, 2, 3$. Therefore, by Proposition 2.3 at least one of these exponents must be strictly negative. But the proof of Theorem 1.2 of [Lekili 2011] shows that when one of the a_i 's is negative, any stabilization of (S', Φ') to an open book with page a four-holed sphere is not right-veering. This would contradict Lemma 2.4, therefore we conclude that $(S, \Phi_{h,k})$ cannot be destabilizable. \square

Note added in proof: after the submission of the present paper the author was informed of unpublished work of A. Wand containing, in particular, a different proof of Proposition 2.3.

Acknowledgements

I wish to thank Yanki Lekili for pointing out to me his paper [Lekili 2011]. The present work is part of the author's activities within CAST, a Research Network Program of the European Science Foundation.

References

- [Arıkan 2008] M. F. Arıkan, "Planar contact structures with binding number three", pp. 90–124 in *Proceedings of Gökova Geometry-Topology Conference 2007*, edited by S. Akbulut et al., GGT, Gökova, 2008. MR 2010h:57038 Zbl 1193.57013
- [Arıkan and Durusoy 2012] M. F. Arıkan and S. Durusoy, "On the classification of certain planar contact structures", *Acta Math. Hungar.* **134**:4 (2012), 529–542. Zbl 06028279
- [Ding and Geiges 2004] F. Ding and H. Geiges, "A Legendrian surgery presentation of contact 3-manifolds", *Math. Proc. Cambridge Philos. Soc.* **136**:3 (2004), 583–598. MR 2005m:57038 Zbl 1069.57015
- [Etnyre 2003] J. B. Etnyre, "Introductory lectures on contact geometry", pp. 81–107 in *Topology and geometry of manifolds* (Athens, GA, 2001), edited by G. Matić and C. McCrory, Proc. Sympos. Pure Math. **71**, Amer. Math. Soc., Providence, RI, 2003. MR 2005b:53139 Zbl 1045.57012
- [Etnyre 2006] J. B. Etnyre, "Lectures on open book decompositions and contact structures", pp. 103–141 in *Floer homology, gauge theory, and low-dimensional topology*, edited by D. A. Ellwood et al., Clay Math. Proc. **5**, Amer. Math. Soc., Providence, RI, 2006. MR 2007g:57042 Zbl 1108.53050
- [Geiges 2008] H. Geiges, *An introduction to contact topology*, Cambridge Studies in Advanced Mathematics **109**, Cambridge University Press, 2008. MR 2008m:57064 Zbl 1153.53002
- [Giroux 2000] E. Giroux, "Structures de contact en dimension trois et bifurcations des feuilletages de surfaces", *Invent. Math.* **141**:3 (2000), 615–689. MR 2001i:53147 Zbl 1186.53097

- [Giroux 2002] E. Giroux, “Géométrie de contact: de la dimension trois vers les dimensions supérieures”, pp. 405–414 in *Proceedings of the International Congress of Mathematicians, Vol. II* (Beijing, 2002), edited by T. Li, Higher Ed. Press, Beijing, 2002. MR 2004c:53144 Zbl 1015.53049
- [Goodman 2005] N. Goodman, “Overtwisted open books from sobering arcs”, *Algebr. Geom. Topol.* **5** (2005), 1173–1195. MR 2006h:57022 Zbl 1090.57020
- [Honda 2000] K. Honda, “On the classification of tight contact structures. I”, *Geom. Topol.* **4** (2000), 309–368. MR 2001i:53148 Zbl 0980.57010
- [Honda et al. 2007] K. Honda, W. H. Kazez, and G. Matic, “Right-veering diffeomorphisms of compact surfaces with boundary”, *Invent. Mathem.* **169**:2 (2007), 427–449. MR 2008e:57028 Zbl 1167.57008
- [Honda et al. 2009] K. Honda, W. H. Kazez, and G. Matic, “On the contact class in Heegaard Floer homology”, *J. Differential Geom.* **83**:2 (2009), 289–311. MR 2011f:57050 Zbl 1186.53098
- [Lekili 2011] Y. Lekili, “Planar open books with four binding components”, *Algebr. Geom. Topol.* **11**:2 (2011), 909–928. MR 2012f:57059 Zbl 1220.57017
- [Lisca and Stipsicz 2004] P. Lisca and A. I. Stipsicz, “Ozsváth–Szabó invariants and tight contact three-manifolds. I”, *Geom. Topol.* **8** (2004), 925–945. MR 2005e:57069 Zbl 1059.57017
- [Ozsváth and Szabó 2005] P. Ozsváth and Z. Szabó, “Heegaard Floer homology and contact structures”, *Duke Math. J.* **129**:1 (2005), 39–61. MR 2006b:57043 Zbl 1083.57042

Received July 26, 2011.

PAOLO LISCA
DIPARTIMENTO DI MATEMATICA “L. TONELLI”
UNIVERSITÀ DI PISA
LARGO BRUNO PONTECORVO, 5
56127 PISA
ITALY
lisca@dm.unipi.it

WEAKLY KRULL DOMAINS AND THE COMPOSITE NUMERICAL SEMIGROUP RING $D + E[\Gamma^*]$

JUNG WOOK LIM

Let $D \subseteq E$ be an extension of integral domains, Γ a numerical semigroup with $\Gamma \subsetneq \mathbb{N}_0$, $\Gamma^* = \Gamma \setminus \{0\}$ and $R = D + E[\Gamma^*]$. In this paper, we completely characterize when R is a weakly Krull domain, an AWFD or a GWFD. We also prove that R is never a WFD.

Introduction

We first review some preliminaries. Let D be an integral domain with quotient field $qf(D)$ and let $\mathbf{F}(D)$ denote the set of nonzero fractional ideals of D . Recall that the v -operation on D is a star-operation on $\mathbf{F}(D)$ defined by $I \mapsto I_v := (I^{-1})^{-1}$, where $I^{-1} = \{x \in qf(D) \mid xI \subseteq D\}$. The t -operation on D is a star-operation defined by $I \mapsto I_t := \bigcup \{J_v \mid J \subseteq I \text{ with } J \in \mathbf{F}(D) \text{ finitely generated}\}$. An $I \in \mathbf{F}(D)$ is said to be a v -ideal if $I_v = I$, and a t -ideal if $I_t = I$. A v -ideal I is said to be of *finite type* if $I = J_v$ for some finitely generated fractional ideal J of D . A t -ideal M of D is called a *maximal t -ideal* if M is maximal among proper integral t -ideals of D . It is well known that maximal t -ideals are prime ideals. Let $t\text{-Max}(D)$ be the set of maximal t -ideals of D . Then $t\text{-Max}(D) \neq \emptyset$ if D is not a field. An $I \in \mathbf{F}(D)$ is said to be t -invertible if $(II^{-1})_t = D$; equivalently, $II^{-1} \not\subseteq M$ for each $M \in t\text{-Max}(D)$. Let $T(D)$ be the abelian group of t -invertible fractional t -ideals of D under the t -multiplication $I * J = (IJ)_t$, and let $\text{Inv}(D)$ and $\text{Prin}(D)$ be the subgroups of $T(D)$ consisting respectively of invertible fractional ideals of D and nonzero principal fractional ideals of D . Then it is clear that $\text{Prin}(D) \subseteq \text{Inv}(D) \subseteq T(D)$. The t -class group of D is an abelian group $\text{Cl}(D) = T(D)/\text{Prin}(D)$ and the *Picard group* $\text{Pic}(D) = \text{Inv}(D)/\text{Prin}(D)$ is a subgroup of $\text{Cl}(D)$. The *local t -class group* $G(D)$ of D is defined by $G(D) = \text{Cl}(D)/\text{Pic}(D)$.

Let $X^1(D)$ stand for the set of height-one prime ideals of D . We say that D is a *weakly Krull domain* if $D = \bigcap_{P \in X^1(D)} D_P$ and this intersection has finite character, i.e., each nonzero element $d \in D$ is a unit in D_P for all but a finite number of P 's in $X^1(D)$; D is a *weakly factorial domain* (WFD) if every nonzero nonunit element of D is a product of primary elements; D is an *almost weakly factorial domain*

MSC2010: primary 13A15, 13G05; secondary 13A02, 13B25, 13F05.

Keywords: numerical semigroup, $D + E[\Gamma^*]$, weakly Krull domain.

(AWFD) if for each nonzero nonunit element $d \in D$, there exists a positive integer $n = n(d)$ such that d^n is a product of primary elements; and D is a *generalized weakly factorial domain* (GWFD) if each nonzero prime ideal of D contains a primary element. (Recall that a nonzero nonunit $d \in D$ is called a *primary element* of D if (d) is a primary ideal of D .) It is well known that

$$\text{WFD} \Rightarrow \text{AWFD} \Rightarrow \text{GWFD} \Rightarrow \text{weakly Krull domain}$$

and a weakly Krull domain has t -dimension one. (The t -dimension of D , abbreviated $t\text{-dim}(D)$, is the supremum of lengths of chains of prime t -ideals of D . Hence $t\text{-dim}(D) = 1$ if and only if each maximal t -ideal of D has height-one.) Also, it was shown in [Anderson and Zafrullah 1990, Theorem] that a weakly Krull domain D is a WFD if and only if $\text{Cl}(D) = 0$, and in [Anderson et al. 1992, Theorem 3.4] that a weakly Krull domain D is an AWFD if and only if $\text{Cl}(D)$ is torsion. We note that $t\text{-dim}(D[\Gamma]) = t\text{-dim}(D[X])$ for any numerical semigroup Γ [Chang et al. 2012, Theorem 1.5].

Let \mathbb{N}_0 (resp., \mathbb{Z}) be the set of nonnegative integers (resp., integers). A semigroup Γ is called a *numerical semigroup* if Γ is a subset of \mathbb{N}_0 containing 0 and generates \mathbb{Z} as a group. It is known that if Γ is a numerical semigroup, then Γ is finitely generated and $\mathbb{N}_0 \setminus \Gamma$ is a finite set. Hence there exists the largest nonnegative integer which is not contained in Γ . This number is called the *Frobenius number* of Γ and is denoted by $F(\Gamma)$.

Throughout this article, $D \subseteq E$ denotes an extension of integral domains, $qf(D)$ (resp., $qf(E)$) is the quotient field of D (resp., E), \bar{D} means the integral closure of D , X is an indeterminate over E , Γ is a numerical semigroup with $\Gamma \subsetneq \mathbb{N}_0$ and $D[\Gamma]$ is the numerical semigroup ring of Γ over D . Note that each element $f \in D[\Gamma]$ is uniquely expressible in the form $f = a_1 X^{\alpha_1} + \cdots + a_k X^{\alpha_k}$, where $a_i \in D$ and $\alpha_i \in \Gamma$ with $\alpha_1 < \cdots < \alpha_k$. Let $\Gamma^* = \Gamma \setminus \{0\}$, $R = D + E[\Gamma^*]$, $T = D + XE[X]$ and $T_n = D + X^n E[X]$ for integers $n \geq 2$, i.e., $R = \{f \in E[\Gamma] \mid f(0) \in D\}$, $T = \{f \in E[X] \mid f(0) \in D\}$ and $T_n = R$ when $\Gamma = \{0\} \cup \{m \in \mathbb{N}_0 \mid m \geq n\}$. Then $D[\Gamma] \subseteq R \subseteq E[\Gamma]$ and $T_{F(\Gamma)+1} \subseteq R \subsetneq T \subseteq E[X]$. For an $f \in qf(D)[\Gamma]$, $c(f)$ means the fractional ideal of D generated by the coefficients of f . If I is an ideal of $D[\Gamma]$, then $c(I)$ denotes the ideal of D generated by the coefficients of all the polynomials in I .

In multiplicative ideal theory, the $D + E[\Gamma^*]$ construction has been extensively studied by several authors for its interest in constructing examples with prescribed properties. As a special kind of pullbacks, this has become so important that in recent years there have been many papers devoted to ring- and ideal-theoretic properties in this construction.

Anderson et al. [2003a; 2006] (see also [Anderson and Chang 2007]) studied when the domains $D[X^2, X^3]$, $D + XE[X]$ and $D + X^2E[X]$ are weakly Krull

domains, WFDs, AWFDs or GWFDs. In fact, they showed that $D[X^2, X^3]$ is a weakly Krull domain if and only if D is a weakly Krull UMT-domain [Anderson et al. 2003a, Proposition 2.7]; if $\text{char}(D) \neq 0$, then $D[X^2, X^3]$ is an AWFD if and only if $D[X^2, X^3]$ is a GWFD [Anderson and Chang 2007, Corollary 2.11]; $D + XE[X]$ is a weakly Krull domain if and only if $D + X^2E[X]$ is a weakly Krull domain [Anderson et al. 2006, Theorem 4.3]; and $D + XE[X]$ is an AWFD if and only if $D + XE[X]$ is a GWFD [Anderson and Chang 2007, Corollary 2.10]. The main purpose of this paper is to determine how certain properties of D , E and Γ influence those of R , and vice versa. This also extends the results for the domains $D[X^2, X^3]$, $D + XE[X]$ and $D + X^2E[X]$ to any composite numerical semigroup ring $D + E[\Gamma^*]$.

In Section 1, we investigate weakly Krull domains, AWFDs and GWFDs in the context of numerical semigroup rings $D[\Gamma]$ which coincide with the domains $R = D + E[\Gamma^*]$ when $D = E$. We prove that $D[\Gamma]$ is a weakly Krull domain if and only if D is a weakly Krull UMT-domain, and that if $\text{char}(D) \neq 0$, then $D[\Gamma]$ is an AWFD if and only if $D[\Gamma]$ is a GWFD, if and only if D is an almost weakly factorial quasi-AGCD-domain, if and only if D is a generalized weakly factorial quasi-AGCD-domain.

In Section 2, we study when the domain $R = D + E[\Gamma^*]$ is a weakly Krull domain, an AWFD or a GWFD, where $D \subsetneq E$. We show that R is a weakly Krull domain if and only if $T = D + XE[X]$ is a weakly Krull domain, and that if $\text{char}(E) \neq 0$, then R is an AWFD if and only if R is a GWFD, if and only if T is an AWFD, if and only if R is a GWFD. We also prove that R is never a WFD.

1. Weakly Krull domains as numerical semigroup rings

In this section, we characterize when the numerical semigroup ring $D[\Gamma]$ is a weakly Krull domain, an AWFD or a GWFD.

The first two lemmas are well known for the general semigroup rings, but we include their proofs for the convenience of the reader.

Lemma 1.1 [El Baghdadi et al. 2002, Lemma 2.3]. *Let D be an integral domain and Γ be a numerical semigroup. The following statements hold for an $I \in \mathbf{F}(D)$:*

- (1) $(ID[\Gamma])^{-1} = I^{-1}D[\Gamma]$.
- (2) $(ID[\Gamma])_v = I_vD[\Gamma]$.
- (3) $(ID[\Gamma])_t = I_tD[\Gamma]$.

Proof. (1) Since $(ID[\Gamma])(I^{-1}D[\Gamma]) \subseteq D[\Gamma]$, $I^{-1}D[\Gamma] \subseteq (ID[\Gamma])^{-1}$. Conversely, let $f \in (ID[\Gamma])^{-1}$. Then $fID[\Gamma] \subseteq D[\Gamma]$ and hence $c(f)I \subseteq D$. Hence $c(f) \subseteq I^{-1}$, and therefore $f \in c(f)D[\Gamma] \subseteq I^{-1}D[\Gamma]$. Thus the equality holds.

(2) By (1), $(ID[\Gamma])_v = ((ID[\Gamma])^{-1})^{-1} = (I^{-1}D[\Gamma])^{-1} = I_vD[\Gamma]$.

(3) Let f_1, \dots, f_n be nonzero elements of $ID[\Gamma]$. Then we have

$$\begin{aligned} ((f_1, \dots, f_n)D[\Gamma])_v &\subseteq ((c(f_1), \dots, c(f_n))D[\Gamma])_v \\ &= (c(f_1), \dots, c(f_n))_v D[\Gamma] \\ &\subseteq I_t D[\Gamma] \end{aligned}$$

by (2), i.e., $(ID[\Gamma])_t \subseteq I_t D[\Gamma]$. For the reverse inclusion, let J be a nonzero finitely generated subideal of I . Then $J_v D[\Gamma] = (JD[\Gamma])_v \subseteq (ID[\Gamma])_t$ by (2). Hence $I_t D[\Gamma] \subseteq (ID[\Gamma])_t$. Thus we have the desired equality. \square

Lemma 1.2 [Anderson and Chang 2005, Corollary 2.3]. *Let D be an integral domain, Γ be a numerical semigroup and let Q be a maximal t -ideal of $D[\Gamma]$ such that $Q \cap D \neq (0)$. Then $Q = (Q \cap D)D[\Gamma]$. In particular, $Q \cap D$ is a maximal t -ideal of D .*

Proof. The containment $(Q \cap D)D[\Gamma] \subseteq Q$ is obvious. For the converse, it suffices to show that $c(Q) \subseteq Q$. Suppose to the contrary that $c(Q) \not\subseteq Q$. Then

$$Q \subsetneq c(Q)D[\Gamma].$$

Since Q is a maximal t -ideal of $D[\Gamma]$, $(c(Q)D[\Gamma])_t = D[\Gamma]$. Therefore $c(Q)_t = D$ by Lemma 1.1(3), and hence $c(f)_v = D$ for some $f \in Q$. Let $0 \neq d \in Q \cap D$ and choose $0 \neq g \in (d, f)^{-1}$. Then $gd \in D[\Gamma]$ and hence $g \in qf(D)[\Gamma]$. Also, we have $fg \in D[\Gamma]$. Hence it follows from [Gilmer 1992, Theorem 28.1] that

$$c(g) \subseteq c(g)_v = (c(f)^{m+1}c(g))_v = (c(f^m)c(fg))_v = c(fg)_v \subseteq D,$$

where m is the degree of g . So $g \in c(g)D[\Gamma] \subseteq D[\Gamma]$, which implies that $(d, f)^{-1} = D[\Gamma]$. This contradicts the fact that Q is a maximal t -ideal of $D[\Gamma]$. Therefore $c(Q) \subseteq Q$, and thus $Q \subseteq (Q \cap D)D[\Gamma]$. The second assertion is an immediate consequence of Lemma 1.1(3). \square

An integral domain B is said to be a *UMT-domain* if every upper to zero (a nonzero prime ideal of $B[X]$ which contracts to zero in B) Q of $B[X]$ is a maximal t -ideal (equivalently, is t -invertible). Now, we give the numerical semigroup ring version of [Anderson et al. 1993, Proposition 4.11].

Theorem 1.3. *Let D be an integral domain and Γ be a numerical semigroup with $\Gamma \subsetneq \mathbb{N}_0$. Then the following assertions are equivalent.*

- (1) $D[\Gamma]$ is a weakly Krull domain.
- (2) $D[X]$ is a weakly Krull domain.
- (3) D is a weakly Krull UMT-domain.

Proof. (1) \Rightarrow (3) Assume $D[\Gamma]$ is a weakly Krull domain. Then $t\text{-dim}(D[\Gamma]) = 1$ [Anderson et al. 1992, Lemma 2.1]. Let P be a prime t -ideal of D . Then $PD[\Gamma]$ is a prime t -ideal of $D[\Gamma]$ by Lemma 1.1(3); so $\text{ht}_D(P) = \text{ht}_{D[\Gamma]}(PD[\Gamma]) = 1$; so $t\text{-dim}(D) = 1$. Since $t\text{-dim}(D[\Gamma]) = 1$, we have $t\text{-dim}(D[X]) = 1$ by [Chang et al. 2012, Theorem 1.5]. Therefore every upper to zero in $D[X]$ is a maximal t -ideal, and thus D is a UMT-domain. Note that

$$D = \bigcap_{P \in X^1(D)} D_P$$

by [Kang 1989, Proposition 2.9]. To show that this intersection has finite character, let $d \in D \setminus \{0\}$. Since $D[\Gamma]$ is a weakly Krull domain, d belongs to only finitely many height-one prime ideals of $D[\Gamma]$, and hence there exists only a finite number of height-one prime ideals of D containing d . Thus D is a weakly Krull domain.

(3) \Rightarrow (1) Assume that D is a weakly Krull UMT-domain and let Q be a maximal t -ideal of $D[\Gamma]$ with $Q \cap D \neq (0)$. By Lemma 1.2, $Q = (Q \cap D)D[\Gamma]$ and $Q \cap D$ is a maximal t -ideal of D . Since $t\text{-dim}(D) = 1$ [Anderson et al. 1992, Lemma 2.1], $\text{ht}_D(Q \cap D) = 1$; so $\text{ht}_{D[\Gamma]}Q \leq 2$ (cf. [Kaplansky 1970, Theorem 37]). If $\text{ht}_{D[\Gamma]}Q = 2$, then there exists a nonzero prime ideal $P \subsetneq Q$ which contracts to zero in D . Note that $P = M \cap D[\Gamma]$ for some prime ideal M of $D[X]$ [Chang et al. 2012, Proposition 1.1]. Since $M \cap D = (0)$ and D is a UMT-domain, M is a maximal t -ideal of $D[X]$. Hence P is a maximal t -ideal of $D[\Gamma]$ [Chang et al. 2012, Theorem 1.4]. This contradicts the choice of P . Thus $t\text{-dim}(D[\Gamma]) = 1$. By [Kang 1989, Proposition 2.9], we have $D[\Gamma] = \bigcap_{Q \in X^1(D[\Gamma])} D[\Gamma]_Q$. We claim that this intersection has finite character. Let $f \in D[\Gamma] \setminus \{0\}$ and set

$$\mathcal{S} = \{Q \in X^1(D[\Gamma]) \mid f \in Q\},$$

$$\mathcal{S}_1 = \{Q \in \mathcal{S} \mid Q \cap D \in X^1(D)\}, \text{ and}$$

$$\mathcal{S}_2 = \{Q \in \mathcal{S} \mid Q \cap D = (0)\}.$$

Then $\mathcal{S} = \mathcal{S}_1 \cup \mathcal{S}_2$. If \mathcal{S}_1 is an infinite set, then $c(f)$ belongs to infinitely many height-one prime ideals of D by Lemma 1.2. This is absurd, because D is a weakly Krull domain. Hence \mathcal{S}_1 is a finite set. Note that $qf(D)[\Gamma]$ is a one-dimensional Noetherian domain; so $qf(D)[\Gamma]$ is a weakly Krull domain. Hence \mathcal{S}_2 is also a finite set. Therefore \mathcal{S} is a finite set. Thus $D[\Gamma]$ is a weakly Krull domain.

(2) \Leftrightarrow (3) See [Anderson et al. 1993, Proposition 4.11]. \square

Recall that if $D \subseteq E$ is an extension of integral domains, then E is said to be a *root extension* of D if for each $z \in E$, there is a positive integer $n = n(z)$ such that $z^n \in D$. A domain B is called an *almost Prüfer v -multiplication domain* (AP v MD) (resp., *almost GCD-domain* (AGCD-domain)) if for each $0 \neq a, b \in B$, there exists a positive integer $n = n(a, b)$ such that $(a^n, b^n)_v$ is t -invertible (resp., principal).

It is known that B is a weakly Krull PvMD if and only if $B[X]$ is weakly Krull and B is integrally closed [Anderson et al. 1993, Corollary 4.13]. We weaken the hypothesis and obtain the following result.

Corollary 1.4. *Let D be an integral domain and Γ be a numerical semigroup.*

- (1) *D is a weakly Krull APvMD if and only if $D[\Gamma]$ is a weakly Krull domain and $D \subseteq \bar{D}$ is a root extension.*
- (2) *D is an almost weakly factorial AGCD-domain if and only if $D[\Gamma]$ is a weakly Krull domain, $\text{Cl}(D)$ is torsion and $D \subseteq \bar{D}$ is a root extension.*

Proof. (1) By [Li 2012, Theorem 3.8], a domain B is an APvMD if and only if B is a UMT-domain and $B \subseteq \bar{B}$ is a root extension. Thus the result follows from Theorem 1.3.

(2) By [Li 2012, Theorem 3.1], a domain B is an AGCD-domain if and only if B is an APvMD and $\text{Cl}(B)$ is torsion. Also, by [Anderson et al. 1992, Theorem 3.4], B is an AWFD if and only if B is a weakly Krull domain and $\text{Cl}(B)$ is torsion. Thus the result is an immediate consequence of Theorem 1.3 and (1). \square

Let S be a saturated multiplicative subset of a domain B and let $N(S) = \{0 \neq b \in B \mid (b, s)_v = B \text{ for all } s \in S\}$ be the m -complement of S . We say that S is an *almost splitting set* if for each $0 \neq b \in B$, there exists a positive integer $n = n(b)$ such that $b^n = st$ for some $s \in S$ and $t \in N(S)$. Following [Anderson and Chang 2007], B is called a *quasi-AGCD-domain* if $B \setminus \{0\}$ is an almost splitting set in $B[X]$. It was shown that if B is integrally closed, then the notion of quasi-AGCD-domains coincides with that of AGCD-domains [Chang 2005, Proposition 2.6]. The next corollary characterizes when the numerical semigroup ring $D[\Gamma]$ is an AWFD or a GWFD.

Corollary 1.5. *Let D be an integral domain with $\text{char}(D) \neq 0$ and Γ be a numerical semigroup with $\Gamma \subsetneq \mathbb{N}_0$. Then the following conditions are equivalent.*

- (1) *$D[\Gamma]$ is an AWFD.*
- (2) *$D[\Gamma]$ is a GWFD.*
- (3) *$D[X]$ is an AWFD.*
- (4) *$D[X]$ is a GWFD.*
- (5) *D is an almost weakly factorial quasi-AGCD-domain.*
- (6) *D is a generalized weakly factorial quasi-AGCD-domain.*
- (7) *D is a weakly Krull quasi-AGCD-domain.*

Proof. Let $\text{char}(D) = p$.

(1) \Rightarrow (2) This is well known.

(1) \Leftrightarrow (3) By [Anderson et al. 1992, Theorem 3.4], an integral domain B is an AWFD if and only if B is a weakly Krull domain and $\text{Cl}(B)$ is torsion, and by Theorem 1.3, $D[\Gamma]$ is a weakly Krull domain if and only if $D[X]$ is a weakly Krull domain. By [Chang et al. 2012, Lemma 2.7], $\text{Pic}(qf(D)[\Gamma])$ is torsion if and only if $\text{char}(D) \neq 0$. Since $\text{Cl}(D[\Gamma]) = \text{Cl}(D[X]) \oplus \text{Pic}(qf(D)[\Gamma])$ [Anderson and Chang 2004, Theorem 5], $\text{Cl}(D[\Gamma])$ is torsion if and only if $\text{Cl}(D[X])$ is torsion and $\text{char}(D) \neq 0$. Thus this equivalence follows from these facts.

(4) \Rightarrow (2) By [Anderson et al. 2003b, Theorem 2.2], a domain B is a GWFD if and only if $t\text{-dim}(B) = 1$ and for each $P \in X^1(B)$, $P = \sqrt{bB}$ for some $b \in B$. Assume that $D[X]$ is a GWFD and let $P \in X^1(D[\Gamma])$. Since $t\text{-dim}(D[\Gamma]) = t\text{-dim}(D[X]) = 1$ [Chang et al. 2012, Theorem 1.5], it suffices to show that $P = \sqrt{fD[\Gamma]}$ for some $f \in D[\Gamma]$. If $P \cap D \neq (0)$, then $P = (P \cap D)D[\Gamma]$ by Lemma 1.2. Since $D[X]$ is a GWFD, $(P \cap D)D[X] = \sqrt{dD[X]}$ for some $d \in P \cap D$. It is easy to see that $P = \sqrt{dD[\Gamma]}$. Next, suppose that $P \cap D = (0)$. Then there exists a prime t -ideal Q of $D[X]$ such that $P = Q \cap D[\Gamma]$ [Chang et al. 2012, Theorem 1.5]. Since $D[X]$ is a GWFD, $Q = \sqrt{fD[X]}$ for some $f \in D[X]$. Also, since $\text{char}(D) = p > 0$, there exists a positive integer n such that $f^{p^n} \in D[\Gamma]$. An easy calculation shows that $P = \sqrt{f^{p^n}D[\Gamma]}$. Thus $D[\Gamma]$ is a GWFD.

(2) \Rightarrow (4) This direction is an easy modification of the proof of (4) \Rightarrow (2).

(2) \Rightarrow (5) See [Anderson and Chang 2007, Corollary 2.9].

(5) \Rightarrow (6) \Rightarrow (7) These implications are obvious.

(7) \Rightarrow (1) Assume that D is a weakly Krull quasi-AGCD-domain. Then D is a UMT-domain and $\text{Cl}(D[X])$ is torsion [Anderson and Chang 2007, Theorem 2.4]. Hence $D[\Gamma]$ is a weakly Krull domain by Theorem 1.3. Also, it follows from [Anderson and Chang 2004, Theorem 5; Chang et al. 2012, Lemma 2.7] that $\text{Cl}(D[\Gamma])$ is torsion. Thus $D[\Gamma]$ is an AWFD [Anderson et al. 1992, Theorem 3.4]. \square

We end this section by noting that $D[\Gamma]$ is never a WFD. We also show that $D[\Gamma]$ need not be an AWFD if $\text{char}(D) = 0$.

Remark 1.6. (1) Let B be an integral domain with quotient field K . In [Gilmer and Martin 1990, Theorem 7], Gilmer and Martin showed that if B is a seminormal domain and $B + X^n B[X] \subseteq B[\Gamma]$, then $\text{Pic}(B[\Gamma]) = \text{Pic}(B) \oplus (W_n/L)$, where $L \subseteq W_n$ are the subgroups of the group $U(B[X]/X^n B[X])$ of units of $B[X]/X^n B[X]$ defined by $W_n = \{1 + Xf + X^n B[X] \mid f \in B[X]\}$ and $L = \{1 + Xf + X^n B[X] \mid 1 + Xf \in B[\Gamma]\}$. Note that $\text{Cl}(B[\Gamma]) = \text{Cl}(B[X]) \oplus \text{Pic}(K[\Gamma])$ [Anderson and Chang 2004, Theorem 5] and that B is a WFD if and only if B is a weakly Krull domain and $\text{Cl}(B) = 0$ [Anderson and Zafrullah 1990, Theorem]. If $D[\Gamma]$ is a WFD, then $\text{Cl}(D[\Gamma]) = 0$, and hence $\text{Pic}(qf(D)[\Gamma]) = 0$. Therefore $W_n = L$;

so $1 + X + X^n qf(D)[X] \in L$, which implies that $1 \in \Gamma$. Thus, if Γ is a proper numerical semigroup, then $D[\Gamma]$ is never a WFD.

(2) If $D[\Gamma]$ is an AWFD, then $\text{Cl}(D[\Gamma])$ is torsion [Anderson et al. 1992, Theorem 3.4]; so $\text{Pic}(qf(D)[\Gamma])$ is torsion [Anderson and Chang 2004, Theorem 5]. Hence $\text{char}(D) \neq 0$ [Chang et al. 2012, Lemma 2.7]. This shows that the condition that $\text{char}(D) \neq 0$ is essential in Corollary 1.5.

(3) It is known that a generalized unique factorization domain (GUFD) is a weakly factorial GCD-domain [Anderson et al. 1995, Theorem 7], and hence integrally closed. (See [Anderson et al. 1995] for the definition and some characterizations of a GUFD.) Thus, if Γ is a numerical semigroup with $\Gamma \subsetneq \mathbb{N}_0$, then $D[\Gamma]$ is not a GUFD by (1). In fact, $D[\Gamma]$ is not integrally closed; so $D[\Gamma]$ is never a GUFD.

2. Weakly Krull domains and the ring $D + E[\Gamma^*]$ when $D \subsetneq E$

For a domain A , $\text{Spec}(A)$ stands for the set of prime ideals of A . Assume that $D \subsetneq E$ is an extension of integral domains, Γ is a numerical semigroup with $\Gamma \subsetneq \mathbb{N}_0$ and let $R = D + E[\Gamma^*]$, $T = D + XE[X]$, $T_n = D + X^n E[X]$ and $\Delta_n = \{0\} \cup \{m \in \mathbb{N}_0 \mid m \geq n\}$ for integers $n \geq 2$. Note that $D[\Gamma] \subsetneq R \subsetneq T$ and $T_n \subsetneq T$. In this section, we characterize when the domains R and T_n are weakly Krull domains, AWFDs or GWFDs. To do this, we need two lemmas.

Lemma 2.1. *Let $R = D + E[\Gamma^*]$ and $T = D + XE[X]$. If Q is a prime ideal of R , then there exists a unique prime ideal of T lying over Q . Thus the natural map $\phi : \text{Spec}(T) \rightarrow \text{Spec}(R)$, given by $P \mapsto P \cap R$, is an order-preserving bijection. In particular, $\text{ht}_T(XE[X]) = \text{ht}_R(E[\Gamma^*])$.*

Proof. Let Q be a prime ideal of R . Since T is an integral extension of R , there exists a prime ideal P of T such that $Q = P \cap R$ [Kaplansky 1970, Theorem 44]. Note that $E[\Gamma^*] \subseteq Q$ if and only if $XE[X] \subseteq P$. If $E[\Gamma^*] \subseteq Q$, then P is the unique prime ideal of T lying over Q because $R/XE[X] \cong D \cong R/E[\Gamma^*]$. If $E[\Gamma^*] \not\subseteq Q$, then $X^{F(\Gamma)+1} f \notin Q$ for some $f \in E[X]$; so

$$g = \frac{X^{F(\Gamma)+1} fg}{X^{F(\Gamma)+1} f} \in R_Q$$

for any $g \in T$. Hence $T_{QR_Q \cap T} = R_Q$. Thus $QR_Q \cap T$ is the unique prime ideal of T lying over Q . □

Let n be an integer ≥ 2 . Then it is clear that if $\Gamma = \Delta_n$, then $R = T_n$. Hence Lemma 2.1 also shows that $\text{ht}_T(XE[X]) = \text{ht}_{T_n}(X^n E[X])$.

Remark 2.2. Let $\Gamma = \{\alpha_1, \dots, \alpha_n\} \cup \Delta_{F(\Gamma)+1}$ with $1 < \alpha_1 < \dots < \alpha_n < F(\Gamma) + 1$ and $R = D + E[\Gamma^*]$.

(1) Let $g \in (R : E[\Gamma^*])$. Then $gE[\Gamma^*] \subseteq R$; hence for each $\alpha \in \Gamma^*$, $gX^\alpha = a_\alpha + f_\alpha$ for some $a_\alpha \in D$ and $f_\alpha \in E[\Gamma^*]$. Therefore $gX^{\alpha+F(\Gamma)} = (a_\alpha + f_\alpha)X^{F(\Gamma)} \in R$, which means that $a_\alpha = 0$. Hence $gX^\alpha = f_\alpha \in E[\Gamma^*]$, and so $g \in \bigcap_{\alpha \in \Gamma^*} \{\frac{1}{X^\alpha} f \mid f \in E[\Gamma^*]\}$. The reverse containment is obvious. Thus we have

$$(R : E[\Gamma^*]) = \bigcap_{\alpha \in \Gamma^*} \left\{ \frac{1}{X^\alpha} f \mid f \in E[\Gamma^*] \right\}.$$

(2) It is clear that $E[\Gamma] \subsetneq (R : E[\Gamma^*])$ because $X^{F(\Gamma)} \in (R : E[\Gamma^*]) \setminus E[\Gamma]$. Let $g \in (R : E[\Gamma^*])$. Then $X^{F(\Gamma)+1}g \in R$; so we can write

$$X^{F(\Gamma)+1}g = \sum_{i=0}^n g_i X^{\alpha_i} + X^{F(\Gamma)+1}h$$

for some $g_i \in E$ and $h \in E[X]$. (For the sake of convenience, set $\alpha_0 = 0$.) Fix a $k \in \{1, \dots, n\}$. Then we have $X^{2F(\Gamma)-\alpha_k+1}g = \sum_{i=0}^{k-1} g_i X^{F(\Gamma)+\alpha_i-\alpha_k} + g_k X^{F(\Gamma)} + X^{F(\Gamma)+1}(\sum_{i=k+1}^n g_i X^{\alpha_i-\alpha_k-1} + h) \in R$; so $g_k = 0$ for all $k = 1, \dots, n$. Also, we have $X^{F(\Gamma)+2}g = g_0 X + X^{F(\Gamma)+2}h \in R$; so $g_0 = 0$. Therefore $X^{F(\Gamma)+1}g = X^{F(\Gamma)+1}h$, and hence $g = h \in E[X]$. Thus $E[\Gamma] \subsetneq (R : E[\Gamma^*]) \subseteq E[X]$. In particular, if $\Gamma = \Delta_{F(\Gamma)+1}$, then $E[X] \subseteq (R : E[\Gamma^*])$; so $(R : E[\Gamma^*]) = E[X]$.

(3) Lemma 4.2 of [Anderson et al. 2006] cannot be extended to any proper numerical semigroup, i.e., it may happen that $(R : E[\Gamma^*]) \subsetneq E[X]$ for some $\Gamma \subsetneq \mathbb{N}_0$. For instance, if $\Gamma = \{2\} \cup \Delta_4$, then $X \in E[X] \setminus (R : E[\Gamma^*])$.

Lemma 2.3. *The following statements hold for $R = D + E[\Gamma^*]$.*

- (1) $E[\Gamma^*]$ is a prime t -ideal of R .
- (2) $E[\Gamma^*]$ is a maximal t -ideal of R if and only if $qf(D) \cap E = D$.

Proof. (1) Let $\Gamma = \{\alpha_1, \dots, \alpha_k\} \cup \Delta_{F(\Gamma)+1}$ such that $0 < \alpha_1 < \dots < \alpha_k < F(\Gamma) + 1$. Since $R/E[\Gamma^*] \cong D$, $E[\Gamma^*]$ is a prime ideal of R . It suffices to show that $E[\Gamma^*]$ is a v -ideal of R , because each v -ideal is a t -ideal.

Case 1. $\{\alpha_1, \dots, \alpha_k\}$ is empty. In this case, $(R : E[\Gamma^*]) = E[X]$ by Remark 2.2(2); so we need to show that $(R : E[X]) = E[\Gamma^*]$. It is clear that $E[\Gamma^*] \subseteq (R : E[X])$. For the converse, let $f \in (R : E[X])$. Then $fE[X] \subseteq R$. Since $1 \in E[X]$, $f \in R$. Also, since $X \in E[X]$, $f(0) = 0$; so $f \in E[\Gamma^*]$.

Case 2. $\{\alpha_1, \dots, \alpha_k\}$ is nonempty. Deny the conclusion, and then there exists a polynomial $g = g_0 + \sum_{i=1}^k g_{\alpha_i} X^{\alpha_i} + \sum_{i=F(\Gamma)+1}^l g_i X^i \in (E[\Gamma^*])_v \setminus E[\Gamma^*]$. Hence $g(R : E[\Gamma^*]) \subseteq R$. Let $f \in (R : E[\Gamma^*])$. Then $f \in E[X]$ by Remark 2.2(2); so we can write $f = \sum_{i=0}^m f_i X^i$. Note that

$$fg = f_0 g_0 + g_0 \sum_{i=1}^{\alpha_1-1} f_i X^i + (f_0 g_{\alpha_1} + f_{\alpha_1} g_0) X^{\alpha_1} + X^{\alpha_1+1} h_1$$

for some $h_1 \in E[X]$. Since $fg \in R$ and $g_0 \neq 0$, $f_1 = \dots = f_{\alpha_1-1} = 0$; so $f = f_0 + \sum_{i=\alpha_1}^m f_i X^i$. Note that $2\alpha_1 \in \Gamma^*$; so $2\alpha_1 \geq F(\Gamma) + 1$ or $2\alpha_1 = \alpha_p$ for some $p \in \{2, \dots, k\}$. If $2\alpha_1 \geq F(\Gamma) + 1$, then we have

$$fg = f_0g_0 + (f_0g_{\alpha_1} + f_{\alpha_1}g_0)X^{\alpha_1} + g_0 \sum_{i=\alpha_1+1}^{\alpha_2-1} f_i X^i + (f_0g_{\alpha_2} + f_{\alpha_2}g_0)X^{\alpha_2} + X^{\alpha_2+1}h_2$$

for some $h_2 \in E[X]$. Again, since $fg \in R$, $f_{\alpha_1+1} = \dots = f_{\alpha_2-1} = 0$. By repeating this process, we have $f_i = 0$ for all $i \in \mathbb{N}_0 \setminus \Gamma$, and hence $f \in R$. Therefore $(R : E[\Gamma^*]) = R$. However, this is impossible because $X^{F(\Gamma)} \in (R : E[\Gamma^*]) \setminus R$. If $2\alpha_1 = \alpha_p$ for some $p \in \{2, \dots, k\}$, a simple modification of the proof of the previous case leads to the same conclusion because $2\alpha_l \geq F(\Gamma) + 1$ for some $l \leq k$. In either case, $E[\Gamma^*]$ is a v -ideal, and thus $E[\Gamma^*]$ is a t -ideal of R .

(2) This appears in [Lim 2012, Lemma 1.2]. □

Now, we are ready to give a necessary and sufficient condition for the domain R to be a weakly Krull domain.

Theorem 2.4. *Let $R = D + E[\Gamma^*]$, $T = D + XE[X]$, $T_n = D + X^nE[X]$ and $\Delta_n = \{0\} \cup \{m \in \mathbb{N}_0 \mid m \geq n\}$ for integers $n \geq 2$. Then the following statements are equivalent.*

- (1) R is a weakly Krull domain.
- (2) T is a weakly Krull domain.
- (3) T_n is a weakly Krull domain.
- (4) $X^nE[X]$ is a height-one maximal t -ideal of T_n and $E[\Delta_n]$ is a weakly Krull domain.
- (5) $E_{D \setminus \{0\}}$ is a field, $qf(D) \cap E = D$ and $E[X]$ is a weakly Krull domain.

Proof. (2) \Rightarrow (1) Let T be a weakly Krull domain. Let $\Gamma = \{\alpha_1, \dots, \alpha_k\} \cup \Delta_{F(\Gamma)+1}$ be such that $0 < \alpha_1 < \dots < \alpha_k < F(\Gamma) + 1$. Then $T = \bigcap_{P \in X^1(T)} T_P$ and this intersection has finite character. Note that $XE[X]$ is a height-one prime ideal of T [Anderson et al. 2006, Theorem 3.4]; so $E[\Gamma^*]$ is a height-one prime ideal of R by Lemma 2.1. We claim that $R = \bigcap_{P \in X^1(R)} R_{P \cap R}$, where P ranges over all height-one prime ideals of T . Suppose to the contrary that there exists an element f in $\bigcap_{P \in X^1(R)} R_{P \cap R} \setminus R$. Note that $f \in T$, and hence we can write $f = \sum_{i=0}^m f_i X^i$. Then there exists a polynomial $g \in R \setminus E[\Gamma^*]$ such that $fg \in R$. Since $g(0) \neq 0$, the same argument as in the proof of Lemma 2.3(1) shows that $f \in R$, which contradicts the choice of f . Thus the equality holds. Since $T = \bigcap_{P \in X^1(T)} T_P$ has finite character, it is clear that the intersection $R = \bigcap_{P \in X^1(R)} R_{P \cap R}$ also has finite character. Thus R is a weakly Krull domain.

(2) \Rightarrow (3) This implication was already shown in the proof of (2) \Rightarrow (1).

(3) \Rightarrow (4) Assume that T_n is a weakly Krull domain. Then $t\text{-dim}(T_n) = 1$ [Anderson et al. 1992, Lemma 2.1]; so $X^n E[X]$ is a maximal t -ideal of T_n by Lemma 2.3(1).

Let $S = \{X^m \mid m \in \Delta_n\}$. Then $E[\Delta_n]_S = E[X, X^{-1}] = (T_n)_S$ is a weakly Krull domain [Anderson et al. 1993, Proposition 4.7]. Note that $XE[X]$ is a height-one prime ideal of $E[X]$; so $X^n E[X]$ is a height-one prime ideal of $E[\Delta_n]$ [Chang et al. 2012, Proposition 1.1]; so $E[\Delta_n]_{X^n E[X]}$ is a one-dimensional quasi-local domain. Hence $E[\Delta_n]_{X^n E[X]}$ is a weakly Krull domain. We claim that $E[\Delta_n] = E[\Delta_n]_S \cap E[\Delta_n]_{X^n E[X]}$. Let $f = f_0 + \sum_{i=n}^{k_1} f_i X^i$ and $h = h_0 + \sum_{i=n}^{k_2} h_i X^i$ be nonzero elements of $E[\Delta_n]$ with $h(0) \neq 0$ and let $g = \sum_{i=0}^{k_3} g_i X^i \in E[X] \setminus \{0\}$ with $g(0) \neq 0$ satisfying $\frac{g}{X^m} = \frac{f}{h} \in E[\Delta_n]_S \cap E[\Delta_n]_{X^n E[X]}$ for some nonnegative integer m . Then $X^m f = gh$; so $m = 0$. By comparing coefficients of f and gh , it is easy to see that $g_i = 0$ for all $i = 1, \dots, n-1$. Hence $\frac{g}{X^m} \in E[\Delta_n]$. The reverse inclusion is clear. Thus $E[\Delta_n]$ is a weakly Krull domain.

(4) \Rightarrow (5) By [Zafrullah 2003, Lemma 2.6], $\text{ht}_T(XE[X]) = \dim(E_{D \setminus \{0\}}[X])$. By (4), $\text{ht}_{T_n}(X^n E[X]) = 1$; so the comment before Remark 2.2 establishes that

$$\dim(E_{D \setminus \{0\}}[X]) = 1.$$

Thus $E_{D \setminus \{0\}}$ is a field. Also, since $X^n E[X]$ is a maximal t -ideal of T_n , $qf(D) \cap E = D$ by Lemma 2.3(2). Finally, it follows directly from Theorem 1.3 that $E[X]$ is a weakly Krull domain.

(5) \Rightarrow (2) [Anderson et al. 2006, Theorem 3.4].

(1) \Rightarrow (2) In the proof of (2) \Leftrightarrow (4), the integer $n \geq 2$ was arbitrary; so it suffices to show that $X^{F(\Gamma)+1} E[X]$ is a height-one maximal t -ideal of $T_{F(\Gamma)+1}$ and $E[\Delta_{F(\Gamma)+1}]$ is a weakly Krull domain. Assume that R is a weakly Krull domain. Since $t\text{-dim}(R) = 1$ [Anderson et al. 1992, Lemma 2.1], $E[\Gamma^*]$ is a height-one maximal t -ideal of R by Lemma 2.3(1); so $X^{F(\Gamma)+1} E[X]$ is a height-one maximal t -ideal of $T_{\Delta_{F(\Gamma)+1}}$ by Lemma 2.1 and the remark before Remark 2.2. Let $S_1 = \{X^\alpha \mid \alpha \in \Delta_{F(\Gamma)+1}\}$ and $S_2 = \{X^\alpha \mid \alpha \in \Gamma\}$. Then $E[\Delta_{F(\Gamma)+1}]_{S_1} = R_{S_2}$ is a weakly Krull domain [Anderson et al. 1993, Proposition 4.7]. Also, $E[\Delta_{F(\Gamma)+1}]_{X^{F(\Gamma)+1} E[X]}$ is a weakly Krull domain because it is one-dimensional quasi-local. Note that $E[\Delta_{F(\Gamma)+1}] = E[\Delta_{F(\Gamma)+1}]_{S_1} \cap E[\Delta_{F(\Gamma)+1}]_{X^{F(\Gamma)+1} E[X]}$ as in the proof of (3) \Rightarrow (4). Thus $E[\Delta_{F(\Gamma)+1}]$ is a weakly Krull domain. \square

Corollary 2.5. *Let $R = D + E[\Gamma^*]$, $T = D + XE[X]$, $T_n = D + X^n E[X]$ and $\Delta_n = \{0\} \cup \{m \in \mathbb{N}_0 \mid m \geq n\}$ for integers $n \geq 2$. If $\text{char}(E) \neq 0$, then the following statements are equivalent.*

- (1) R is an AWFD.
- (2) R is a GWFD.
- (3) T is an AWFD.

- (4) T is a GWFD.
- (5) T_n is an AWFD.
- (6) T_n is a GWFD.
- (7) $X^n E[X]$ is a maximal t -ideal of T_n , $E[\Delta_n]$ is an AWFD and for each $0 \neq e \in E$, there exist an integer $m = m(e) \geq 1$ and a unit u of E such that $ue^m \in D$.
- (8) $X^n E[X]$ is a maximal t -ideal of T_n , $E[\Delta_n]$ is a GWFD and for each $0 \neq e \in E$, there exist an integer $m = m(e) \geq 1$ and a unit u of E such that $ue^m \in D$.
- (9) $qf(D) \cap E = D$, $E[X]$ is an AWFD and for each $0 \neq e \in E$, there exist an integer $m = m(e) \geq 1$ and a unit u of E such that $ue^m \in D$.
- (10) $qf(D) \cap E = D$, $E[X]$ is a GWFD and for each $0 \neq e \in E$, there exist an integer $m = m(e) \geq 1$ and a unit u of E such that $ue^m \in D$.

Proof. (1) \Rightarrow (2) and (5) \Rightarrow (6) Their definitions lead to these implications.

(3) \Leftrightarrow (9) [Anderson et al. 2006, Theorem 3.5].

(4) \Leftrightarrow (10) [Anderson and Chang 2007, Corollary 2.10].

(7) \Leftrightarrow (8) and (9) \Leftrightarrow (10) See Corollary 1.5.

(7) \Leftrightarrow (9) This equivalence follows from Corollary 1.5 and Lemma 2.3(2).

(3) \Rightarrow (1) Assume that T is an AWFD. Then T is a weakly Krull domain [Anderson et al. 1992, Theorem 3.4]. Hence $E[X]$ is a weakly Krull domain by Theorem 2.4. Let $S = \{X^m \mid m \in \mathbb{N}_0\}$. Since X is a prime element of $E[X]$, $\text{Cl}(E[X]) = \text{Cl}(T_S)$ is torsion [Anderson et al. 1993, Corollary 4.9]; so $E[X]$ is an AWFD [Anderson et al. 1992, Theorem 3.4]. Let $f \in R \setminus \{0\}$. Then there exists an integer $m \geq 1$ such that $f^m = X^l f_1 \cdots f_r$ for some nonnegative positive integer l and primary elements f_1, \dots, f_r of $E[X]$ with nonzero constant terms. Also, since $\text{char}(E) \neq 0$, there exists an integer $k \geq F(\Gamma) + 1$ such that $f_i^k \in E[\Gamma]$ for all $i = 1, \dots, r$; so $f^{mk} = X^{lk} f_1^k \cdots f_r^k \in E[\Gamma]$. Fix an $i \in \{1, \dots, r\}$, and we claim that $\sqrt{f_i^k E[\Gamma]}$ is a prime ideal of $E[\Gamma]$ [Anderson et al. 2003b, Lemma 2.1]. Note that $\sqrt{f_i^k E[X]} = \sqrt{f_i^k E[X]}$. If $\sqrt{f_i^k E[X]} = XE[X]$, then an easy calculation using a similar method as in the proof of (2) \Rightarrow (1) in Theorem 2.4 shows that $\sqrt{f_i^k E[\Gamma]} = E[\Gamma^*]$ is a prime ideal. Assume that $\sqrt{f_i^k E[X]} \neq XE[X]$. Since $f_i(0) \neq 0$, $f_i^k E[X, X^{-1}]$ is a primary ideal of $E[X, X^{-1}]$; so $f_i^k E[X, X^{-1}] \cap E[\Gamma]$ is primary in $E[\Gamma]$. It is easy to see that $\sqrt{f_i^k E[X, X^{-1}] \cap E[\Gamma]} = \sqrt{f_i^k E[\Gamma]}$. Hence $\sqrt{f_i^k E[\Gamma]}$ is a prime ideal. Therefore we may assume that f_1, \dots, f_r are primary elements of $E[\Gamma]$ with nonzero constant terms and write $f^m = X^l f_1 \cdots f_r$ as above. Note that for each $i = 1, \dots, r$, there exist a unit u_i of E and an integer $a_i \geq F(\Gamma) + 1$ such that

$u_i f_i(0)^{a_i} \in D$ as in the proof of (3) \Leftrightarrow (9); so $u_i f_i^{a_i} \in R$. Let

$$a = a_1 \cdots a_r, \quad \hat{a}_i = \frac{a}{a_i}, \quad \text{and} \quad u = u_1^{\hat{a}_1} \cdots u_r^{\hat{a}_r}.$$

Then $u f^{am} = X^{al} (u_1 f_1^{a_1})^{\hat{a}_1} \cdots (u_r f_r^{a_r})^{\hat{a}_r}$ and $\sqrt{(u_i f_i^{a_i})^{\hat{a}_i} E[\Gamma]} = \sqrt{f_i E[\Gamma]}$ for each $i = 1, \dots, r$. Since $t\text{-dim}(E[\Gamma]) = 1$, $(u_i f_i^{a_i})^{\hat{a}_i} E[\Gamma]$ is a primary ideal of $E[\Gamma]$ [Anderson et al. 2003b, Lemma 2.1] for each $1 \leq i \leq r$.

Claim. For each $1 \leq i \leq r$, $(u_i f_i^{a_i})^{\hat{a}_i} R$ is a primary ideal of R .

Proof. Note that $(u_i f_i^{a_i})^{\hat{a}_i} \in R$ and fix an $i \in \{1, \dots, r\}$. We also note that $t\text{-dim}(R) = 1$ because R is a weakly Krull domain by Theorem 2.4. Hence, by [Anderson et al. 2003b, Lemma 2.1], it suffices to show that $\sqrt{(u_i f_i^{a_i})^{\hat{a}_i} R}$ is a prime ideal of R . If $\sqrt{(u_i f_i^{a_i})^{\hat{a}_i} E[\Gamma]} = E[\Gamma^*]$, then it is easy to see that $\sqrt{(u_i f_i^{a_i})^{\hat{a}_i} R} = E[\Gamma^*]$ is a prime ideal of R . Assume that $\sqrt{(u_i f_i^{a_i})^{\hat{a}_i} E[\Gamma]} \neq E[\Gamma^*]$. Then $(u_i f_i(0)^{a_i})^{\hat{a}_i} \neq 0$. Now, we show that $(u_i f_i^{a_i})^{\hat{a}_i} E[X, X^{-1}] \cap R = (u_i f_i^{a_i})^{\hat{a}_i} R$. Let $h \in (u_i f_i^{a_i})^{\hat{a}_i} E[X, X^{-1}] \cap R$. Note that we have

$$\begin{aligned} (u_i f_i^{a_i})^{\hat{a}_i} E[X, X^{-1}] \cap R &\subseteq (u_i f_i^{a_i})^{\hat{a}_i} E[X, X^{-1}] \cap E[\Gamma] \\ &= (u_i f_i^{a_i})^{\hat{a}_i} E[\Gamma] \end{aligned}$$

by adapting the proof of (2) \Rightarrow (1) in Theorem 2.4. So, we can write $h = (u_i f_i^{a_i})^{\hat{a}_i} g$ for some $g \in E[\Gamma]$. Then

$$g(0) = \frac{(u_i f_i(0)^{a_i})^{\hat{a}_i}}{h(0)} \in qf(D) \cap E = D$$

by Theorem 2.4; so $g \in R$. Therefore $h \in (u_i f_i^{a_i})^{\hat{a}_i} R$, and hence

$$(u_i f_i^{a_i})^{\hat{a}_i} E[X, X^{-1}] \cap R \subseteq (u_i f_i^{a_i})^{\hat{a}_i} R.$$

The reverse inclusion is clear, and hence $(u_i f_i^{a_i})^{\hat{a}_i} E[X, X^{-1}] \cap R = (u_i f_i^{a_i})^{\hat{a}_i} R$. Since $(u_i f_i^{a_i})^{\hat{a}_i} E[\Gamma]$ is a primary ideal of $E[\Gamma]$, $(u_i f_i^{a_i})^{\hat{a}_i} E[X, X^{-1}]$ is a primary ideal of $E[X, X^{-1}]$. Therefore $\sqrt{(u_i f_i^{a_i})^{\hat{a}_i} R} = \sqrt{(u_i f_i^{a_i})^{\hat{a}_i} E[X, X^{-1}] \cap R}$ is a prime ideal of R , and thus $(u_i f_i^{a_i})^{\hat{a}_i} R$ is a primary ideal of R . The claim is proved. \square

If $l = 0$, then $u f(0)^{am} = (u_1 f_1(0)^{a_1})^{\hat{a}_1} \cdots (u_r f_r(0)^{a_r})^{\hat{a}_r} \in D$; so u is a unit of D because u is a unit of E . If $l \geq 1$, then $f^{am} = u^{-1} X^{al} (u_1 f_1^{a_1})^{\hat{a}_1} \cdots (u_r f_r^{a_r})^{\hat{a}_r}$. Since $u^{-1} X^{al} E[\Gamma]$ is a primary ideal of $E[\Gamma]$, $u^{-1} X^{al} R$ is a primary ideal of R by imitating the previous proof. Hence f^{am} is a product of primary elements of R , and thus R is an AWFD.

(2) \Rightarrow (8) Assume that R is a GWFD and fix an integer $n \geq 2$. Then R is a weakly Krull domain [Anderson et al. 2003b, Corollary 2.3]; so $X^n E[X]$ is a height-one maximal t -ideal of T_n by Theorem 2.4.

Next, we claim that $E[\Delta_n]$ is a GWFD. Let $S_1 = \{X^m \mid m \in \Delta_n\}$ and $S_2 = \{X^m \mid m \in \Gamma\}$. Then $E[\Delta_n]_{S_1} = E[X, X^{-1}] = R_{S_2}$ is a GWFD. Let Q be a nonzero prime ideal of $E[\Delta_n]$. If $Q \cap S_1 \neq \emptyset$, then Q contains a primary element X^n of $E[\Delta_n]$. If $Q \cap S_1 = \emptyset$, then $QE[\Delta_n]_{S_1}$ is a prime ideal of $E[\Delta_n]_{S_1}$; so $QE[\Delta_n]_{S_1}$ contains a primary element $f \in E[X, X^{-1}]$. Note that X is a unit of $E[X, X^{-1}]$ and $f^k \in E[\Delta_n]$ for some integer $k \geq 1$ because $\text{char}(E) \neq 0$; so we may assume that $f \in E[\Delta_n]$ with $f(0) \neq 0$. Then

$$fE[\Delta_n] \subseteq fE[\Delta_n]_{S_1} \cap E[\Delta_n] \subseteq QE[\Delta_n]_{S_1} \cap E[\Delta_n] = Q;$$

so Q contains a primary element f . Hence $E[\Delta_n]$ is a GWFD.

In order to check the final condition, let $e \in E \setminus \{0\}$. If e is a unit of E , then we have nothing to prove. So, we assume that e is not a unit of E and let $h = e + X \in E[X]$. Since $c(h)_v = E$, $hE[X] = hqf(E)[X] \cap E[X]$ [Anderson and Chang 2007, Lemma 2.1(1)]; so $hE[X]$ is a height-one prime ideal. Let $P = hE[X] \cap R$. Since e is not a unit of E , $X^{F(\Gamma)+1} \notin P$; so $X^\alpha \notin P$ for all $\alpha \in \Gamma$. Therefore $hE[X, X^{-1}] = PR_{S_2} \subsetneq R_{S_2}$, and hence $\text{ht}_R(P) = 1$. Since R is a GWFD, $P = \sqrt{gR}$ for some primary element $g \in R$ [Anderson et al. 2003b, Theorem 2.2]. Suppose to the contrary that $g(0) = 0$. Since $E_{D \setminus \{0\}}$ is a field by Theorem 2.4, $\frac{1}{e} = \frac{e'}{d}$ for some $0 \neq d \in D$ and $e' \in E$; so $e'h = d + e'X \in T$. Since $\text{char}(E) \neq 0$, $(e'h)^k \in hE[X] \cap R = P$ for some integer $k \geq 1$. Hence $(e'h)^{kl} \in gR$ for some integer $l \geq 1$. However, this is impossible because $e \neq 0$. Therefore $g(0) \neq 0$. It is clear that gR_{S_2} is a primary ideal of R_{S_2} , $gR_{S_2} \cap E[X] = gE[X]$, $PR_{S_2} = \sqrt{gR_{S_2}}$ and $PR_{S_2} \cap E[X] = hE[X]$. Hence $gE[X]$ is a $hE[X]$ -primary ideal. Therefore $g = uh^m$ for some $u \in qf(E)$ and some integer $m \geq 1$; so $ue^m = g(0) \in D$. Thus u is a unit of E .

(3) \Rightarrow (5) and (6) \Rightarrow (8) These implications can be obtained by applying $\Gamma = \Delta_n$ to the proofs of (3) \Rightarrow (1) and (2) \Rightarrow (8), respectively. \square

We are closing this paper by showing that $R = D + E[\Gamma^*]$ is never a WFD and the assumption “ $\text{char}(E) = 0$ ” is essential in Corollary 2.5.

Remark 2.6. Assume that $R = D + E[\Gamma^*]$ is a WFD or an AWFD. Let $h = 1 + X \in E[X]$, $P = hE[X] \cap R$ and let M be a maximal t -ideal of R . If $M = E[\Gamma^*]$, then $PR_M = R_M$ because $1 + (-1)^{F(\Gamma)} X^{F(\Gamma)+1} \in P \setminus E[\Gamma^*]$. Assume that $M \neq E[\Gamma^*]$. Since $c(h)_v = E$, $hqf(E)[X] \cap E[X] = hE[X]$ [Anderson and Chang 2007, Lemma 2.1(1)]. Let $S = \{X^m \mid m \in \Gamma\}$. Then $PE[X, X^{-1}] = hE[X, X^{-1}]$; so $PR_M = hR_M$ is principal. Hence P is t -locally principal, and thus P is t -invertible [Anderson et al. 1992, Lemma 2.2].

(1) If R is a WFD, then $P = gR$ for some $g \in R$ with $g(0) \neq 0$ [Anderson and Zafrullah 1990, Theorem]. Note that $hE[X, X^{-1}] = gE[X, X^{-1}]$; so $g = uh$ for some unit u of E . Hence $uh \in R$, which is impossible. Thus R is not a WFD.

(2) Assume that R is an AWFD. Then $P^m = gR$ for some integer $m \geq 1$ and $g \in R$ with $g(0) \neq 0$ [Anderson et al. 1992, Theorem 3.4]. We note that

$$h^m E[X, X^{-1}] = gE[X, X^{-1}];$$

so $uh^m = g$ for some unit u of E . Hence $uh^m \in R$. However, this can not happen if $\text{char}(E) = 0$. Thus R is never an AWFD whenever $\text{char}(E) = 0$.

Acknowledgements

The author would like to thank the referee for several valuable comments and suggestions which resulted in an improved version of the paper. Also, the author is deeply grateful to professor Gyu Whan Chang for pointing out Remark 1.6(1).

This work was supported by the Brain Korea 21 Project Team to Nurture the Next Generation of First-class Mathematical Scientists by the Korean Government.

References

- [Anderson and Chang 2004] D. F. Anderson and G. W. Chang, “The class group of $D[\Gamma]$ for D an integral domain and Γ a numerical semigroup”, *Comm. Algebra* **32**:2 (2004), 787–792. MR 2005g:13021 Zbl 1092.13017
- [Anderson and Chang 2005] D. F. Anderson and G. W. Chang, “Homogeneous splitting sets of a graded integral domain”, *J. Algebra* **288**:2 (2005), 527–544. MR 2006g:13001 Zbl 1084.13001
- [Anderson and Chang 2007] D. F. Anderson and G. W. Chang, “Almost splitting sets in integral domains. II”, *J. Pure Appl. Algebra* **208**:1 (2007), 351–359. MR 2007i:13004 Zbl 1171.13300
- [Anderson and Zafrullah 1990] D. D. Anderson and M. Zafrullah, “Weakly factorial domains and groups of divisibility”, *Proc. Amer. Math. Soc.* **109**:4 (1990), 907–913. MR 90k:13015 Zbl 0704.13008
- [Anderson et al. 1992] D. D. Anderson, J. L. Mott, and M. Zafrullah, “Finite character representations for integral domains”, *Boll. Un. Mat. Ital. B (7)* **6**:3 (1992), 613–630. MR 93k:13001 Zbl 0773.13004
- [Anderson et al. 1993] D. D. Anderson, E. G. Houston, and M. Zafrullah, “ t -linked extensions, the t -class group, and Nagata’s theorem”, *J. Pure Appl. Algebra* **86**:2 (1993), 109–124. MR 94e:13036 Zbl 0777.13002
- [Anderson et al. 1995] D. D. Anderson, D. F. Anderson, and M. Zafrullah, “A generalization of unique factorization”, *Boll. Un. Mat. Ital. A (7)* **9**:2 (1995), 401–413. MR 96d:13027 Zbl 0919.13001
- [Anderson et al. 2003a] D. F. Anderson, G. W. Chang, and J. Park, “ $D[X^2, X^3]$ over an integral domain D ”, pp. 1–14 in *Commutative ring theory and applications* (Fez, 2001), Lecture Notes in Pure and Appl. Math. **231**, Dekker, New York, 2003. MR 2004k:13011 Zbl 1080.13510
- [Anderson et al. 2003b] D. F. Anderson, G. W. Chang, and J. Park, “Generalized weakly factorial domains”, *Houston J. Math.* **29**:1 (2003), 1–13. MR 2004e:13003 Zbl 1029.13012
- [Anderson et al. 2006] D. F. Anderson, G. W. Chang, and J. Park, “Weakly Krull and related domains of the form $D + M$, $A + XB[X]$ and $A + X^2B[X]$ ”, *Rocky Mountain J. Math.* **36**:1 (2006), 1–22. MR 2007h:13027 Zbl 1133.13022

- [Chang 2005] G. W. Chang, “Almost splitting sets in integral domains”, *J. Pure Appl. Algebra* **197**:1-3 (2005), 279–292. MR 2005j:13005 Zbl 1091.13001
- [Chang et al. 2012] G. W. Chang, H. Kim, and J. W. Lim, “Numerical semigroup rings and almost Prüfer v -multiplication domains”, preprint, 2012. Accepted in *Comm. Algebra*.
- [El Baghdadi et al. 2002] S. El Baghdadi, L. Izelgue, and S. Kabbaj, “On the class group of a graded domain”, *J. Pure Appl. Algebra* **171**:2-3 (2002), 171–184. MR 2003d:13012 Zbl 1058.13006
- [Gilmer 1992] R. Gilmer, *Multiplicative ideal theory*, Queen’s Papers in Pure and Applied Mathematics **90**, Queen’s University, Kingston, ON, 1992. MR 93j:13001 Zbl 0804.13001
- [Gilmer and Martin 1990] R. Gilmer and M. B. Martin, “On the Picard group of a class of nonseminormal domains”, *Comm. Algebra* **18**:10 (1990), 3263–3293. MR 91g:13017
- [Kang 1989] B. G. Kang, “Prüfer v -multiplication domains and the ring $R[X]_{N_v}$ ”, *J. Algebra* **123**:1 (1989), 151–170. MR 90e:13017 Zbl 0668.13002
- [Kaplansky 1970] I. Kaplansky, *Commutative rings*, Allyn and Bacon, Boston, 1970. Reprinted Polygonal Publishing House, Washington, 1994.
- [Li 2012] Q. Li, “On almost Prüfer v -multiplication domains”, *Algebra Colloq.* (2012).
- [Lim 2012] J. W. Lim, “Generalized Krull domains and the composite semigroup ring $D + E[\Gamma^*]$ ”, *J. Algebra* **357** (2012), 20–25.
- [Zafrullah 2003] M. Zafrullah, “Various facets of rings between $D[X]$ and $K[X]$ ”, *Comm. Algebra* **31**:5 (2003), 2497–2540. MR 2004d:13029 Zbl 1052.13003

Received September 17, 2011. Revised April 10, 2012.

JUNG WOOK LIM
DEPARTMENT OF MATHEMATICS
SOGANG UNIVERSITY
SEOUL 121-742
SOUTH KOREA
lovemath@postech.ac.kr
jwlim@sogang.ac.kr

ARITHMETICITY OF COMPLEX HYPERBOLIC TRIANGLE GROUPS

MATTHEW STOVER

Complex hyperbolic triangle groups, originally studied by Mostow in building the first nonarithmetic lattices in $\mathrm{PU}(2, 1)$, are a natural generalization of the classical triangle groups. A theorem of Takeuchi states that there are only finitely many Fuchsian triangle groups that determine an arithmetic lattice in $\mathrm{PSL}_2(\mathbb{R})$, so triangle groups are generically nonarithmetic. We prove similar finiteness theorems for complex hyperbolic triangle groups that determine an arithmetic lattice in $\mathrm{PU}(2, 1)$.

1. Introduction

In a seminal paper [1980], Mostow constructed lattices in $\mathrm{PU}(2, 1)$ generated by three complex reflections. He not only gave a new geometric method for building lattices acting on the complex hyperbolic plane, but gave the first examples of nonarithmetic lattices in $\mathrm{PU}(2, 1)$. Complex reflection groups are a generalization of groups generated by reflections through hyperplanes in constant curvature spaces, and Mostow's groups are a natural extension to the complex hyperbolic plane of the classical triangle groups. They are often called *complex hyperbolic triangle groups*. We introduce these groups in Section 2. See also [Goldman and Parker 1992; Schwartz 2002], which, along with [Mostow 1980], inspired much of the recent surge of activity surrounding these groups.

Around the same time, Takeuchi [1977] classified the Fuchsian triangle groups that determine arithmetic lattices in $\mathrm{PSL}_2(\mathbb{R})$. In particular, he proved that there are finitely many and gave a complete list. Since there are infinitely many triangle groups acting on the hyperbolic plane discretely with finite covolume, triangle groups are generically nonarithmetic. The purpose of this paper is to give analogous finiteness results for complex hyperbolic triangle groups that determine an arithmetic lattice in $\mathrm{PU}(2, 1)$.

A particular difficulty with complex hyperbolic triangle groups is that the complex triangle is not uniquely determined by its angles. One must also consider the

Partially supported by NSF RTG grant DMS 0602191.

MSC2010: 11F06, 20H10, 22E40.

Keywords: complex hyperbolic geometry, arithmetic lattices, complex hyperbolic triangle groups.

so-called *angular invariant* $\psi \in [0, 2\pi)$. See Section 2. In particular, there is a 1-dimensional deformation space of complex triangles with fixed triple of angles. The typical assumption is that ψ is a rational multiple of π , in which case the angular invariant is called *rational*. We call it *irrational* otherwise.

When a complex hyperbolic triangle group is also an arithmetic lattice, we will call it an arithmetic complex hyperbolic triangle group. Note that this immediately implies discreteness. Our first result is for nonuniform arithmetic complex hyperbolic triangle groups. We prove the following in Section 6.

Theorem 1.1. *There are finitely many nonuniform arithmetic complex hyperbolic triangle groups with rational angular invariant. If Γ is a nonuniform arithmetic complex hyperbolic triangle group with irrational angular invariant ψ , then $e^{i\psi}$ is contained in a biquadratic extension of \mathbb{Q} .*

We emphasize that complex reflection groups are allowed to have generators of arbitrary finite order. A usual assumption is that all generators have the same order, a restriction that we avoid. See Theorem 6.1 for a more precise formulation of Theorem 1.1. Proving that a candidate is indeed a lattice is remarkably difficult, as evidenced in [Mostow 1980; Deraux et al. 2011], so we do not give a definitive list. One consequence of the proof (see Theorem 1.5(1) below) is the following.

Corollary 1.2. *Suppose that Γ is a nonuniform lattice in $U(2, 1)$. If Γ contains a complex reflection of order 5 or at least 7, then Γ is nonarithmetic.*

In the cocompact setting, the arithmetic is much more complicated. Arithmetic subgroups of $U(2, 1)$ come in two types, defined in Section 3, often called first and second. In Section 4 we prove the following auxiliary result, generalizing a well-known fact for hyperbolic reflection groups.

Theorem 1.3. *Let $\Gamma < U(2, 1)$ be a lattice containing a complex reflection. Then Γ contains a Fuchsian subgroup stabilizing the wall of the reflection in $\mathbf{H}_{\mathbb{C}}^2$.*

We also give a generalization to higher-dimensional complex reflection groups. Theorem 1.3 leads to the following, which we also prove in Section 4.

Theorem 1.4. *Let $\Gamma < U(2, 1)$ be a lattice, and suppose that Γ is commensurable with a lattice Λ containing a complex reflection. Then Γ is either arithmetic of first type or nonarithmetic.*

In particular, when considering a complex reflection group as a candidate for a nonarithmetic lattice, one must only show that it is not of the first type. Fortunately, this is the case where the arithmetic is simplest to understand.

The effect of the angular invariant is a particular sticking point in generalizing Takeuchi's methods. In Section 5, the technical heart of the paper, we study the interdependence between the geometric invariants of the triangle and the arithmetic

of the lattice. We collect the most useful of these facts as the following. See §§2-3 for our notation.

Theorem 1.5. *Suppose that Γ is an arithmetic complex hyperbolic triangle group. Suppose that for $j = 1, 2, 3$ the generators have reflection factors η_j , the complex angles of the triangle are θ_j , and that the angular invariant is ψ . Let E be the totally imaginary quadratic extension of the totally real field F that defines Γ as an arithmetic lattice. Then:*

- (1) $\eta_j \in E$ for all j ;
- (2) $\cos^2 \theta_j \in F$ for all j ;
- (3) $e^{2i\psi} \in E$ and $\cos^2 \psi \in F$;
- (4) If $\theta_j \leq \pi/3$ for all j , then

$$\cos^2 \psi \in \mathbb{Q}(\cos^2 \theta_1, \cos^2 \theta_2, \cos^2 \theta_3, \cos \theta_1 \cos \theta_2 \cos \theta_3);$$

- (5) $E \subseteq \mathbb{Q}(\cos^2 \theta_1, \cos^2 \theta_2, \cos^2 \theta_3, e^{i\psi} \cos \theta_1 \cos \theta_2 \cos \theta_3)$;
- (6) If ψ is rational, then E is a subfield of a cyclotomic field.

In Section 6, we use the results from Section 5 to prove finiteness results for cocompact arithmetic complex hyperbolic triangle groups with rational angular invariant. We also give restrictions for irrational angular invariants, though it is unknown whether such a lattice exists. When the complex triangle is a right triangle, we prove the following.

Theorem 1.6. *Suppose that Γ is an arithmetic complex hyperbolic triangle group for which the associated complex triangle is a right triangle. Then the angles of the triangle are the angles of an arithmetic Fuchsian triangle group. There are finitely many such Γ with rational angular invariant.*

Finally, we consider equilateral triangles at the end of Section 6. This is the case which has received the most attention, in particular from Mostow [1980] and, in the ideal case, by Goldman and Parker [1992] and Schwartz [2002]. See also [Deraux 2006]. Here we cannot explicitly bound orders of generators, angles, or angular invariants because our proof relies on asymptotic number theory for which we do not know precise constants. Nevertheless, we obtain finiteness in the situation that has received the greatest amount of attention since Mostow's original paper. See [Parker 2008; Parker and Paupert 2009; Paupert 2010; Deraux et al. 2011] and references therein for more recent examples of lattices and restrictions on discreteness.

Theorem 1.7. *There are finitely many arithmetic complex hyperbolic equilateral triangle groups with rational angular invariant.*

2. Complex hyperbolic triangle groups

We assume some basic knowledge of complex hyperbolic geometry, e.g., the first three chapters of [Goldman 1999]. Let V be a three-dimensional complex vector space, equipped with a hermitian form h of signature $(2, 1)$. Complex hyperbolic space $\mathbf{H}_{\mathbb{C}}^2$ is the space of h -negative lines in V . The metric on $\mathbf{H}_{\mathbb{C}}^2$ is defined via h as in [Goldman 1999, Chapter 3], and the action of $U(2, 1)$ on $\mathbf{H}_{\mathbb{C}}^2$ by isometries descends from its action on V and factors through projection onto $PU(2, 1)$. Its ideal boundary $\partial\mathbf{H}_{\mathbb{C}}^2$ is the space of h -isotropic lines, and we set $\overline{\mathbf{H}}_{\mathbb{C}}^2 = \mathbf{H}_{\mathbb{C}}^2 \cup \partial\mathbf{H}_{\mathbb{C}}^2$.

A *complex reflection* is a diagonalizable linear map $R : V \rightarrow V$ with one eigenvalue of multiplicity 2 (or, more generally, multiplicity $n - 1$ when $\dim(V) = n$). We assume that R has finite order, so the third eigenvalue is a root of unity η . We call η the *reflection factor* of R . Decompose $V = V_1 \oplus V_{\eta}$ into the 1- and η -eigenspaces, and choose $v_{\eta} \in V$ so that $V_{\eta} = \text{Span}_{\mathbb{C}}\{v_{\eta}\}$. We begin with an elementary lemma that will be of use later, keeping in mind that every complex reflection has 1 as an eigenvalue.

Lemma 2.1. *Let $A \in GL_n(\mathbb{C})$ be a diagonalizable linear transformation. Let $E \subseteq \mathbb{C}$ be a subfield, and suppose that E^n has a basis consisting of eigenvectors for A . Furthermore, suppose that A has at least one eigenvalue in E and that there exists $x \in \mathbb{C}^{\times}$ so that $x A \in GL_n(E)$. Then all eigenvalues of A are in E .*

Proof. Let $v_1, \dots, v_n \in E^n$ be a basis of eigenvectors for A , and let λ_j be the eigenvalue associated with v_j , $1 \leq j \leq n$. Without loss of generality, $\lambda_1 \in E$. Then $x A$ also has eigenvectors v_1, \dots, v_n , and $x A v_j = x \lambda_j v_j \in E^n$ for all j , since $x A \in GL_n(E)$. Then $x \lambda_j \in E$, $1 \leq j \leq n$. Since $\lambda_1 \in E$, it follows that $x \in E$, which implies that $\lambda_j \in E$ for all j . □

Assume that $R \in U(2, 1)$. Then the fixed point set of R acting on $\mathbf{H}_{\mathbb{C}}^2$ is the subset of h -negative lines in V_1 . This is a totally geodesic holomorphic embedding of the hyperbolic plane if and only if V_{η} is an h -positive line. These subspaces are called *complex hyperbolic lines*. Following [Goldman 1999, §3.1], we call v_{η} a *polar vector* for R .

When V_{η} is h -negative, the fixed set of R on $\mathbf{H}_{\mathbb{C}}^2$ is a point, and R is sometimes called a reflection through that point. The complex reflections in this paper will always be through complex hyperbolic lines. That is, the η -eigenspace will always be an h -positive line.

Let W be the complex hyperbolic line in $\mathbf{H}_{\mathbb{C}}^2$ fixed by R . We call this the *wall* of R . If v_{η} is a polar vector, then R is the linear transformation

$$(1) \quad z \mapsto z + (\eta - 1) \frac{h(z, v_{\eta})}{h(v_{\eta}, v_{\eta})} v_{\eta}.$$

We refrain from normalizing the polar vector to have h -norm one, since we will often choose a polar vector with coordinates in a subfield E of \mathbb{C} , and $E^3 \subset V$ might not contain an h -norm one representative for the given line of polar vectors.

Now, consider three complex reflections $R_1, R_2, R_3 \in U(2, 1)$ with respective distinct walls W_1, W_2, W_3 in $\mathbf{H}_{\mathbb{C}}^2$. If v_j is a polar vector for R_j , then W_j and W_{j+1} (with cyclic indices) meet in $\mathbf{H}_{\mathbb{C}}^2$ if and only if

$$(2) \quad h(W_j, W_{j+1}) = \frac{|h(v_j, v_{j+1})|^2}{h(v_j, v_j)h(v_{j+1}, v_{j+1})} < 1.$$

The two walls meet at a point z_j stabilized by the subgroup of $U(2, 1)$ generated by R_j and R_{j+1} . The *complex angle* θ_j between W_j and W_{j+1} , the minimum angle between the two walls, satisfies $\cos^2 \theta_j = h(W_j, W_{j+1})$.

The walls W_j and W_{j+1} meet at a point p_j in $\partial\mathbf{H}_{\mathbb{C}}^2$ if and only if

$$(3) \quad \frac{|h(v_j, v_{j+1})|^2}{h(v_j, v_j)h(v_{j+1}, v_{j+1})} = 1,$$

so we say that the complex angle is zero. The group generated by R_j and R_{j+1} fixes p_j , so it is contained in a parabolic subgroup of $U(2, 1)$. See [Goldman 1999, §3.3.2].

Let $\{R_j\}$ be reflections through walls $\{W_j\}$, $j = 1, 2, 3$. When the pairwise intersections of the walls are nontrivial in $\overline{\mathbf{H}_{\mathbb{C}}^2}$, they determine a *complex triangle* in $\mathbf{H}_{\mathbb{C}}^2$, possibly with ideal vertices. The subgroup $\Delta(R_1, R_2, R_3)$ of $U(2, 1)$ generated by the R_j s is called a *complex hyperbolic triangle group*.

A complex hyperbolic triangle group is sometimes defined as one with order two generators, and groups with higher order generators are called *generalized triangle groups*. We avoid this distinction and do not make the usual assumption that all generators have the same order.

Unlike Fuchsian triangle groups, the complex angles $\{\theta_1, \theta_2, \theta_3\}$ do not suffice to determine $\Delta(R_1, R_2, R_3)$ up to $\text{Isom}(\mathbf{H}_{\mathbb{C}}^2)$ -equivalence. We also need to consider Cartan’s *angular invariant*

$$(4) \quad \psi = \arg(h(v_1, v_2)h(v_2, v_3)h(v_3, v_1)).$$

A complex triangle is uniquely determined up to complex hyperbolic isometry by the complex angles between the walls and the angular invariant. See [Brehm 1990] and [Pratoussevitch 2005, Proposition 1]. Up to the action of complex conjugation on $\mathbf{H}_{\mathbb{C}}^2$, we can assume $\psi \in [0, \pi]$.

We call the angular invariant *rational* if $\psi = s\pi/t$ for some (relatively prime) $s, t \in \mathbb{Z}$. In other words, the angular invariant is rational if and only if $e^{i\psi}$ is a root of unity.

Let $\Delta(R_1, R_2, R_3)$ be a complex hyperbolic triangle group in $U(2, 1)$ with reflection factors η_j , complex angles θ_j , polar vectors v_j , $j = 1, 2, 3$, and angular invariant ψ . Suppose that $\{v_1, v_2, v_3\}$ is a basis for V . Then $\Delta(R_1, R_2, R_3)$ preserves the hermitian form

$$(5) \quad h_{\Delta(R_1, R_2, R_3)} = \begin{pmatrix} 1 & e^{i\psi} \cos \theta_1 & e^{i\psi} \cos \theta_3 \\ e^{-i\psi} \cos \theta_1 & 1 & e^{i\psi} \cos \theta_2 \\ e^{-i\psi} \cos \theta_3 & e^{-i\psi} \cos \theta_2 & 1 \end{pmatrix}.$$

We denote this by h_Δ when the generators are clear.

3. Arithmetic subgroups of $U(2, 1)$

Let F be a totally real number field, E a totally imaginary quadratic extension, and \mathcal{D} a central simple E -algebra of degree d . Let $\tau : \mathcal{D} \rightarrow \mathcal{D}$ be an involution, that is, an antiautomorphism of order two. Then τ is of *second kind* if $\tau|_E$ is the Galois involution of E/F . There are two cases of interest.

- (1) If $\mathcal{D} = E$ (i.e., $d = 1$), then τ is the Galois involution.
- (2) If $d = 3$, then \mathcal{D} is a cubic division algebra with center E .

See [Knus et al. 1998] for more on algebras with involution.

For d as above, let $r = 3/d$. A form $h : \mathcal{D}^r \rightarrow \mathcal{D}$ is called *hermitian* or τ -*hermitian* if it satisfies the usual definition of a hermitian form with τ in place of complex conjugation. If $d = 1$, then h is a hermitian form on E^3 as usual. If $d = 3$, then there exists an element $x \in \mathcal{D}^*$ such that $\tau(x) = x$ and $h(y_1, y_2) = \tau(y_1)xy_2$ for all $y_1, y_2 \in \mathcal{D}$.

This determines an algebraic group \mathcal{G} , the group of elements in $GL_r(\mathcal{D})$ preserving h . For every embedding $\iota : F \rightarrow \mathbb{R}$, we obtain an embedding of \mathcal{G} into the real Lie group $U(\iota(h))$. Let $\overline{\mathcal{G}}$ be the associated projective unitary group.

If \mathcal{O} is a order in \mathcal{D}^r , then the subgroup $\Gamma_{\mathcal{O}}$ of $GL_r(\mathcal{O})$ preserving h embeds as a discrete subgroup of

$$\mathcal{G}(\mathbb{R}) = \prod_{\iota: F \rightarrow \mathbb{R}} U(\iota(h)).$$

If $\overline{\Gamma}_{\mathcal{O}}$ is the image of $\Gamma_{\mathcal{O}}$ in $\overline{\mathcal{G}}$, then $\overline{\Gamma}_{\mathcal{O}}$ is a discrete subgroup of the associated product of projective unitary groups.

The projection of $\Gamma_{\mathcal{O}}$ onto any factor of $\mathcal{G}(\mathbb{R})$ is discrete if and only if the kernel of the projection of $\mathcal{G}(\mathbb{R})$ onto that factor is compact. Therefore, we obtain a discrete subgroup of $U(2, 1)$ if and only if $U(\iota(h))$ is noncompact for exactly one real embedding of F .

Then $\overline{\Gamma}_{\mathcal{O}}$ is a lattice in $PU(2, 1)$ by the well-known theorem of Borel and Harish-Chandra. An arithmetic lattice in $PU(2, 1)$ is any lattice $\Gamma < PU(2, 1)$ which is commensurable with $\overline{\Gamma}_{\mathcal{O}}$ for some \mathcal{G} as above and an order \mathcal{O} in \mathcal{D} .

Since arithmeticity only requires commensurability with $\Gamma_{\mathbb{C}}$, studying an arbitrary Γ in the commensurability class of $\Gamma_{\mathbb{C}}$ requires great care. The image of any element $\gamma \in \Gamma$ in $\text{PU}(2, 1)$ does, however, have a representative in $\text{GL}_3(E)$, that is, there exists $x \in \mathbb{C}^\times$ so $x\gamma \in \text{GL}_3(E)$. This follows from the fact, due to Vinberg [1971], that Γ is F -defined over the *adjoint form* $\overline{\mathcal{G}}$, i.e.,

$$\mathbb{Q}(\text{Tr Ad } \Gamma) = F.$$

This important fact also follows from [Platonov and Rapinchuk 1994, Proposition 4.2].

4. Proofs of Theorems 1.3 and 1.4

We require some elementary results from the theory of discrete subgroups of Lie groups. The primary reference is [Raghunathan 1972]. Let G be a second countable, locally compact group and $\Gamma < G$ a lattice. Recall that G/Γ carries a finite G -invariant measure and Γ is *uniform* in G if G/Γ is compact. For a subgroup $H < G$, we let $Z_G(H)$ denote the centralizer of H in G . We need the following two results.

Lemma 4.1 [Raghunathan 1972, Lemma 1.14]. *Let G be a second countable locally compact group, $\Gamma < G$ a lattice, $\Delta \subset \Gamma$ a finite subset, and $Z_G(\Delta)$ the centralizer of Δ in G . Then, $Z_G(\Delta)\Gamma$ is closed in G .*

Theorem 4.2 [Raghunathan 1972, Theorem 1.13]. *Let G be a second countable locally compact group, $\Gamma < G$ be a uniform lattice, and $H < G$ be a closed subgroup. Then $H\Gamma$ is closed in G if and only if $H \cap \Gamma$ is a lattice in H .*

Proof of Theorem 1.3. Assume that Γ is a cocompact arithmetic lattice in $\text{U}(2, 1)$ containing a complex reflection and that Δ is the subgroup of Γ generated by this reflection. The centralizer of Δ in $\text{U}(2, 1)$ is isomorphic to the extension of $\text{U}(1, 1)$ by the center of $\text{U}(2, 1)$, and is the stabilizer in $\text{U}(2, 1)$ of the wall of the reflection that generates Δ . It follows from Lemma 4.1 and Theorem 4.2 that $\Gamma \cap \text{U}(1, 1)$ is a lattice. Since any sublattice of an arithmetic lattice is arithmetic, Γ contains a totally geodesic arithmetic Fuchsian subgroup. □

Proof of Theorem 1.4. A totally geodesic arithmetic Fuchsian group comes from a subalgebra of \mathcal{D}^r , with notation as in Section 3. When Γ is of second type, \mathcal{D} is a cubic division algebra. The totally geodesic Fuchsian group would correspond to a quaternion subalgebra of \mathcal{D} , which is impossible. When Γ is of first type, this quaternion subalgebra corresponds to rank 2 subspaces of E^3 on which h has signature $(1, 1)$. Therefore, Γ contains complex reflections if and only if Γ is of first type. □

Remark. One can also prove Theorem 1.4 using the structure of unit groups of division algebras.

We now briefly describe how these results generalize to reflections acting on higher-dimensional complex hyperbolic spaces. If $\Gamma < \mathrm{U}(n, 1)$ is a lattice, an element $R \in \Gamma$ is a *codimension s reflection* if it stabilizes a totally geodesic embedded $\mathbf{H}_{\mathbb{C}}^{n-s}$ and acts by an element of the unitary group of the normal bundle to the wall. If Γ is arithmetic, the associated algebraic group is constructed via a hermitian form on \mathfrak{D}^r , where \mathfrak{D} is a division algebra of degree d with involution of the second kind over a totally imaginary field E , and where $rd = n + 1$.

Theorem 4.3. *Suppose $\Gamma < \mathrm{U}(n, 1)$ is a cocompact arithmetic lattice with associated algebraic group coming from a hermitian form on \mathfrak{D}^r , where \mathfrak{D} is a central simple algebra with involution of the second kind. If Γ contains a codimension s reflection, then Γ contains a cocompact lattice in $\mathrm{U}(n - s, 1)$. Also, $n - s + 1 = \ell d$ for some $1 < \ell \leq r$ and the associated algebraic subgroup comes from a hermitian form on \mathfrak{D}^{ℓ} .*

Corollary 4.4. *Let $\Gamma < \mathrm{U}(n, 1)$ be an arithmetic lattice generated by complex reflections through totally geodesic walls isometric to $\mathbf{H}_{\mathbb{C}}^{n-1}$. Then Γ is of so-called first type, i.e., the associated algebraic group is the automorphism group of a hermitian form on E^{n+1} , where E is some totally imaginary quadratic extension of a totally real field.*

5. Arithmetic data for complex hyperbolic triangle groups

In this section, we relate the geometric invariants of a complex triangle to the arithmetic invariants of the complex reflection group. It is the technical heart of the paper.

Let $\Gamma = \Delta(R_1, R_2, R_3)$ be a complex hyperbolic triangle group with reflection factors η_j , complex angles θ_j , and angular invariant ψ . Assume that Γ is an arithmetic lattice in $\mathrm{U}(2, 1)$. By Theorem 1.4, Γ is of first type, so there is an associated hermitian form h over a totally imaginary field E . Let F be the totally real quadratic subfield of E .

Lemma 5.1. *We can choose polar vectors v_j for the reflection R_j so that $v_j \in E^3$.*

Proof. Associated with each reflection is an arithmetic Fuchsian subgroup of Γ . When Γ is a uniform lattice, this follows from Theorem 1.3. For the nonuniform case, see [Holzapfel 1998, Chapter 5]. Arithmetic Fuchsian subgroups stabilizing a complex hyperbolic line come from the h -orthogonal complement of an h -positive line in E^3 . (To be more precise, this line is h -positive over the unique real embedding of F at which h is indefinite.) Any vector in E^3 representing this line is a polar vector for R_j . \square

This leads us to the following important fact.

Lemma 5.2. *Each reflection factor η_j is contained in E .*

Proof. It follows from Proposition 4.2 in [Platonov and Rapinchuk 1994] that there exists an $x_j \in \mathbb{C}$ so that $x_j R_j \in \text{GL}_3(E)$ (see the end of Section 3 above). By Lemma 5.1, and because the h -orthogonal complement to a polar vector evidently has an E -basis, E^3 has a basis of eigenvectors for R_j . The lemma follows from Lemma 2.1. \square

Now we turn to the complex angles and the angular invariant.

Lemma 5.3. *For each j , $\cos^2 \theta_j \in F$ and $e^{2i\psi} \in E$.*

Proof. Choose polar vectors $v_j \in E^3$. The terms in Equations (2) and (3) resulting from these choices of polar vectors are all contained in E . Hence $\cos^2 \theta_j \in F$. One can also prove this using $\text{Tr Ad}(R_1 R_2)$ and Lemma 5.2.

Similarly, consider

$$\delta = h(v_1, v_2)h(v_2, v_3)h(v_3, v_1) = r e^{i\psi} \in E$$

from (4). Note that $e^{i\psi} \in E$ if and only if $r \in E$. Either way, when $\delta \neq 0$, we have $\delta/\bar{\delta} = e^{2i\psi} \in E$. This completes the proof. \square

Combining the above, we see that

$$\mathbb{Q}(\eta_1, \eta_2, \eta_3, \cos^2 \theta_1, \cos^2 \theta_2, \cos^2 \theta_3, e^{2i\psi}) \subseteq E.$$

We can also bound E from above using the fact that $E \subseteq \mathbb{Q}(\text{Tr } \Gamma)$. Using well-known computations of traces for products of reflections (e.g., [Mostow 1980, §4] or [Pratoussevitch 2005]), we have

$$\mathbb{Q}(\text{Tr } \Gamma) = \mathbb{Q}(\eta_1, \eta_2, \eta_3, \cos^2 \theta_1, \cos^2 \theta_2, \cos^2 \theta_3, e^{i\psi} \cos \theta_1 \cos \theta_2 \cos \theta_3).$$

Similarly,

$$\begin{aligned} \mathbb{Q}(\text{Re } \eta_1, \text{Re } \eta_2, \text{Re } \eta_3, \cos^2 \theta_1, \cos^2 \theta_2, \cos^2 \theta_3, \cos^2 \psi) &\subseteq F \\ &\subseteq \mathbb{Q}(\text{Re } \eta_1, \text{Re } \eta_2, \text{Re } \eta_3, \cos^2 \theta_1, \cos^2 \theta_2, \cos^2 \theta_3, \cos \psi \cos \theta_1 \cos \theta_2 \cos \theta_3). \end{aligned}$$

This gives the following.

Corollary 5.4. *Let Γ be a complex hyperbolic triangle group and an arithmetic lattice in $\text{U}(2, 1)$. If the angular invariant of the triangle associated with Γ is rational, then the fields that define Γ as an arithmetic lattice are subfields of a cyclotomic field.*

Let h_Δ be as in (5) and consider h_Δ as a hermitian form on the extension

$$E_\Delta = \mathbb{Q}(\eta_1, \eta_2, \eta_3, \cos \theta_1, \cos \theta_2, \cos \theta_3, e^{i\psi}),$$

of E . It follows from [Mostow 1980, §2] that h and h_Δ are equivalent over E_Δ . Consequently, h_Δ is indefinite over exactly one complex conjugate pair of places of E . This implies that there are precisely $[E_\Delta : E]$ conjugate pairs of places of E_Δ over which h_Δ is indefinite.

Let H be a hermitian form in 3 variables over the complex numbers for which there is a vector with positive H -norm. Then H is indefinite if and only if $\det H < 0$. Since any polar vector has positive h_Δ -norm by definition, we have the following.

Proposition 5.5. *There are exactly $[E_\Delta : E]$ complex conjugate pairs of Galois automorphisms τ of $E_\Delta \subset \mathbb{C}$ under which $\tau(\det h_\Delta)$ is negative. All such automorphisms act trivially on E .*

This has the following consequence for the relationship between the geometry of the triangle and the arithmetic of the lattice.

Corollary 5.6. *If Γ is a complex hyperbolic triangle group and an arithmetic lattice, then the reflection factors of Γ are restricted by the geometry of the triangle. In particular,*

$$E_\Delta = \mathbb{Q}(\cos \theta_1, \cos \theta_2, \cos \theta_3, e^{i\psi}).$$

Proof. Since $\det h_\Delta$ is independent of the reflection factors, for each Galois automorphism of

$$E_\Delta/\mathbb{Q}(\cos \theta_1, \cos \theta_2, \cos \theta_3, e^{i\psi})$$

we obtain a new complex conjugate pair of embeddings of E_Δ into \mathbb{C} such that $\det h_\Delta$ is negative. Any such automorphism necessarily acts nontrivially on some reflection factor η_j . These embeddings of E_Δ lie over different places of E by Lemma 5.2. This contradicts Proposition 5.5. \square

We also obtain the following dependence between the angular invariant and the angles of the triangle.

Proposition 5.7. *If Γ is a complex hyperbolic triangle group and an arithmetic lattice. If Γ has rational angular invariant and $\theta_j \leq \pi/3$ for $j = 1, 2, 3$, then*

$$\cos^2 \psi \in F' = \mathbb{Q}(\cos^2 \theta_1, \cos^2 \theta_2, \cos^2 \theta_3, \cos \theta_1 \cos \theta_2 \cos \theta_3).$$

Proof. If ψ is rational, then E_Δ is a subfield of a cyclotomic field $K_N = \mathbb{Q}(\zeta_N)$, where ζ_N is a primitive N -th root of unity. Therefore the Galois automorphisms of E_Δ are induced by $\zeta_N \mapsto \zeta_N^m$ for some m relatively prime to N .

Consider the stabilizer S of F' in $\text{Gal}(K_N/\mathbb{Q})$. It acts on the roots of unity in E_Δ as a group of rotations along with complex conjugation. By definition of E_Δ , every nontrivial element of S acts nontrivially on $e^{i\psi}$. In particular, if $\cos^2 \psi \notin \mathbb{Q}$ and S contains a rotation through an angle other than an integral multiple of π , then

the orbit of $e^{i\psi}$ under S contains two non-complex conjugate points with distinct negative real parts.

Let τ be any such automorphism of E_Δ . Then, since $\tau(\cos \theta_j) = \cos \theta_j$ for all j by definition of S ,

$$\tau(\det h_\Delta) = 1 - \sum_{j=1}^3 \cos^2 \theta_j + 2\tau(\cos \psi) \prod_{j=1}^3 \cos \theta_j.$$

Furthermore, $1 - \sum \cos^2 \theta_j \leq 0$ for any triple of angles $\theta_j = \pi/r_j$ that are the angles of a hyperbolic triangle with each $r_j \geq 3$. Since $\tau(\cos \psi) < 0$ and $\cos \theta_j > 0$, it follows that $\tau(\det h_\Delta) < 0$. Since τ acts nontrivially on $e^{2i\psi} \in E$, this contradicts Proposition 5.5. Therefore, S is generated by complex conjugation and rotation by π , so $\cos^2 \psi \in F'$. □

Remark. For several of the lattices in [Mostow 1980], $F' = F$ (with notation as above) and $\cos \psi \notin F'$. Thus Proposition 5.7 is the strongest possible constraint on rational angular invariants.

6. Finiteness results

We are now prepared to collect facts from Section 5 to prove Theorem 1.1. A more precise version is the following.

Theorem 6.1. *Suppose that Γ is a complex hyperbolic triangle group and a non-uniform arithmetic lattice in $U(2, 1)$. Then:*

- (1) *Each generator has order 2, 3, 4, or 6.*
- (2) *Each complex angle θ_j of the triangle comes from the set*

$$\{\pi/2, \pi/3, \pi/4, \pi/6, 0\}.$$

- (3) *If ψ is the angular invariant, then $e^{i\psi}$ lies in a biquadratic extension of \mathbb{Q} .*
- (4) *If ψ is rational, then $\psi = s\pi/t$ for*

$$t \in \{2, 3, 4, 6, 8, 12\}.$$

Proof. Since Γ is a nonuniform arithmetic lattice, the associated field E is imaginary quadratic. For (1), we apply Lemma 5.2 to E . For (2) and (3), we apply Lemma 5.3. Then (4) follows from determining those integers m so that $\varphi(m) = 2$ or 4 and $e^{2i\psi}$ is at most quadratic over \mathbb{Q} , where φ is Euler’s totient function. □

See [Paupert 2010; Deraux et al. 2011] for the known nonuniform arithmetic complex hyperbolic triangle groups. We now determine the right triangle groups that can determine an arithmetic lattice in $SU(2, 1)$.

Proof of Theorem 1.6. Suppose that Γ is an arithmetic complex hyperbolic triangle group with $\theta_1 = \pi/2$. The hermitian form h_Δ associated with the triangle has determinant

$$1 - \cos^2 \theta_2 - \cos^2 \theta_3.$$

By Lemma 5.3, this is an element of the totally real field F that defines Γ as an arithmetic lattice. Consequently, there is no Galois automorphism of F over \mathbb{Q} under which this expression remains negative.

This is precisely Takeuchi’s condition that determines whether or not the triangle in the hyperbolic plane with angles $\pi/2, \theta_2, \theta_3$ determines an arithmetic Fuchsian group. The theorem follows from Takeuchi’s classification of arithmetic Fuchsian right triangle groups, Lemma 5.3, and Corollary 5.6. \square

There are 41 such right triangles in \mathbf{H}^2 . We now finish the paper with finiteness for arithmetic complex hyperbolic triangle groups with equilateral complex triangle and rational angular invariant.

Proof of Theorem 1.7. Let Γ be an arithmetic complex hyperbolic triangle group with equilateral triangle of angles π/n and angular invariant ψ . By Proposition 5.7, we can assume that $\psi = s\pi/12n$ for some integer s . Indeed, $F' = \mathbb{Q}(\cos \pi/n)$, and the assertion follows from an easy Galois theory computation.

Then

$$(6) \quad \det h_\Delta = 1 - 3 \cos^2(\pi/n) + 2 \cos(s\pi/12n) \cos^3(\pi/n),$$

so we want to find a nontrivial Galois automorphism of F_Δ whose restriction to F is nontrivial and such that the image of (6) under this automorphism is negative. Let p be the smallest rational prime not dividing $12n$. This determines a nontrivial Galois automorphism τ_p of F_Δ under which

$$(7) \quad \tau_p(\det h_\Delta) = 1 - 3 \cos^2(p\pi/n) + 2 \cos(ps\pi/12n) \cos^3(p\pi/n).$$

It is nontrivial on F by definition. If we show that $\tau_p(\det h_\Delta) < 0$ for n sufficiently large, this, along with Corollary 5.6, suffices to prove the theorem.

First, notice that the function

$$D(x, y) = 1 - 3 \cos^3 x + 2 \cos y \cos^3 x$$

is an increasing function of $x \in (0, \pi/2)$ for any fixed y . In our language, this implies that if y is the angular invariant of an equilateral complex triangle in $\mathbf{H}^2_{\mathbb{C}}$ with angle x , then it remains an angular invariant for a complex triangle with angle x' for any $x' < x$. Similarly, if we know that $\pi/12n$ is an angular invariant for a triangle with angles $p\pi/n$, then we know that $ps\pi/n$ (more precisely, a representative modulo 2π) is the angular invariant of an equilateral triangle in $\mathbf{H}^2_{\mathbb{C}}$ with angles $p\pi/n$. Therefore, it is enough to show that $\pi/12n$ is the angular

invariant of a triangle having angles $p\pi/n$ for all sufficiently large n , where p is the smallest not prime dividing $12n$.

From the above, we conclude further that it suffices to show that there exists a function $q(n)$ such that $p < q(n)$ and

$$(8) \quad 1 - 3 \cos^2(q(n)\pi/n) + 2 \cos(\pi/12n) \cos^3(q(n)\pi/n) < 0$$

for all sufficiently large n . To prove this, we consider the function $j(n)$, defined in [Jacobsthal 1961]. For any integer n , $j(n)$ is the smallest integer such that any $j(n)$ consecutive integers must contain one that is relatively prime to n . Clearly $p \leq j(12n)$.

Iwaniec [1978] proved that

$$j(n) \ll (\log n)^2.$$

Therefore, for any $\epsilon > 0$, there is an n_ϵ so that the first prime number coprime to $12n$ is at most $(\log 12n)^{2+\epsilon}$ for every $n \geq n_\epsilon$. Consider the function

$$f_\epsilon(x) = 1 - 3 \cos^2(\log(12/x)^{2+\epsilon} \pi x) + 2 \cos(\pi x/12) \cos^3(\log(12/x)^{2+\epsilon} \pi x).$$

Then $\lim_{x \rightarrow 0} f_\epsilon(x)$ exists and equals 0 for all $\epsilon > 0$. Further, $x = 0$ is a local maximum of f_ϵ , so $f_\epsilon(1/n) < 0$ for all sufficiently large n .

Taking $q(n) = (\log n)^{2+\epsilon}$ for any small ϵ shows that (8) holds for all sufficiently large n . This implies that (7) is negative for all large n . This proves the theorem. \square

Unfortunately, the proof of Theorem 1.7 isn't effective, so we cannot list the angles that can possibly determine an arithmetic lattice. In particular, we don't know which n makes the bound from [Iwaniec 1978] effective for any $\epsilon > 0$. If this bound is less than $n = 10^5$ for some ϵ , which computer experiments show is extraordinarily likely, then we obtain $n < 5,000,000$. We expect the actual bound to be quite a bit smaller, especially given that the smallest equilateral triangle in \mathbf{H}^2 that defines an arithmetic Fuchsian group has angles $\pi/15$.

Acknowledgments

I thank the referee for several helpful comments.

References

- [Brehm 1990] U. Brehm, "The shape invariant of triangles and trigonometry in two-point homogeneous spaces", *Geom. Dedicata* **33**:1 (1990), 59–76. MR 91c:53048 Zbl 0695.53038
- [Deraux 2006] M. Deraux, "Deforming the \mathbb{R} -Fuchsian $(4, 4, 4)$ -triangle group into a lattice", *Topology* **45**:6 (2006), 989–1020. MR 2007m:32015 Zbl 1120.20052
- [Deraux et al. 2011] M. Deraux, J. R. Parker, and J. Paupert, "Census of the complex hyperbolic sporadic triangle groups", *Exp. Math.* **20**:4 (2011), 467–486. MR 2859902

- [Goldman 1999] W. M. Goldman, *Complex hyperbolic geometry*, Oxford University Press, New York, 1999. MR 2000g:32029 Zbl 0939.32024
- [Goldman and Parker 1992] W. M. Goldman and J. R. Parker, “Complex hyperbolic ideal triangle groups”, *J. Reine Angew. Math.* **425** (1992), 71–86. MR 93c:20076 Zbl 0739.53055
- [Holzapfel 1998] R.-P. Holzapfel, *Ball and surface arithmetics*, Aspects of Mathematics **E29**, Vieweg, Braunschweig, 1998. MR 2000d:14044 Zbl 0980.14026
- [Iwaniec 1978] H. Iwaniec, “On the problem of Jacobsthal”, *Demonstratio Math.* **11**:1 (1978), 225–231. MR 80h:10053 Zbl 0378.10029
- [Jacobsthal 1961] E. Jacobsthal, “Über Sequenzen ganzer Zahlen, von denen keine zu n teilerfremd ist”, *Norske Vid. Selsk. Forh. Trondheim* **33** (1961), 117–124. MR 23 #A2354 Zbl 0096.26002
- [Knus et al. 1998] M.-A. Knus, A. Merkurjev, M. Rost, and J.-P. Tignol, *The book of involutions*, American Mathematical Society Colloquium Publications **44**, American Mathematical Society, Providence, RI, 1998. MR 2000a:16031 Zbl 0955.16001
- [Mostow 1980] G. D. Mostow, “On a remarkable class of polyhedra in complex hyperbolic space”, *Pacific J. Math.* **86**:1 (1980), 171–276. MR 82a:22011 Zbl 0456.22012
- [Parker 2008] J. R. Parker, “Unfaithful complex hyperbolic triangle groups, I: Involutions”, *Pacific J. Math.* **238**:1 (2008), 145–169. MR 2009h:20056
- [Parker and Paupert 2009] J. R. Parker and J. Paupert, “Unfaithful complex hyperbolic triangle groups, II: Higher order reflections”, *Pacific J. Math.* **239**:2 (2009), 357–389. MR 2009h:20057
- [Paupert 2010] J. Paupert, “Unfaithful complex hyperbolic triangle groups, III: Arithmeticity and commensurability”, *Pacific J. Math.* **245**:2 (2010), 359–372. MR 2011d:20094
- [Platonov and Rapinchuk 1994] V. Platonov and A. Rapinchuk, *Algebraic groups and number theory*, Pure and Applied Mathematics **139**, Academic Press, Boston, MA, 1994. MR 95b:11039 Zbl 0841.20046
- [Pratoussevitch 2005] A. Pratoussevitch, “Traces in complex hyperbolic triangle groups”, *Geom. Dedicata* **111** (2005), 159–185. MR 2006d:32036 Zbl 1115.32015
- [Raghunathan 1972] M. S. Raghunathan, *Discrete subgroups of Lie groups*, Ergebnisse der Math. **68**, Springer, New York, 1972. MR 58 #22394a
- [Schwartz 2002] R. E. Schwartz, “Complex hyperbolic triangle groups”, pp. 339–349 in *Proceedings of the International Congress of Mathematicians* (Beijing, 2002), vol. 2, edited by T. T. Li et al., Higher Ed. Press, Beijing, 2002. MR 2004b:57002 Zbl 1022.53034
- [Takeuchi 1977] K. Takeuchi, “Arithmetic triangle groups”, *J. Math. Soc. Japan* **29**:1 (1977), 91–106. MR 55 #2754 Zbl 0344.20035
- [Vinberg 1971] È. B. Vinberg, “Rings of definition of dense subgroups of semisimple linear groups.”, *Izv. Akad. Nauk SSSR Ser. Mat.* **35** (1971), 45–55. MR 43 #4929

Received October 12, 2011. Revised November 30, 2011.

MATTHEW STOVER
 DEPARTMENT OF MATHEMATICS
 UNIVERSITY OF MICHIGAN
 530 CHURCH STREET
 ANN ARBOR, MI 48109
 UNITED STATES
 stoverm@umich.edu

Guidelines for Authors

Authors may submit manuscripts at msp.berkeley.edu/pjm/about/journal/submissions.html and choose an editor at that time. Exceptionally, a paper may be submitted in hard copy to one of the editors; authors should keep a copy.

By submitting a manuscript you assert that it is original and is not under consideration for publication elsewhere. Instructions on manuscript preparation are provided below. For further information, visit the web address above or write to pacific@math.berkeley.edu or to Pacific Journal of Mathematics, University of California, Los Angeles, CA 90095–1555. Correspondence by email is requested for convenience and speed.

Manuscripts must be in English, French or German. A brief abstract of about 150 words or less in English must be included. The abstract should be self-contained and not make any reference to the bibliography. Also required are keywords and subject classification for the article, and, for each author, postal address, affiliation (if appropriate) and email address if available. A home-page URL is optional.

Authors are encouraged to use \LaTeX , but papers in other varieties of \TeX , and exceptionally in other formats, are acceptable. At submission time only a PDF file is required; follow the instructions at the web address above. Carefully preserve all relevant files, such as \LaTeX sources and individual files for each figure; you will be asked to submit them upon acceptance of the paper.

Bibliographical references should be listed alphabetically at the end of the paper. All references in the bibliography should be cited in the text. Use of $\text{Bib}\TeX$ is preferred but not required. Any bibliographical citation style may be used but tags will be converted to the house format (see a current issue for examples).

Figures, whether prepared electronically or hand-drawn, must be of publication quality. Figures prepared electronically should be submitted in Encapsulated PostScript (EPS) or in a form that can be converted to EPS, such as GnuPlot, Maple or Mathematica. Many drawing tools such as Adobe Illustrator and Aldus FreeHand can produce EPS output. Figures containing bitmaps should be generated at the highest possible resolution. If there is doubt whether a particular figure is in an acceptable format, the authors should check with production by sending an email to pacific@math.berkeley.edu.

Each figure should be captioned and numbered, so that it can float. Small figures occupying no more than three lines of vertical space can be kept in the text (“the curve looks like this:”). It is acceptable to submit a manuscript with all figures at the end, if their placement is specified in the text by means of comments such as “Place Figure 1 here”. The same considerations apply to tables, which should be used sparingly.

Forced line breaks or page breaks should not be inserted in the document. There is no point in your trying to optimize line and page breaks in the original manuscript. The manuscript will be reformatted to use the journal’s preferred fonts and layout.

Page proofs will be made available to authors (or to the designated corresponding author) at a website in PDF format. Failure to acknowledge the receipt of proofs or to return corrections within the requested deadline may cause publication to be postponed.

PACIFIC JOURNAL OF MATHEMATICS

Volume 257 No. 1 May 2012

Energy and volume of vector fields on spherical domains	1
FABIANO G. B. BRITO, ANDRÉ O. GOMES and GIOVANNI S. NUNES	
Maps on 3-manifolds given by surgery	9
BOLDIZSÁR KALMÁR and ANDRÁS I. STIPSICZ	
Strong solutions to the compressible liquid crystal system	37
YU-MING CHU, XIAN-GAO LIU and XIAO LIU	
Presentations for the higher-dimensional Thompson groups nV	53
JOHANNA HENNIG and FRANCESCO MATUCCI	
Resonant solutions and turning points in an elliptic problem with oscillatory boundary conditions	75
ALFONSO CASTRO and ROSA PARDO	
Relative measure homology and continuous bounded cohomology of topological pairs	91
ROBERTO FRIGERIO and CRISTINA PAGLIANTINI	
Normal enveloping algebras	131
ALEXANDRE N. GRISHKOV, MARINA RASSKAZOVA and SALVATORE SICILIANO	
Bounded and unbounded capillary surfaces in a cusp domain	143
YASUNORI AOKI and DAVID SIEGEL	
On orthogonal polynomials with respect to certain discrete Sobolev inner product	167
FRANCISCO MARCELLÁN, RAMADAN ZEJNULLAHU, BUJAR FEJZULLAHU and EDMUNDO HUERTAS	
Green versus Lempert functions: A minimal example	189
PASCAL THOMAS	
Differential Harnack inequalities for nonlinear heat equations with potentials under the Ricci flow	199
JIA-YONG WU	
On overtwisted, right-veering open books	219
PAOLO LISCA	
Weakly Krull domains and the composite numerical semigroup ring $D + E[\Gamma^*]$	227
JUNG WOOK LIM	
Arithmeticity of complex hyperbolic triangle groups	243
MATTHEW STOVER	



0030-8730(201205)257:1;1-7