

# CURRENT EVENTS BULLETIN

Friday, January 15, 2010, 1:00 PM to 5:00 PM

Joint Mathematics Meetings, San Francisco, CA

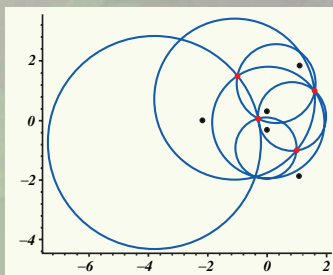
Organized by David Eisenbud, University of California, Berkeley



1:00 PM

**Ben Green**

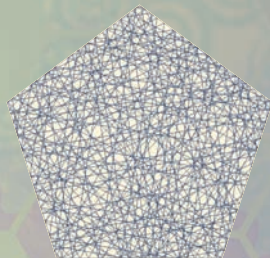
Approximate groups and their applications:  
Work of Bourgain, Gamburd, Helfgott and Sarnak



2:00 PM

**David G. Wagner**

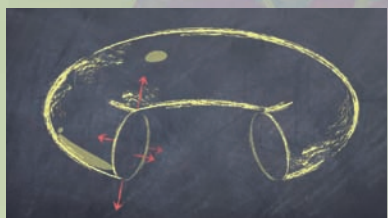
Multivariate stable polynomials: Theory and applications



3:00 PM

**Laura DeMarco**

The conformal geometry of billiards



4:00 PM

**Michael Hopkins**

On the Kervaire Invariant Problem

William Browder, Mark Mahowald and Peter Goerss report on the history of the problem

## Introduction to the Current Events Bulletin

Will the Riemann Hypothesis be proved this week? What is the Geometric Langlands Conjecture about? How could you best exploit a stream of data flowing by too fast to capture? I love the idea of having an expert explain such things to me in a brief, accessible way. I think we mathematicians are provoked to ask such questions by our sense that underneath the vastness of mathematics is a fundamental unity allowing us to look into many different corners -- though we couldn't possibly work in all of them. And I, like most of us, love common-room gossip.

The Current Events Bulletin Session at the Joint Mathematics Meetings, begun in 2003, is an event where the speakers do not report on their own work, but survey some of the most interesting current developments in mathematics, pure and applied. The wonderful tradition of the Bourbaki Seminar is an inspiration, but we aim for more accessible treatments and a wider range of subjects. I've been the organizer of these sessions since they started, but a broadly constituted advisory committee helps select the topics and speakers. Excellence in exposition is a prime consideration.

A written exposition greatly increases the number of people who can enjoy the product of the sessions, so speakers are asked to do the hard work of producing such articles. These are made into a booklet distributed at the meeting. Speakers are then invited to submit papers based on them to the *Bulletin of the AMS*, and this has led to many fine publications.

I hope you'll enjoy the papers produced from these sessions, but there's nothing like being at the talks -- don't miss them!

David Eisenbud, Organizer  
University of California, Berkeley  
de@msri.org

For PDF files of talks given in prior years, see  
<http://www.ams.org/ams/current-events-bulletin.html>.  
The list of speakers/titles from prior years may be found at the end of this booklet.



# APPROXIMATE GROUPS AND THEIR APPLICATIONS: WORK OF BOURGAIN, GAMBURD, HELFGOTT AND SARNAK

BEN GREEN

ABSTRACT. This is a survey of several exciting recent results in which techniques originating in the area known as additive combinatorics have been applied to give results in other areas, such as group theory, number theory and theoretical computer science. We begin with a discussion of the notion of an approximate group and also that of an approximate field, describing key results of Freiman-Ruzsa, Bourgain-Katz-Tao, Helfgott and others in which the structure of such objects is elucidated. We then move on to the applications. In particular we will look at the work of Bourgain and Gamburd on expansion properties of Cayley graphs on  $SL_2(\mathbb{F}_p)$  and at its application in the work of Bourgain, Gamburd and Sarnak on nonlinear sieving problems.

## 1. INTRODUCTION

The subject of *additive combinatorics* has grown enormously over the last ten years and now comprises a large collection of tools with many applications in number theory and elsewhere, for example in group theory and theoretical computer science. It has often been thought a little difficult to specify to an outsider exactly what the subject *is*<sup>1</sup>. However the following point of view seems to be gradually crystallising: additive combinatorics is the study of *approximate mathematical structures* such as approximate groups, rings, fields, polynomials and homomorphisms. It is interested in what the right definitions of these approximate structures are, what can be said about them, and what applications this has to other parts of mathematics.

This article has three main aims. Firstly, we wish to introduce the above point of view to a general audience, focussing in particular on the basic theory of approximate groups and approximate fields. Secondly, we wish to sketch some beautiful applications of these ideas. One of them has to do with the beautiful picture on the cover (for which we thank Cliff Reiter) of an Apollonian circle packing. It is classical that the radii of these circles are all reciprocals of integers. We will describe work of Bourgain, Gamburd and Sarnak giving upper bounds for the number of circles at “depth  $n$ ” which have radius the reciprocal of a prime. Thirdly, we wish to hint at the extraordinary variety of different areas of mathematics which have started to interact with additive combinatorics: geometric group theory, analytic number theory, model theory and point-set topology are just the ones we shall mention here.

---

2000 *Mathematics Subject Classification*. Primary .

This article was written while the author was a fellow at the Radcliffe Institute at Harvard. It is a pleasure to thank the institute for its support and excellent working conditions.

<sup>1</sup>See, for example, my own attempt in the opening remarks of [26].

What we offer here is merely a taste of this viewpoint of additive combinatorics as the theory of approximate structure and of its applications. We do not touch on the theory of approximate polynomials (a.k.a. the theory of Gowers norms) or say much at all about approximate homomorphisms, or anything about the many applications of these two notions. These topics will be covered in detail in forthcoming lecture notes of the author [29].

## 2. APPROXIMATE GROUPS

Before we can define an approximate group, we need to recall what an exact one is. We shall be concerned with finite groups, and we shall be working inside some ambient group  $G$ , so that it makes sense to talk about multiplication of elements and taking inverses. If  $A \subseteq G$  is a finite set then we shall write  $A \cdot A := \{a_1 a_2 : a_1, a_2 \in A\}$  and  $A \cdot A^{-1} := \{a_1 a_2^{-1} : a_1, a_2 \in A\}$ . Later on we shall see more general notations such as  $A \cdot A \cdot A$  and  $A \cdot B$  whose meaning, we hope, will be evident. The following proposition, whose proof is an exercise in undergraduate group theory, gives various criteria for  $A$  to be a subgroup or something very closely related.

**Proposition 2.1.** *Let  $A$  be a finite subset of some ambient group  $G$ . Then we have the following statements<sup>2</sup>:*

- (1)  $|A \cdot A^{-1}| \geq |A|$ , with equality if and only if  $A = Hx$  for some subgroup  $H \leq G$  and some element  $x \in G$ ;
- (2)  $|A \cdot A| \geq |A|$ , with equality if and only if  $A = Hx$  for some subgroup  $H \leq G$  and some element  $x$  in the normaliser  $N_G(H)$ ;
- (3) The number of quadruples  $(a_1, a_2, a_3, a_4) \in A^4$  with  $a_1 a_2^{-1} = a_3 a_4^{-1}$  is at most  $|A|^3$ , with equality if and only if  $A = Hx$  for some subgroup  $H \leq G$  and some element  $x \in G$ ;
- (4) The number of quadruples  $(a_1, a_2, a_3, a_4) \in A^4$  with  $a_1 a_2 = a_3 a_4$  is at most  $|A|^3$ , with equality if and only if  $A = Hx$  for some subgroup  $H \leq G$  and some element  $x \in N_G(H)$ ;
- (5)  $\mathbb{P}(a_1 a_2 \in A | a_1, a_2 \in A) \leq 1$ , with equality if and only if  $A = H$  for some subgroup  $H \leq G$ ;
- (6)  $\mathbb{P}(a_1 a_2^{-1} \in A | a_1, a_2 \in A) \leq 1$ , with equality if and only if  $A = H$  for some subgroup  $H \leq G$ .

This would be a rather odd proposition to see formulated in an algebra text. However each of the statements (1) – (6) has been constructed as an inequality in such a way that one may ask when equality approximately holds. Before we can talk about such approximate variants, however, we need to know how approximate they will be. For this purpose we introduce a parameter  $K \geq 1$ ; larger values of  $K$  will indicate more approximate, and thus less structured, objects<sup>3</sup>.

---

<sup>2</sup>Statements (5) and (6) look “probabilistic” but this is just a notation. By  $\mathbb{P}(a_1 a_2 \in A | a_1, a_2 \in A)$  we mean simply the proportion of all pairs  $a_1, a_2 \in A$  for which  $a_1 a_2$  also lies in  $A$ .

<sup>3</sup>In practice the theory when  $K \approx 1$  is very different from the theory when, for example,  $K \sim 100$ . In the former setting, these approximate notions of subgroup constitute very small perturbations of the exact characterisations of Proposition 2.1, and it turns out (though is not always trivial to prove) that the approximate objects so defined are small perturbations of the exact objects characterised by Proposition 2.1. In conversation Tao and I tend to refer to this regime as “the 99% world”, an expression I would not be averse to popularising. In this paper  $K$  will be much larger, causing the theory to become much richer. Tao and I call this the “1%

Consider, then, the following list of properties that a finite set  $A \subseteq G$  might enjoy.

- (1)  $|A \cdot A^{-1}| \leq K|A|$ ;
- (2)  $|A \cdot A| \leq K|A|$ ;
- (3) The number of quadruples  $(a_1, a_2, a_3, a_4) \in A^4$  with  $a_1a_2^{-1} = a_3a_4^{-1}$  is at least  $|A|^3/K$ ;
- (4) The number of quadruples  $(a_1, a_2, a_3, a_4) \in A^4$  with  $a_1a_2 = a_3a_4$  is at least  $|A|^3/K$ ;
- (5)  $\mathbb{P}(a_1a_2 \in A | a_1, a_2 \in A) \geq 1/K$ ;
- (6)  $\mathbb{P}(a_1a_2^{-1} \in A | a_1, a_2 \in A) \geq 1/K$ .

Now these are by no means as closely equivalent as the properties (1) – (6) in Proposition 2.1. Let us give an example in which the ambient group is  $\mathbb{Z}$ , and where we use additive rather than multiplicative notation. Take  $A = \{1, \dots, n\} \cup \{2^{n+1}, 2^{n+2}, \dots, 2^{2n}\}$ . Then it is easy to check that (3) and (4) are both satisfied with any  $K > 12$ , as  $n \rightarrow \infty$ , this being because there are  $\frac{2}{3}n^3(1 + o(1))$  solutions to  $a_1 + a_2 = a_3 + a_4$  with  $a_1, a_2, a_3, a_4 \in \{1, \dots, n\}$ . On the other hand the sumset  $A + A$  contains the numbers  $2^{n+i} + j$  for each pair  $i, j$  with  $0 < i, j \leq n$ . Since these numbers are all distinct, we have  $|A + A| \geq n^2 = |A|^2/2$ , which means that if  $n$  is sufficiently large depending on  $K$  then (2) is not satisfied at all.

Rather remarkably, however, there is a sense in which the concepts (1) – (6) are all *roughly* the same. To say what we mean by that, we introduce the following notion of rough equivalence<sup>4</sup>.

**Definition 2.2** (Rough Equivalence). Suppose that  $A$  and  $B$  are two finite sets in some ambient group and that  $K \geq 1$  is a parameter. Then we write  $A \sim_K B$  to mean that there is some  $x$  in the ambient group such that  $|A \cap Bx| \geq \max(|A|, |B|)/K$ . We say that  $A$  and  $B$  are roughly equivalent (with parameter  $K$ ).

The remarkable fact alluded to above is the following. For every choice of  $j, j' \in \{1, \dots, 6\}$ , suppose that some set  $A$  satisfies condition (j) in the list above with parameter  $K$ . Then there is a set  $B$  satisfying condition (j') with parameter  $K' = \text{poly}(K)$  (some polynomial in  $K$ ) such that  $A \sim_{K'} B$ . Of particular note is the fact that the weak “statistical” properties (3) – (6) imply the apparently more structured properties (1) and (2). The proof of this is not at all trivial and the main content of it is the so-called Balog-Szemerédi-Gowers theorem [22], generalised to the nonabelian setting in the fundamental paper of Tao [59], as well as a collection of “sumset estimates” which, in the abelian case, I refer to collectively as *Ruzsa calculus* [28]. These estimates of Ruzsa have such a classical role in the theory that we record two of them, in the abelian setting, here: we will mention these two again later on.

---

world” although the parameter  $K$  could be anything between 2 (say) and some small power of  $|A|$ .

<sup>4</sup>The fact that we have written  $Bx$  rather than  $xB$  is a little arbitrary. The notion of rough equivalence will, in this survey, be applied to classes of sets (such as (1) – (6) here) which are invariant under conjugation, in which case whether we multiply on the left or on the right in the definition makes little difference.

**Theorem 2.3** (Ruzsa). *Suppose that  $A_1, A_2$  and  $A_3$  are finite sets in some ambient abelian group. Then  $|A_1||A_2 - A_3| \leq |A_1 - A_2||A_1 - A_3|$  and  $|A_1 + A_1| \leq |A_1 - A_1|^3/|A_1|^2$ .*

The original paper [47], the book [63] or the notes [28] may be consulted for more details. The first estimate is true in general groups but adapting the second requires care: see [59].

It might be remarked that for many pairs  $(j)$  and  $(j')$  the correspondence between the relevant properties is a little tighter than mere rough equivalence, and often this can be useful. We shall not dwell on this point here. In the paper of Tao just mentioned one finds what has become the “standard” notion of an approximate group.

**Definition 2.4** (Approximate group). Suppose that  $A$  is a finite subset of some ambient group and that  $K \geq 1$  is a parameter. Then we say that  $A$  is a  $K$ -approximate group if it is symmetric (that is, if  $a \in A$  then  $a^{-1} \in A$ , and the identity lies in  $A$ ) and if there is a set  $X$  in the ambient group with  $|X| \leq K$  and such that  $A \cdot A \subseteq X \cdot A$ .

This notion, it turns out, is roughly equivalent to (1) - (6) above. It has certain advantages over (1) - (6), for example as regards its behaviour under homomorphisms. It is also clear that an approximate group in this sense enjoys good control of iterated sumsets. Thus, for example,  $A \cdot A \cdot A \subseteq X \cdot X \cdot A$ , which means that  $|A^3| = |A \cdot A \cdot A| \leq K^2|A|$ , and similarly  $|A^n| \leq K^{n-1}|A|$  where  $A^n$  denotes the set of all products  $a_1 \dots a_n$  with  $a_1, \dots, a_n \in A$ . From now on, when we speak of an approximate group, we will be referring primarily to Definition 2.4.

With this discussion in mind, we can introduce what might be termed the rough classification problem of approximate group theory.

**Question 2.5.** Consider the collection  $\mathcal{C}$  of all  $K$ -approximate groups  $A$  in some ambient group  $G$ . Is there some “highly structured” subcollection  $\mathcal{C}'$  such that every  $A \in \mathcal{C}$  is roughly equivalent to some set  $B \in \mathcal{C}'$  with parameter  $K'$ , where  $K'$  depends only on  $K$ ?

This question has been addressed in a great many different contexts, starting with the Freiman-Ruzsa theorem [20, 48], which gives an answer for subsets of  $\mathbb{Z}$ . Here, it is possible to take  $\mathcal{C}'$  to consist of the so-called *generalised arithmetic progressions*, that is to say sets  $B$  of the form

$$B := \{l_1x_1 + \dots + l_dx_d : l_i \in \mathbb{Z}, |l_i| \leq L_i\},$$

where  $x_1, \dots, x_d \in \mathbb{R}$ , the quantities  $L_1, \dots, L_d$  are “lengths” and  $d \leq K$ . Note in particular that, even in the highly abelian setting of the integers  $\mathbb{Z}$ , approximate groups are a more general kind of object than genuine subgroups. That is, the theory of approximate groups, even up to rough equivalence, is a little richer than the theory of finite subgroups of  $\mathbb{Z}$  (which is in fact a rather trivial theory). The remarkable feature of the Freiman-Ruzsa theorem is that the theory is not *much* richer, in the sense that generalised progressions remain highly “algebraic” objects. Here is a list of other contexts in which the question has been at least partially answered:

- abelian groups [30];

- nilpotent and solvable groups [10, 11, 19, 50, 61];
- free groups [46];
- linear groups  $\mathrm{SL}_2(\mathbb{R})$  [16],  
 $\mathrm{SL}_2(\mathbb{C})$  [13, 33, 34],  
 $\mathrm{SL}_3(\mathbb{Z})$  [13],  
 $\mathrm{SL}_3(\mathbb{C})$  (sketched in [34]),  
“bounded” subsets of  $\mathrm{SL}_n(\mathbb{C})$  including  $\mathrm{U}_n(\mathbb{C})$  [12],  
 $\mathrm{SL}_2(\mathbb{F}_p)$  [33],  
 $\mathrm{SL}_3(\mathbb{F}_p)$  [34]  
and  $\mathrm{SL}_2(\mathbb{Z}/q\mathbb{Z})$  for various other  $q$  (cf. [4]).

It is generally felt that approximate groups in quite general contexts can be controlled by objects built up from genuine subgroups and nilpotent objects; this has been found in all of the examples just mentioned and is suggested by the famous theorem of Gromov on groups with polynomial growth [32] and the recent quantitative formulation of it due to Shalom and Tao [55]. Quite precise suggestions along these lines have been made by Helfgott, Lindenstrauss and others: more information on this can be found on Tao’s blog [62].

Before leaving this subject, we remark that even (perhaps *especially*) in the abelian case the issue of the dependence of  $K'$  on  $K$  is far from being resolved. No examples are known to rule out the possibility that, with the right definition of the highly-structured class  $\mathcal{C}'$ ,  $K'$  can be taken to be polynomial in  $K$ . In particular this is conjectured when the ambient group  $G$  is  $\mathbb{F}_2^{\mathbb{Z}}$ , the countable infinite vector space over the field of two elements, and  $\mathcal{C}'$  consists of (finite) subgroups of  $G$ . This assertion<sup>5</sup> is known as the *polynomial Freïman-Ruzsa conjecture* [49], see also [25, 27]. It is equivalent to the following question which, for many years, I have tried to advertise to those for whom the word cohomology holds no fear.

**Question 2.6.** Suppose that  $\phi : \mathbb{F}_2^n \rightarrow \mathbb{F}_2^{\mathbb{Z}}$  is a map such that  $\phi(x+y) - \phi(x) - \phi(y)$  takes on at most  $K$  different values as  $x, y$  range over  $\mathbb{F}_2^n$ . Is it true that  $\phi = \tilde{\phi} + \eta$ , where  $\tilde{\phi}$  is linear and  $\eta$  takes on at most  $K' = \mathrm{poly}(K)$  different values?

It is a very easy exercise to obtain such a statement with  $K' = 2^K$  but, so far as I know, no serious improvement of this bound has ever been obtained<sup>6</sup>.

### 3. APPROXIMATE RINGS AND FIELDS

Fortified by the experiences of the last section, one might attempt to come up with a sensible notion of an approximate *ring*. A natural one, based perhaps on (2) in the previous section, is as follows: if  $A$  is a finite subset of some ambient ring  $R$ , we say that it is a  $K$ -approximate ring if  $|A+A| \leq K|A|$  and  $|A \cdot A| \leq K|A|$ . Here, of course,  $A+A := \{a_1+a_2 : a_1, a_2 \in A\}$  and  $A \cdot A = \{a_1a_2 : a_1, a_2 \in A\}$  as before. If  $R = \mathbb{F}$  is actually a field (or an integral domain, which embeds into its field of fractions) then we refer to  $A$  as an approximate field, noting that approximate closure under division is essentially automatic in view of the rough equivalence of the notions (1) and (2) of approximate group.

<sup>5</sup>There are variants of this conjecture over other groups, such as  $\mathbb{Z}$ ; see [23, 31].

<sup>6</sup>I would be very interested to see even a bound of the form  $2^{o(K)}$ .



The study of approximate rings and fields was initiated in a paper of Erdős and Szemerédi [17] who proved (though not in this language!) that a  $K$ -approximate subfield of  $\mathbb{Z}$  must have size  $\text{poly}(K)$ . They in fact conjectured that the right bound is  $C_\epsilon K^{1+\epsilon}$  for any  $\epsilon > 0$ , but this is so far unresolved; the best exponent so far obtained is  $3 + \epsilon$ , a result of Solymosi [57]. Note that this is equivalent to, and more usually stated as, the lower bound

$$\max(|A + A|, |A \cdot A|) \geq c_\epsilon |A|^{4/3+\epsilon}$$

for all finite sets  $A \subseteq \mathbb{Z}$ . In a different paper [56], Solymosi generalised the Erdős-Szemerédi result to show that every  $K$ -approximate subfield of  $\mathbb{C}$  has size at most  $2^{12}K^4$ .

The general theory of approximate fields can be said to have started with the papers of Bourgain-Katz-Tao [8] and Bourgain-Glibichuk-Konyagin [7, 9], where<sup>7</sup> the following result is established.

**Theorem 3.1.** *Let  $p$  be a prime and let  $K \geq 2$ . Then every  $K$ -approximate subfield of  $\mathbb{F}_p$  has size at most  $K^C$  or at least  $K^{-C}p$ , for some absolute constant  $C$ .*

The arguments on page 384 of [7], though they are phrased in a more limited context, essentially prove that every approximate subfield (in an arbitrary ambient field) must be roughly equivalent to a genuine finite subfield. This unifies the results of Erdős-Szemerédi and Solymosi with Theorem 3.1. In fact something similar is true for approximate rings, at least provided the ambient ring  $R$  does not have “too many” zero divisors. These issues are comprehensively explored in an interesting paper [60] of Tao, which also has a very comprehensive collection of references.

Suppose that  $A$  is a  $K$ -approximate field in some ambient field  $\mathbb{F}$ , that is to say both  $|A \cdot A|$  and  $|A + A|$  are bounded by  $K|A|$ . We are going to sketch a proof that  $\mathbb{F}$  must contain a genuine subfield  $B$  which is “close” to  $A$ . The first step is to prove the *Katz-Tao lemma*, which asserts that  $A$  (or, more precisely, a large subset  $A' \subseteq A$ ) behaves in a manner which more strongly resembles that of a field: that is to say,  $A$  is almost closed under both addition/subtraction and multiplication/division *simultaneously*. To give a (relevant) example, the set

$$\overline{A} := \left\{ a_5 + a_6 \frac{a_1 - a_3}{a_4 - a_2} : a_1, \dots, a_6 \in A \right\}$$

has size  $\overline{K}|A|$ , where  $\overline{K} = \text{poly}(K)$ .

A slick proof of the Katz-Tao lemma is given in [60, Section 2.5] and we shall say little more about it here other than to remark that it involves a combination of Ruzsa’s sumset calculus and clever elementary arguments. Personally, I regard it as part of the “basic” theory of approximate fields as opposed to the “structural theory”, to be regarded on the same level as the arguments used to show that definitions (1) – (6) of an approximate group are roughly equivalent (namely, Ruzsa’s sumset calculus and the Balog-Szemerédi-Gowers theorem). In other words one might argue that the smallness of  $\overline{A}$ , or of similar objects, might be taken as an alternative *definition* of approximate field.

<sup>7</sup>The original paper [8] of Bourgain, Katz and Tao did not quite classify the very small (smaller than  $p^\delta$ ) approximate subrings of  $\mathbb{F}_p$ ; this restriction was removed in [7, 9]. Very often the approximate fields under consideration in a given setting will have size at least  $p^\delta$ , and for this reason one often refers to the Bourgain-Katz-Tao theorem.

Suppose that  $A$  is known to have this property, that is to say  $|\overline{A}| \leq \overline{K}|A|$ . Then it is possible to establish an intriguing dichotomy: if  $\xi \in \mathbb{F}^\times$  then either

$$(3.1) \quad |A + A\xi| = |A|^2$$

or

$$(3.2) \quad |A + A\xi| \leq \overline{K}|A|.$$

Here,  $A + A\xi$  refers to the set of all  $a_1 + a_2\xi$  with  $a_1, a_2 \in A$ . To see why this is so, note that  $|A + A\xi| \leq |A|^2$  and that equality occurs if and only if the elements  $a_1 + a_2\xi$  are all distinct. If equality does not occur then we may find a nontrivial solution to  $a_1 + a_2\xi = a_3 + a_4\xi$ , which means that  $\xi = \frac{a_1 - a_3}{a_4 - a_2}$ . But then every element of  $A + A\xi$  has the form

$$a_5 + a_6\xi = a_5 + a_6 \frac{a_1 - a_3}{a_4 - a_2},$$

and thus lies in  $\overline{A}$ .

On the other hand, it is not hard to see using Ruzsa calculus<sup>8</sup> that if  $\xi_1, \xi_2$  satisfy (3.2) then for  $\xi = \xi_1 + \xi_2, \xi_1 - \xi_2, \xi_1\xi_2, \xi_1\xi_2^{-1}$  we have

$$|A + \xi \cdot A| \leq \overline{K}^C |A| \leq K^{C'} |A|$$

for absolute constants  $C, C'$ . If  $K$  is a sufficiently small power of  $|A|$  then this means that (3.1) cannot hold, forcing us to conclude that (3.2) holds for  $\xi$ . In this way we identify the set<sup>9</sup> of all  $\xi$  satisfying (3.2) as a genuine subfield of  $\mathbb{F}$ . Straightforward additional arguments allow one to show that this subfield and  $\mathbb{F}$  are roughly equivalent.

The original argument of [8] is different and specific to  $\mathbb{F}_p$  but rather fun and, given the preceding discussion, it is not hard to say a few meaningful words about it. Suppose for the sake of illustration that  $A \subseteq \mathbb{F}_p$  is a  $K$ -approximate subfield of size  $\sim p^{1/10}$ ; our task is to derive a contradiction if (say)  $K = p^{o(1)}$ . Suppose that the Katz-Tao lemma has already been applied, so that  $\overline{A}$ , as defined above, is known to be small. The sets  $\overline{\overline{A}}, \overline{\overline{\overline{A}}}, \dots$  arising from (boundedly many) successive applications of this operation may also be shown to be small. Now simple averaging arguments (using nothing more than the fact that  $|A| = p^{1/10}$ ) show that  $\mathbb{F}_p$  has dimension at most 100 (say) as a “vector space” over  $A$ ; that is, there exist  $x_1, \dots, x_{100} \in \mathbb{F}_p$  such that

$$(3.3) \quad \mathbb{F}_p = Ax_1 + \dots + Ax_{100}.$$

Now  $x_1, \dots, x_{100}$  cannot be a “basis” for  $\mathbb{F}_p$  over  $A$  since otherwise we would have  $p = |A|^{100}$ , contrary to the assumption that  $p$  is prime. Thus there must exist some  $x \in \mathbb{F}_p$  which is representable in two different ways as

$$x = a_1x_1 + \dots + a_{100}x_{100} = a'_1x_1 + \dots + a'_{100}x_{100}$$

<sup>8</sup>In addition to the bounds of Theorem 2.3 one requires an inequality controlling  $|A_1 + A_2 + A_3|$  in terms of the  $|A_i + A_j|$ .

<sup>9</sup>Note that this set may be identified with  $\frac{A-A}{(A-A)^\times}$ .

with  $a_1, \dots, a_{100}, a'_1, \dots, a'_{100} \in A$ . Suppose, without loss of generality, that  $a_{100} \neq a'_{100}$ . Then

$$x_{100} = \frac{(a_1 - a'_1)x_1 + \dots + (a_{99} - a'_{99})x_{99}}{a'_{100} - a_{100}}.$$

By substituting this expression for  $x_{100}$  into (3.3), we see that

$$\mathbb{F}_p = \overline{A}x_1 + \dots + \overline{A}x_{99}.$$

Repeating the argument gives (without loss of generality)

$$\mathbb{F}_p = \overline{\overline{A}}x_1 + \dots + \overline{\overline{A}}x_{98},$$

and we may continue in this fashion to get, eventually,

$$\mathbb{F}_p = \overline{\overline{\overline{A}}}x_1.$$

This is contrary to the fact that none of the sets  $\overline{A}, \overline{\overline{A}}, \dots$  has size much larger than that of  $A$  itself, namely about  $p^{1/10}$ , and a contradiction ensues.

Remarkably, the main “dimension reduction” idea here comes from a paper in point-set topology, namely Edgar and Miller’s solution of the Erdős-Volkmann ring problem [15] (that is, the statement that all Borel subrings of  $\mathbb{R}$  have dimension 0 or 1). See in particular Lemma 1.3 of that paper.

#### 4. HELFGOTT’S RESULTS

In this section we discuss the results of Helfgott [33, 34] concerning approximate subgroups of

$$\mathrm{SL}_2(\mathbb{F}_p) := \left\{ \begin{pmatrix} a & b \\ c & d \end{pmatrix} : a, b, c, d \in \mathbb{F}_p : ad - bc = 1 \right\}.$$

Helfgott proves the following.

**Theorem 4.1** (Helfgott). *Suppose that  $A \subseteq \mathrm{SL}_2(\mathbb{F}_p)$  is a  $K$ -approximate group. Then  $A$  is roughly  $K^C$ -equivalent to an upper-triangular  $K^C$ -approximate subgroup of  $\mathrm{SL}_2(\mathbb{F}_p)$  (that is, an approximate subgroup conjugate to a set of upper-triangular matrices).*

Rather than discuss Helfgott’s result itself, we discuss the analogous question for  $\mathrm{SL}_2(\mathbb{C})$ . Here the answer is rather simpler and is given in [13], based on Helfgott’s work.

**Theorem 4.2.** *Suppose that  $A \subseteq \mathrm{SL}_2(\mathbb{C})$  is a  $K$ -approximate group. Then  $A$  is roughly  $K^C$ -equivalent to an abelian  $K^C$ -approximate subgroup of  $\mathrm{SL}_2(\mathbb{C})$ .*

If desired the abelian approximate group could itself be controlled by a generalised progression using the Freïman-Ruzsa theorem.

We will only sketch a proof of the weaker result that  $A$  is  $K^C$ -equivalent to an upper-triangular  $K^C$ -approximate subgroup, that is to say the direct analogue of Helfgott’s result. In  $\mathrm{SL}_2(\mathbb{C})$ , additional arguments may then be applied to prove

Theorem 4.2; there are no such arguments in  $\mathrm{SL}_2(\mathbb{F}_p)$ , since the upper-triangular “Borel subgroup”

$$\left\{ \begin{pmatrix} \lambda & \mu \\ 0 & \lambda^{-1} \end{pmatrix} : \lambda \in \mathbb{F}_p^*, \mu \in \mathbb{F}_p \right\}$$

is not close to abelian.

The proof of this weak form of Theorem 4.2 is simpler than that of Theorem 4.1 in two major ways. Firstly since  $\mathbb{C}$  is algebraically closed we may talk about eigenvalues, eigenvectors and diagonalization without the need to pass to an extension field, whereas over  $\mathbb{F}_p$  we would have to involve the quadratic extension  $\mathbb{F}_{p^2}$ . Secondly, the structure of  $K$ -approximate subfields of  $\mathbb{C}$  is easy to describe: by the theorem of Solymosi [56] they are all sets of size at most  $2^{12}K^4$ . Theorem 3.1, by contrast, has to allow for those approximate fields which are almost all of  $\mathbb{F}_p$ . Worse still, to handle  $\mathrm{SL}_2(\mathbb{F}_p)$  Helfgott must in fact classify approximate subfields of  $\mathbb{F}_{p^2}$ , and this involves the additional possibility of sets which are close to the subfield  $\mathbb{F}_p$ .

For the sake of exposition, we will assume in the first instance that  $A$  is a genuine finite subgroup of  $\mathrm{SL}_2(\mathbb{C})$ ; our task is to show that  $A$  contains a large upper-triangular subgroup. When we have sketched how Helfgott’s argument looks in this case we will remark on the additional technicalities required to make the argument “robust” enough to apply to  $K$ -approximate groups.

The key idea in Helfgott’s argument, referred to by subsequent authors as *trace amplification*, involves examining the set of traces

$$\mathrm{tr} A := \{\mathrm{tr} a : a \in A\}.$$

We will sketch a proof that a large subset of this set of traces is a  $2^{24}$ -approximate subfield of  $\mathbb{C}$  of size greater than  $2^{108}$ . This contradicts Solymosi’s theorem [56] and so we must be in one of those degenerate situations. Careful analysis of each of them leads to the conclusion that  $A$  is roughly upper-triangular.

The first degenerate situation to analyse is that in which  $\mathrm{tr} A$  is small, an appropriate notion of *small* being  $|\mathrm{tr} A| \leq 2^{111}$ . Now a linear algebra computation (Lemma 4.2 of [6]) confirms that if  $g, h \in A$  are elements without a common eigenvector in  $\mathbb{C}^2$  then the map

$$\mathrm{SL}_2(\mathbb{C}) \rightarrow \mathbb{C}^3 : x \mapsto (\mathrm{tr} x, \mathrm{tr}(gx), \mathrm{tr}(hx))$$

is at most two-to-one. This, or rather the fact that something like this holds, is not at all surprising: indeed knowledge of  $\mathrm{tr}(x), \mathrm{tr}(gx), \mathrm{tr}(hx)$  together with the fact that  $\det(x) = 1$  provides four pieces of information which, generically, ought to more-or-less determine the four entries of the matrix  $x$ . If  $A$  contains two such elements  $g, h$  then it follows that we have

$$|A| \leq 2|\mathrm{tr} A|^3 \leq 2^{334},$$

and so  $|A|$  is also small<sup>10</sup>. If, by contrast,  $A$  does not contain two such elements, and if  $|A| > 3$ , then it is easy to see that there is some  $v \in \mathbb{C}^2$  which is an eigenvector for all of  $A$  simultaneously. With respect to a basis containing  $v$ , every matrix in  $A$  is upper-triangular.

<sup>10</sup>Additive combinatorics has a bad reputation for referring to quantities like  $2^{334}$  as “small”. “Bounded by an absolute constant” might be more appropriate.

Suppose, then, that  $|\operatorname{tr} A| > 2^{111}$ . In particular (!) there is some element  $g \in A$  which is non-parabolic, or in other words  $\operatorname{tr} g \neq \pm 2$ ; such elements have distinct eigenvalues and so are diagonalisable.

Write  $A' \subseteq A$  for the set of non-parabolic elements; then  $|\operatorname{tr} A'| \geq |\operatorname{tr} A| - 2 \geq \frac{1}{2}|\operatorname{tr} A|$ . Now in  $\operatorname{SL}_2(\mathbb{C})$  the trace of a non-parabolic element  $g$  completely determines the conjugacy class of  $g$ . It follows that there is some non-parabolic  $g \in A$  such that the conjugacy class of  $A$  containing  $g$  has size at most  $2|A|/|\operatorname{tr} A|$ . By the orbit-stabiliser theorem, the centraliser<sup>11</sup>

$$T = C_A(g) = \{a \in A : ag = ga\}$$

has size at least  $\frac{1}{2}|\operatorname{tr} A|$ . But by changing basis so that  $g$  is in diagonal form (with distinct diagonal entries) it is not hard to check that  $T$  consists entirely of diagonal matrices. No single trace can arise from more than two of these elements, and so  $|\operatorname{tr} T| \geq \frac{1}{4}|\operatorname{tr} A| > 2^{109}$ . We shall show that the set

$$R := \{\operatorname{tr} a^2 : a \in T\}$$

is a  $2^{24}$ -approximate subfield of  $\mathbb{C}$ . Noting that

$$(4.1) \quad |R| \geq \frac{1}{2}|\operatorname{tr} T| > 2^{108},$$

this is contrary to Solymosi's theorem. In order to do this we play around a little with traces. Such playing around is most productive if, in the basis just selected, there is an element  $a = \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix} \in A$  with  $a_{11}a_{12}a_{21}a_{22} \neq 0$ . The absence of such an element is another degenerate situation to analyse, and once again one can check<sup>12</sup> that  $A$  must be either upper-triangular or else equal to one of the dihedral groups, each of which has an index two abelian subgroup.

Now let us note that

$$(4.2) \quad R \cdot R \subseteq R + R,$$

this being a consequence of the fact that

$$(t_1^2 + t_1^{-2})(t_2^2 + t_2^{-2}) = (t_1^2 t_2^2 + t_1^{-2} t_2^{-2}) + (t_1^2 t_2^{-2} + t_1^{-2} t_2^2).$$

Let us also note that

$$\operatorname{tr} \left( \begin{pmatrix} t_1 t_2 & 0 \\ 0 & t_1^{-1} t_2^{-1} \end{pmatrix} a \begin{pmatrix} t_1 t_2^{-1} & 0 \\ 0 & t_1^{-1} t_2 \end{pmatrix} a^{-1} \right) = \mu(t_1^2 + t_1^{-2}) + \lambda(t_2^2 + t_2^{-2}),$$

where  $\mu := a_{11}a_{22} \neq 0$  and  $\lambda := -a_{12}a_{21} \neq 0$ , which means that

$$\lambda R + \mu R \subseteq \operatorname{tr} A.$$

In particular

$$|R + \frac{\mu}{\lambda}R| = |\lambda R + \mu R| \leq |\operatorname{tr} A| \leq 16|R|,$$

which, by Ruzsa's inequalities (Theorem 2.3, applied with  $A_1 = \frac{\mu}{\lambda}R$  and  $A_2 = A_3 = -R$ ) implies that  $|R + R| \leq 2^{24}|R|$ . This, together with (4.2), implies that  $R$  is a  $2^{24}$ -approximate subring of  $\mathbb{C}$ . By Solymosi's theorem this implies that  $|R| \leq 2^{108}$ , contrary to (4.1).  $\square$

In the above sketch we assumed, of course, that  $A$  was actually a finite subgroup. However the argument was of a type that can be made to work for  $K$ -approximate

<sup>11</sup> $T$  is for *torus*, the word used for such a subgroup in Lie theory.

<sup>12</sup>This is, admittedly, a somewhat tedious check.

groups also. To explain what we mean by this let us remark, rather vaguely, on how one or two of the steps adapt and then offer some general remarks.

*Orbit-Stabiliser theorem.* If  $A$  is a group and if  $x \in A$  then we used the fact that the size of the conjugacy class  $\Sigma(x)$  containing  $x$  and that of the centraliser  $C_A(x)$  are related by  $|\Sigma(x)||C_A(x)| = |A|$ . In fact we only used the inequality  $|C_A(x)| \geq |A|/|\Sigma(x)|$ , giving us an element with large centraliser, and here is a simple way of seeing why this holds: all of the conjugates  $axa^{-1}$ ,  $a \in A$ , lie in  $\Sigma(x)$ , and so by the pigeonhole principle there must be distinct elements  $a_1, \dots, a_k \in A$ ,  $k \geq |A|/|\Sigma(x)|$ , with  $a_1xa_1^{-1} = \dots = a_kxa_k^{-1}$ . But then the elements  $a_i^{-1}a_1$ ,  $i = 1, \dots, k$ , centralise  $x$ . Now if  $A$  is only a  $K$ -approximate group then this argument does not quite work, as there is no well-defined notion of conjugacy class. However a similar pigeonhole argument nonetheless gives us an element with large centraliser, since the conjugates  $axa^{-1}$  are all constrained to lie in  $A^3$ , a set of size at most  $K^2|A|$ .

*Escape from subvarieties.* A more interesting point concerns the location of an element of  $A$  which, in a given basis, has no zero entries. Whilst this might not be *a priori* possible if  $A$  is only an approximate group, it *is* possible to find such an element in  $A^n$  for some bounded  $n$  (independent of the approximation parameter  $K$ ), and this is good enough for Helfgott's purposes. This is a special case of a nice lemma of Eskin, Mozes and Oh [18] called "escape from subvarieties". The point is that the group  $\langle A \rangle$  generated by  $A$ , if it is not almost upper-triangular, contains an element with no zero entries – indeed this fact was used in the above sketch. In other words,  $\langle A \rangle$  is not contained in the subvariety of  $\mathrm{SL}_2(\mathbb{C})$  defined by

$$V := \left\{ \begin{pmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \end{pmatrix} : x_{11}x_{12}x_{21}x_{22} = 0 \right\}.$$

The Eskin-Mozes-Oh result states that in such a situation we can find "evidence" for the non-containment of  $\langle A \rangle$  inside  $V$  by taking just a bounded number, depending only on  $V$ , of products of  $A$ .

It seems, then, that certain types of argument – in some sense those involving "bounded length" computations in the ambient group – adapt very well from the traditional group theory setting to approximate groups. At the moment we do not have anything approaching a precise formulation of this principle and indeed at present the passage from the "exact" to the approximate is as much an art as a science. Nonetheless, there seems to be merit in looking for "bounded length" proofs in traditional group theory which might be adapted to the approximate setting. Perhaps this is as good a place as any to mention the remarkable recent paper of Hrushovski [36] in which tools from model theory have been applied to the study of approximate groups. The ramifications of that paper are not yet completely clear, but it looks as though Theorem 1.3 of that paper together with some structure theory of algebraic groups ought to lead, without too much difficulty, to a proof of the following statement.

**Conjecture 4.3.** Suppose that  $A \subseteq \mathrm{SL}_n(\mathbb{C})$  is a  $K$ -approximate group. Then there is a  $K'$ -approximate group  $B$  which is *nilpotent* and  $K'$ -controls  $A$ , where  $K'$  depends only on  $K$ .

It seems reasonable to conjecture that  $K'$  can be taken to depend polynomially on  $K$ , although in their present form Hrushovski's techniques will not give this.

5. CAYLEY GRAPHS ON  $\mathrm{SL}_2(\mathbb{F}_p)$ 

We move on now to applications of the theory of approximate groups. In this section we discuss the paper [3] of Bourgain and Gamburd. This paper concerns *expander graphs*. For the purposes of this discussion these are  $2k$ -regular graphs  $\Gamma$  on  $n$  vertices for which there is a constant  $c > 0$  such that for any set  $X$  of at most  $n/2$  vertices of  $\Gamma$ , the number of vertices outside  $X$  which are adjacent to  $X$  is at least  $c|X|$ . Expander graphs share many of the properties of random regular graphs, and this is an important reason why they are of great interest in theoretical computer science. There are many excellent articles on expander graphs ranging from the very concise [51] to the seriously comprehensive [35].

A key issue is that of constructing explicit expander graphs, and in particular that of constructing *families* of expanders in which  $k$  and  $c$  are fixed but the number  $n$  of vertices tends to infinity. Many constructions have been given, and several of them arise from Cayley graphs. Let  $G$  be a finite group and suppose that  $S = \{g_1^{\pm 1}, \dots, g_k^{\pm 1}\}$  is a symmetric set of generators for  $G$ . The Cayley graph  $\mathcal{C}(G, S)$  is the  $2k$ -regular graph on vertex set  $G$  in which vertices  $x$  and  $y$  are joined if and only if  $xy^{-1} \in S$ . Such graphs provided some of the earliest examples of expanders [41, 42]. A natural way to obtain a family of such graphs is to take some large “mother” group  $\tilde{G}$  admitting many homomorphisms  $\pi$  from  $\tilde{G}$  to finite groups, a set  $\tilde{S} \subseteq \tilde{G}$ , and then to consider the family of Cayley graphs  $\mathcal{C}(\pi(\tilde{G}), \pi(\tilde{S}))$  as  $\pi$  ranges over a family of homomorphisms. The work under discussion concerns the case  $\tilde{G} = \mathrm{SL}_2(\mathbb{Z})$ , which of course admits homomorphisms  $\pi_p : \mathrm{SL}_2(\mathbb{Z}) \rightarrow \mathrm{SL}_2(\mathbb{F}_p)$  for each prime  $p$ . For certain sets  $\tilde{S} \subseteq \tilde{G}$ , for example

$$\tilde{S} = \left\{ \begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix}^{\pm 1}, \begin{pmatrix} 1 & 0 \\ 1 & 1 \end{pmatrix}^{\pm 1} \right\}$$

or

$$\tilde{S} = \left\{ \begin{pmatrix} 1 & 2 \\ 0 & 1 \end{pmatrix}^{\pm 1}, \begin{pmatrix} 1 & 0 \\ 2 & 1 \end{pmatrix}^{\pm 1} \right\},$$

spectral methods from the theory of automorphic forms may be used to show that  $(\mathcal{C}(\pi_p(\tilde{G}), \pi_p(\tilde{S})))_{p \text{ prime}}$  is a family of expanders. See [40] and the references therein. These methods depend on the fact that the group  $\langle \tilde{S} \rangle$  has finite index in  $\tilde{G} = \mathrm{SL}_2(\mathbb{Z})$  and they fail when this is not the case, for example when

$$(5.1) \quad \tilde{S} = \left\{ \begin{pmatrix} 1 & 3 \\ 0 & 1 \end{pmatrix}^{\pm 1}, \begin{pmatrix} 1 & 0 \\ 3 & 1 \end{pmatrix}^{\pm 1} \right\}.$$

In [40] Lubotzky asked whether the corresponding Cayley graphs in this and other cases might nonetheless form a family of expanders, the particular case of (5.1) being known as his “1-2-3 question”. The paper of Bourgain and Gamburd under discussion answers this quite comprehensively, showing that all that is required is that the group generated by  $\tilde{S}$  is not *virtually abelian* (contains a finite index abelian subgroup). We will sketch the proof in the case that  $\tilde{S}$  generates a nonabelian free subgroup of  $\mathrm{SL}_2(\mathbb{Z})$ . This is essentially the most general case, since the kernel of the natural homomorphism from  $\langle \tilde{S} \rangle$  to  $\mathrm{SL}_2(\mathbb{F}_2) \cong \mathrm{Sym}(3)$  is free and has index at most 6 in  $\langle \tilde{S} \rangle$ .

**Theorem 5.1** (Bourgain – Gamburd). *Let  $\tilde{G} = \mathrm{SL}_2(\mathbb{Z})$  as above and suppose that  $\tilde{S}$  is a finite symmetric set generating a free subgroup of  $\mathrm{SL}_2(\mathbb{Z})$ . Then*

$$(\mathcal{C}(\pi_p(\tilde{G}), \pi_p(\tilde{S})))_p \text{ prime}$$

*is a family of expanders.*

The notation we have introduced here is rather cumbersome, so let us write  $\Gamma_p := \mathcal{C}(\pi_p(\tilde{G}), \pi_p(\tilde{S}))$ . For concreteness we will focus on the special case  $\tilde{S} = \{A, A^{-1}, B, B^{-1}\}$ , where  $A = \begin{pmatrix} 1 & 3 \\ 0 & 1 \end{pmatrix}$  and  $B = \begin{pmatrix} 1 & 0 \\ 3 & 1 \end{pmatrix}$  are the matrices relevant to Lubotzky’s 1-2-3 question. The argument is almost identical in any other case. In this case, then,  $\Gamma_p$  is the graph on vertex set  $\mathrm{SL}_2(\mathbb{F}_p)$  in which  $x$  is joined to  $y$  if and only if  $xy^{-1}$  is one of the elements  $A, A^{-1}, B$  or  $B^{-1}$  considered modulo  $p$ . Supposing that  $p > 3$ , each of these graphs is 4-regular. The number of vertices in  $\Gamma_p$  is  $n := |\mathrm{SL}_2(\mathbb{F}_p)| = p(p^2 - 1)$ .

The reader may be interested to see a proof, using the “ping-pong” technique of Felix Klein, that the subgroup of  $\mathrm{SL}_2(\mathbb{Z})$  generated by these  $A$  and  $B$  is indeed free. Consider the natural action of  $A$  and  $B$  on the projective plane  $\mathbb{P}^1(\mathbb{Q})$ . Write

$$X := \{(\lambda : 1) \in \mathbb{P}^1(\mathbb{Q}) : |\lambda| < 1\}$$

and

$$Y := \{(1 : \lambda) \in \mathbb{P}^1(\mathbb{Q}) : |\lambda| < 1\},$$

and observe that  $X$  and  $Y$  are disjoint and *jovent au ping pong*, that is to say

$$A^n(X) \subseteq Y \quad \text{for all } n \in \mathbb{Z} \setminus \{0\}$$

and

$$B^n(Y) \subseteq X \quad \text{for all } n \in \mathbb{Z} \setminus \{0\}.$$

(The origin of the name should be clear – the “players”  $A$  and  $B$  hit the domains  $X$  and  $Y$  back and forth – as should the preference for the French term rather than the cumbersome “play table tennis with one another”.) If the group generated by  $A$  and  $B$  is not free, then some nontrivial reduced word in  $A$  and  $B$  is equal to the identity, where “reduced word” means a finite word of the form  $\dots A^{m_1} B^{n_1} \dots A^{m_k} B^{n_k} \dots$  with  $m_1, n_1, \dots, m_k, n_k \neq 0$ . The conjugate of such a word by an appropriate power of  $A$  will still be the identity and will now have the form  $w = A^{m_1} B^{n_1} \dots A^{m_k} B^{n_k} A^{m_{k+1}}$  with  $m_i, n_j \neq 0$ . However by repeated application of the ping-pong properties we see that  $w(X) \subseteq Y$ , certainly an impossibility since  $X$  and  $Y$  are disjoint and  $w$  is supposed to be the identity.

Following that slight digression let us focus once again on the Cayley graphs  $\Gamma_p$ , our aim being to prove that they form a family of expanders as  $p$  ranges over the primes. To do this we begin by giving a spectral interpretation of the expansion property which we defined combinatorially above. For each  $p$  we may consider the *Laplacian* of the corresponding Cayley graph, that is to say the operator

$$\Delta : L^2(\mathrm{SL}_2(\mathbb{F}_p)) \rightarrow L^2(\mathrm{SL}_2(\mathbb{F}_p))$$

defined by

$$\Delta f(x) := f(x) - \frac{1}{4}(f(Ax) + f(A^{-1}x) + f(Bx) + f(B^{-1}x)).$$

The eigenvalues of the Laplacian lie in the interval  $[0, 2]$ . Zero is certainly an eigenvalue, since  $\Delta 1 = 0$ . Write the eigenvalues in ascending order as  $0 = \lambda_0 \leq$



$\lambda_1 \leq \dots \leq \lambda_{n-1}$ . It turns out the expansion properties of the graph  $\Gamma_p$  (in fact of *any* regular graph) are intimately connected with the size of the second-smallest eigenvalue  $\lambda_1 = \lambda_1(\Gamma_p)$ . The precise relation between the combinatorial property of expansion and this spectral property is discussed in Section 2 of [35], but for our purposes we need only remark that it suffices to show that the second-smallest eigenvalue  $\lambda_1(\Gamma_p)$  is bounded away from zero independently of  $n$  (in fact, this is also a necessary condition for expansion). The term *spectral gap* is used to describe this property: there is a gap at the bottom of the spectrum in which there are no eigenvalues apart from zero.

To try to show that there is a spectral gap, consider the operator

$$T : L^2(\mathrm{SL}_2(\mathbb{F}_p)) \rightarrow L^2(\mathrm{SL}_2(\mathbb{F}_p))$$

given by  $T := 4(\mathrm{id} - \Delta)$ , that is to say

$$Tf(x) := f(Ax) + f(A^{-1}x) + f(Bx) + f(B^{-1}x).$$

The matrix<sup>13</sup> of  $T$  is same thing as the adjacency matrix of the graph  $\Gamma_p$ , that is to say the matrix whose  $xy$  entry is 1 if  $x \sim y$  and zero otherwise. The eigenvalues of  $T$  are of course  $\mu_i = 4(1 - \lambda_i)$ ,  $i = 0, \dots, n-1$ , and it is a very well-known and easy to establish fact that the  $2m$ th moment  $\sum_{i=0}^{n-1} \mu_i^{2m}$  is equal to  $n$  times  $W_{2m}$ , the number of closed walks of length  $2m$  from the identity to itself. It follows that we have

$$(5.2) \quad W_{2m} = \frac{1}{n} 4^{2m} \left( 1 + \sum_{i=1}^{n-1} (1 - \lambda_i)^{2m} \right).$$

Note in particular that  $W_{2m} \geq \frac{1}{n} 4^{2m}$ , since all the terms are non-negative. At first glance it looks as though the only way to use (5.2) to bound  $\lambda_1$  away from zero would be to get rather precise estimates on  $W_{2m}$ , and in particular one would at the very least want to show that  $W_{2m} < \frac{2}{n} 4^{2m}$ . However a remarkable observation, used earlier in related contexts by Sarnak and Xue [54] and Gamburd [21], comes into play. This is that any eigenspace of the Laplacian is  $\mathrm{SL}_2(\mathbb{F}_p)$ -invariant, where the action of  $\mathrm{SL}_2(\mathbb{F}_p)$  on  $L^2(\mathrm{SL}_2(\mathbb{F}_p))$  is the right-regular one given by  $g \circ f(x) := f(xg)$ . In other words, any such eigenspace has the structure of a representation of  $\mathrm{SL}_2(\mathbb{F}_p)$  and thus, by basic representation theory, decomposes into irreducible representation of  $\mathrm{SL}_2(\mathbb{F}_p)$ . But by a classical theorem of Frobenius all such representations have dimension at least  $(p-1)/2 \sim n^{1/3}$ . This means that  $\lambda_1 = \lambda_2 = \dots = \lambda_l$  for some  $l \sim n^{1/3}$ , and hence from (5.2) we in fact have the bound

$$(5.3) \quad W_{2m} \gg \frac{1}{n^{2/3}} 4^{2m} (1 - \lambda_1)^{2m}.$$

This enables a meaningful spectral gap (lower bound on  $\lambda_1$ ) to be obtained from somewhat weaker upper bounds on  $W_{2m}$ .

The main new content of [3], then, is to obtain those upper bounds on  $W_{2m}$ , the number of walks of length  $2m$  starting and finishing at the identity, for appropriate  $m$ . A nice way of thinking about these walks is in terms of *convolution powers* of the probability measure

$$\nu := \frac{1}{4}(\delta_A + \delta_{A^{-1}} + \delta_B + \delta_{B^{-1}})$$

---

<sup>13</sup>With respect to the basis of  $\mathrm{SL}_2(\mathbb{F}_p)$  consisting of the functions  $\mathbf{1}_t : \mathrm{SL}_2(\mathbb{F}_p) \rightarrow \mathbb{C}$  defined by  $\mathbf{1}_t(x) = 1$  if  $x = t$  and 0 otherwise.

on  $\mathrm{SL}_2(\mathbb{F}_p)$ , where  $\delta_g(x) = n$  if  $x = g$  and 0 otherwise. This measure  $\nu$  is a very singular or “spiky” probability measure, supported on just the four points  $A, A^{-1}, B$  and  $B^{-1}$ . Now the convolution

$$\nu^{(2)} := \nu * \nu(x) := \mathbb{E}_{y \in \mathrm{SL}_2(\mathbb{F}_p)} \nu(xy^{-1})\nu(y)$$

is supported on words of length at most two in  $A, A^{-1}, B$  and  $B^{-1}$ , or alternatively those  $x$  in the graph  $\Gamma_p$  which can be reached from the identity by a path of length two, the value of  $\nu * \nu(x)$  being  $4^{-2}n$  times the number of paths of length two from the identity to  $x$ . Similarly higher convolution powers  $\nu^{(m)}(x) := \nu * \dots * \nu(x)$  give  $4^{-m}n$  times the number of paths of length  $m$  from the identity to  $x$ . The idea of the proof is to examine these convolution powers, showing that they become progressively more “spread out” until, for suitable  $m$ ,  $\nu^{(2m)}$  vaguely resembles the uniform measure  $\mathbf{1}$  which assigns weight one to each point of  $\mathrm{SL}_2(\mathbb{F}_p)$ . Then, in particular,  $\nu^{(2m)}(0) \sim 1$ , meaning that  $W_{2m} \sim 4^{2m}/n$ . Combined with (5.3), this is enough to establish the desired spectral gap.

The notion of a probability measure  $\mu$  on  $\mathrm{SL}_2(\mathbb{F}_p)$  being “spread out” may be quantified using the  $L^2$ -norm

$$\|\mu\|_2 := \left( \mathbb{E}_{x \in \mathrm{SL}_2(\mathbb{F}_p)} \mu(x)^2 \right)^{1/2}.$$

The  $L^2$ -norm of a delta measure  $\delta_g$  is  $n^{1/2}$ , which is huge, whilst that of the uniform measure  $\mathbf{1}$  is equal to one, the smallest value possible by the Cauchy-Schwarz inequality. It is not hard to show that convolution cannot increase the  $L^2$ -norm, and so we have the chain of inequalities

$$(5.4) \quad n^{1/2} = \|\nu^{(1)}\|_2 \geq \|\nu^{(2)}\|_2 \geq \dots$$

The aim is to show that this sequence is, in fact, rather rapidly decreasing. Roughly speaking one shows that

$$(5.5) \quad \|\nu^{(m_1)}\|_2 \approx 1$$

for some  $m_1 \approx C_1 \log p$ ; this  $m_1$  turns out to be an appropriate choice to substitute into (5.3) in order to reach the desired conclusions.

It turns out that this sequence gets off to a rather good start. This is a consequence of an observation of Margulis [43], namely that the freeness of the subgroup of  $\mathrm{SL}_2(\mathbb{Z})$  generated by  $A$  and  $B$  persists to some extent even when reduced modulo  $p$ . Indeed let us take a reduced word  $w = A^{m_1} B^{n_1} \dots A^{m_k} B^{n_k}$  with  $m_1, \dots, m_k, n_1, \dots, n_k \neq 0$  and suppose that this equals the identity when reduced modulo  $p$ , that is to say in  $\mathrm{SL}_2(\mathbb{F}_p)$ . Lifting back up to  $\mathrm{SL}_2(\mathbb{Z})$  we have

$$\tilde{w} = A^{m_1} B^{n_1} \dots A^{m_k} B^{n_k} \equiv \mathrm{id} \pmod{p}.$$

But the freeness of the lifted group means that  $\tilde{w} \neq \mathrm{id}$ , and thus in order to be congruent to the identity mod  $p$  the matrix  $\tilde{w}$  must have at least one entry of size at least  $p - 1$ . But by some simple matrix inequalities this is impossible provided that

$$|m_1| + |n_1| + \dots + |m_k| + |n_k| < c \log p$$

for some absolute constant  $c > 0$ .

It follows that the subgroup of  $\mathrm{SL}_2(\mathbb{F}_p)$  generated by  $A$  and  $B$  is “free up to words of length  $c \log p$ ”. In terms of the Cayley graphs  $\Gamma_p$  this means that up to retracing steps there is a unique walk of length  $m$  between the identity and  $x$  for

any  $x \in \mathrm{SL}_2(\mathbb{F}_p)$  and for any  $m \leq m_0 := c \log p/2$ . This implies that the measures  $\nu^{(m)}$ ,  $m \geq m_0$  are already rather spread out. To quantify this (and in particular to deal with the issue of “retracing steps”) a result of Kesten concerning random walks in the free group may be applied. The conclusion is that

$$(5.6) \quad \|\nu^{(m_0)}\|_2 \ll n^{1/2-\gamma}$$

for some  $\gamma > 0$ . This is good progress on the way to (5.5) and represents a significant improvement on the initial bound  $\|\nu^{(1)}\|_2 = n^{1/2}$ .

It is convenient to imagine, for the rest of the argument, that all probability measures  $\mu$  on  $G$  have the form  $\mu(x) = \frac{n}{|A|} 1_A(x)$  for some set  $A \subseteq G$ , the “support” of  $\mu$ . Whilst this is clearly not true, various (somewhat technical) decompositions into level sets may be used to reduce to this case. For such a measure we have

$$\|\mu\|_2 = (n/|A|)^{1/2},$$

and so the bound (5.6) corresponds to  $|A| \gg n^{2\gamma}$ , certainly a reasonable level of spreadoutness.

The rest of the argument, which constitutes the heart of the paper, involves examining the convolution powers between  $\nu^{(m_0)}$  and  $\nu^{(m_1)}$  for a suitable  $m_1 \sim C_1 \log p$ , the aim being to establish (5.5). An application of the “dyadic pigeonholing argument”, used to great effect by Bourgain in many papers, is employed: if  $\|\nu^{(m_1)}\|_2$  is much larger than 1, this means that the sequence (5.4) cannot decay too rapidly between  $\nu^{(m_0)}$  and  $\nu^{(m_1)}$  and so there must be two convolution powers  $\nu^{(m)}$  and  $\nu^{(2m)}$ ,  $m_0 \leq m < m_1$ , such that  $\|\nu^{(2m)}\|_2 \approx \|\nu^{(m)}\|_2$ . Let us be deliberately vague about the meaning of  $\approx$  here.

Suppose that  $\nu^{(m)}(x) = \frac{n}{|A|} 1_A(x)$  for some set  $A \subseteq G$ . Noting that  $\nu^{(2m)} = \nu^{(m)} * \nu^{(m)}$ , it is not hard to compute that the ratio

$$\|\nu^{(2m)}\|_2^2 / \|\nu^{(m)}\|_2^2$$

is actually equal to  $|A|^{-3}$  times the number of quadruples  $a_1, a_2, a_3, a_4 \in A^4$  with  $a_1 a_2 = a_3 a_4$ . This may be compared with condition (4) in the list of properties which are known to roughly characterise approximate groups. Thus, being even rougher at this point,

$$(5.7) \quad \nu^{(m)} \sim \frac{1}{H} 1_H$$

for some approximate group  $H \subseteq \mathrm{SL}_2(\mathbb{F}_p)$ . Note that the rough equivalence of (4) and other, more flexible definitions such as Definition 2.4 is one of the deeper equivalences mentioned in §2, being reliant on the nonabelian Balog-Szemerédi-Gowers theorem of Tao [59].

If  $H$  is already all of  $\mathrm{SL}_2(\mathbb{F}_p)$  then (5.7) is telling us that  $\nu^{(m)}$  is close to the uniform distribution, in which case so is  $\nu^{(m_1)}$ , hence (5.5) is established and we are done. If not then we apply Helfgott’s result, Theorem 4.1, to conclude that  $H$  is essentially upper-triangular, and hence that  $\nu^{(m_0)}$  has significant mass on an upper-triangular subgroup of  $\mathrm{SL}_2(\mathbb{F}_p)$ .

The support of  $\nu^{(m_0)}$ , however, consists of words of length at most  $m_0$  in the generators  $A, A^{-1}, B$  and  $B^{-1}$  and, as we stated, these elements behave freely up to words of this length. This is highly incompatible with upper-triangularity, which

in particular implies that we always have the commutator relation<sup>14</sup>

$$(5.8) \quad [[g_1, g_2], [g_3, g_4]] = \text{id}.$$

A pleasant group-theoretic argument formalises this incompatibility and allows one to show that any set of words of length at most  $m_0$  in the generators  $A, A^{-1}, B$  and  $B^{-1}$  satisfying (5.8) has size at most  $m_0^6$ . This represents a tiny proportion of the set of all such words, which (counted with multiplicity at least) has cardinality  $4^{m_0}$ . This contradiction finishes the sketch proof of Theorem 5.1.  $\square$

Before moving on, we wish to record, for use in the next section, a further observation concerning the measures  $\nu^{(m)}$ . We sketched a proof that  $\|\nu^{(m_1)}\|_2 \approx 1$  for some  $m_1 \sim C_1 \log p$ , that is to say  $\nu^{(m_1)}$  vaguely resembles the uniform distribution on  $\text{SL}_2(\mathbb{F}_p)$ . By taking further convolutions and using the fact that irreducible representations have large degree once more, this may be bootstrapped to show that  $\nu^{(m)}$  becomes exponentially well uniformly-distributed:

$$(5.9) \quad \nu^{(m)}(x) = 1 + O(ne^{-cm})$$

for some absolute  $c > 0$  and for all  $m$ . Alternatively, such a statement can be deduced directly from the spectral gap property, as is done for example in [6, §3.3].

It is interesting to ask whether the arguments might adapt to deal with Cayley graphs on  $\text{SL}_n(\mathbb{F}_p)$  with  $n \geq 3$ . A recent paper of Bourgain and Gamburd [5] shows that this is the case when  $n = 3$ . The argument is, in large part, quite similar to the above, except of course that Helfgott’s theorem on approximate subgroups of  $\text{SL}_2(\mathbb{F}_p)$  must be replaced by his more difficult result [34] on approximate subgroups of  $\text{SL}_3(\mathbb{F}_p)$ . There is one significant extra difficulty, however, which is that there are proper subgroups of  $\text{SL}_3(\mathbb{F}_p)$  which are not close to upper-triangular, an obvious example being a copy of  $\text{SL}_2(\mathbb{F}_p)$ . To deal with this a deep algebro-geometric result of Nori [45] is brought into play, which states that any proper subgroup of  $\text{SL}_3(\mathbb{F}_p)$ ,  $p$  sufficiently large, must satisfy a non-trivial polynomial equation. To obtain a contradiction, it must be shown that the set of words of length  $m_0$  in the generators  $A$  and  $B$  (say) does not concentrate on the corresponding subvariety of  $\text{SL}_3(\mathbb{C})$ , and here techniques from the theory of random matrix products and a certain amount of “quantitative algebraic geometry” are brought into play.

## 6. NONLINEAR SIEVING PROBLEMS

In this section we discuss work of Bourgain, Gamburd and Sarnak [6]. The goal of *sieve theory*, traditionally viewed as a part of analytic number theory, is to find prime numbers or at least to say something about them. Historically, the sieve arose through work of Brun and Merlin on the twin prime problem, that is to say the problem of finding infinitely many primes  $p$  such that  $p + 2$  is also prime. Whilst this remains a famous open problem, approximations to it have been found. For example, Brun established the following result.

**Theorem 6.1** (Brun). *There are infinitely many integers  $n$  such that  $n(n + 2)$  has at most 9 prime factors.*

---

<sup>14</sup>In other words, upper-triangular subgroups of  $\text{SL}_2(\mathbb{F}_p)$  are 2-step solvable.

Much later, Chen [14] replaced 9 by 3. One way of stating this type of result is as follows: there are infinitely many  $n$  for which both  $n(n+2)$  is a 3-*almost prime*, that is to say a positive integer with at most 3 prime factors.

The aim of [6] is to discover almost primes in more exotic locales, and specifically in *orbits* of linear groups. We will sketch a proof of the following result.

**Theorem 6.2** (Bourgain-Gamburd-Sarnak). *Let  $A$  and  $B$  be two matrices in  $\mathrm{SL}_2(\mathbb{Z})$  generating a free subgroup. Then there is some  $r$  such that this group contains infinitely many  $r$ -almost prime matrices (matrices, the product of whose entries is  $r$ -almost prime).*

Henceforth we shall say “almost prime” instead of “ $r$ -almost prime for some  $r$ ”. We remark that in the specific case we focussed on in the last section, when  $A = \begin{pmatrix} 1 & 3 \\ 0 & 1 \end{pmatrix}$  and  $B = \begin{pmatrix} 1 & 0 \\ 3 & 1 \end{pmatrix}$ , the theorem as stated follows from classical sieve theory of the type used to prove Brun’s theorem. Indeed (for example)  $A^n B A = \begin{pmatrix} 9n+1 & 30n+3 \\ 3 & 10 \end{pmatrix}$ , and the product of the entries here is  $2 \cdot 3^2 \cdot 5 \cdot (9n+1) \cdot (10n+1)$ , which will be almost prime for infinitely many  $n$  by a simple variant of Brun’s analysis. The issue here is that the subgroup generated by  $A$  and  $B$  contains unipotent elements (in this case both  $A$  and  $B$  are themselves unipotent).

We start with a (very) elementary discussion of what a sieve is. Suppose one has a finite set  $X$  of integers and that one wishes to find primes or almost primes in  $X$ . The most naïve way to do this would be to try to adapt the sieve of Eratosthenes, using the inclusion-exclusion principle to compute

$$\#\{\text{primes in } X\} = |X| - |X_2| - |X_3| - |X_5| - \dots + |X_6| + |X_{10}| + |X_{15}| + \dots - |X_{30}| - \dots$$

where  $X_q$  is the set of elements of  $X$  which are divisible by  $q$ . Unfortunately it is well-known that, even when  $X$  is an extremely simple set such as  $\{1, \dots, n\}$ , it is not generally possible to evaluate  $|X_q|$  sufficiently accurately to avoid the error terms in this long sum blowing up. In this simple case just mentioned, for example, we have  $|X_q| = \lfloor n/q \rfloor$ . However the floor function is rather unpleasant and it is tempting to write instead  $|X_q| = n/q + O(1)$ , but then one finds that there are so many  $O(1)$  errors that the sieve of Eratosthenes becomes useless.

By and large, *sieve theory* is concerned with what it is possible to say about primes or almost primes in  $X$  given “reasonably nice” information about the size of the sets  $X_q$ . Although the sieve of Eratosthenes is bad, other sieves fare rather better. These other sieves are generally cleverly weighted versions of the sieve of Eratosthenes, but we will not dwell upon their construction here. A typical example of “reasonably nice” information about  $|X_q|$  would be

$$|X_q| = \beta(q)|X| + r_q$$

for all squarefree  $q \leq |X|^\gamma$ , where  $\beta(q)$  is some pleasant multiplicative function and the error  $r_q$  is small in the sense that  $|r_q| \ll |X|^{1-\gamma}$  for some  $\gamma > 0$ . For example, if  $X = \{1, \dots, n\}$  then this is true with  $\beta(q) = 1/q$  and for any  $\gamma \leq 1$ .

The fundamental theorem of the combinatorial sieve states, roughly speaking, that such information is enough to find almost primes in  $X$ ; in fact, one can even estimate the number of almost primes. What is meant by “almost prime” – that is, how many prime factors these numbers will have – depends on how large we can

take  $\gamma$  as well as on the so-called dimension of the sieve, which has to do with the average size of the quantities  $\beta$ . We will not delay ourselves by expanding upon the details here. Let us instead refer the reader to [6] for the precise formulation convenient to the application there and to the book [38] or the unpublished notes [37] for a more wide-ranging discussion of sieves in general with full proofs.

All we shall take from the preceding discussion is the notion that, given a finite set  $X$  to be sieved in order to locate almost primes, we should be looking for good asymptotics for the size of the sets  $|X_q|$ ,  $q$  squarefree. Returning to Theorem 6.2, the first obvious question to answer is that of what the set  $X$  to be sieved should be. The set in which we wish to find almost primes is

$$\mathcal{A} := \{x_1 x_2 x_3 x_4 : \begin{pmatrix} x_1 & x_2 \\ x_3 & x_4 \end{pmatrix} \in \langle A, B \rangle\}.$$

Now  $\mathcal{A}$  is of course an infinite set of integers. Rather than truncate in the usual way and take  $X = \mathcal{A} \cap \{1, \dots, N\}$ , it is much more natural to truncate in a manner that respects the group structure more. This we do by taking

$$X := \{x_1 x_2 x_3 x_4 : \begin{pmatrix} x_1 & x_2 \\ x_3 & x_4 \end{pmatrix} \in \Sigma_m(A, B)\},$$

where

$$\Sigma_m(A, B) = \{U_1 U_2 \dots U_m : U_i \in \{A, A^{-1}, B, B^{-1}\}\}$$

is the set of words of length  $m$  in  $A, A^{-1}, B$  and  $B^{-1}$  and  $X$  is counted with multiplicity so that  $|X| = 4^m$ .

Suppose that  $p$  is a prime. Then  $|X_p|$  is equal to the number of words  $w \in \Sigma_m(A, B)$ , counted with multiplicity, which, when reduced modulo  $p$ , give rise to a matrix in  $\mathrm{SL}_2(\mathbb{F}_p)$  with at least one zero entry. Writing  $S \subseteq \mathrm{SL}_2(\mathbb{F}_p)$  for the set of such matrices, it is easy to compute that  $|S| = 2(2p-1)(p-1)$ . Now the number of words  $w \in \Sigma_m(A, B)$  which reduce modulo  $p$  to some  $x \in \mathrm{SL}_2(\mathbb{F}_p)$  is, in the notation of the last section, precisely  $\frac{1}{n}|X|\nu^{(m)}(x)$ , and so

$$|X_p| = \frac{1}{n}|X| \sum_{x \in S} \nu^{(m)}(x).$$

However at the end of the last section we saw<sup>15</sup> that  $\nu^{(m)}(x)$  becomes very close to 1. In fact, in (5.9) we noted the bound  $\nu^{(m)}(x) = 1 + O(ne^{-cm})$ . Using this we obtain

$$|X_p| = \beta(p)|X| + r_p$$

where  $\beta(p) := 2(2p-1)/p(p+1)$  and  $|r_p| = |X|^{1-\gamma}$  for some  $\gamma > 0$ .

Thus the expansion property of the Cayley graphs  $(\mathcal{C}(\pi_p(\tilde{G}), \pi_p(\tilde{S})))_{p \text{ prime}}$  gives exactly the kind of information that can be input into the combinatorial sieve!

There is, however, a very major caveat. What we have just said applies only to  $X_p$  when  $p$  is a prime, and for the sieve one must understand  $X_q$  when  $q$  is a general squarefree number. To do this requires the establishment of Theorem 5.1 for the family  $(\mathcal{C}(\pi_q(\tilde{G}), \pi_q(\tilde{S})))_q$ , where now  $q$  ranges over all squarefrees and not just over primes. The broad scheme of the proof is the same, but every single ingredient must be generalised to the more general setting, starting from the classification of

<sup>15</sup>Either as a byproduct of the proof, or a consequence, of the expansion property of the family of Cayley graphs  $\Gamma_p = \mathcal{C}(\pi_p(\mathcal{G}), \pi_p(\tilde{S}))$ .



FIGURE 1. Apollonian circle packing

approximate subrings of  $\mathbb{Z}/q\mathbb{Z}$ . The situation here is more complicated because this ring will in general have many approximate subrings, namely  $\mathbb{Z}/q'\mathbb{Z}$  with  $q'|q$ . One of the main technical results of [6] (occupying some 20 pages) is the statement that, very roughly speaking, these are the only approximate subrings of  $\mathbb{Z}/q\mathbb{Z}$ . Although this is a deeply technical argument of a type that this author would struggle to summarise meaningfully even to an expert audience, it might be compared with the 92-page proof [2] of the corresponding assertion without the squarefree assumption on  $q$ . Thankfully<sup>16</sup> this is not required for the present application. Once the classification of approximate subrings of  $\mathbb{Z}/q\mathbb{Z}$  for  $q$  squarefree is in place a suitable analogue of Helfgott's argument is applied to roughly classify approximate subgroups of  $\mathrm{SL}_2(\mathbb{Z}/q\mathbb{Z})$ . Even the statement of this result (Proposition 4.3 in the paper) is rather technical. Finally, the majority of the argument outlined in the last section in the case  $q$  prime goes over without substantial change.

This concludes our discussion of the proof of Theorem 6.2. To conclude this survey, we wish to mention a beautiful application, mentioned in the original paper [6] and in other articles such as [52], of these nonlinear sieving ideas. This has to do with Apollonian packings such as the one in the attractive image above.

---

<sup>16</sup>This is one of the most extraordinarily long and technical arguments the author has ever seen. The theory of approximate rings when there are many zero-divisors seems to be very difficult.

For a very pleasant and gentle introduction to Apollonian packings, see [1]. Referring to Figure 1, inside each circle is an integer which represents the curvature of that circle, or in other words the reciprocal of the radius. Some of the number theory associated with the integers that arise in this way is discussed in the letter [53] where, for example, it is shown that infinitely many of these curvatures are prime and in fact that there are infinitely many touching pairs of circles with prime curvature.

Now a pleasant exercise in Euclidean geometry gives a theorem of Descartes, namely that the relation between the four integers  $a_1, a_2, a_3, a_4$  inside four mutually touching circles is given by

$$(6.1) \quad 2(a_1^2 + a_2^2 + a_3^2 + a_4^2) = (a_1 + a_2 + a_3 + a_4)^2.$$

Examples of quadruples  $(a_1, a_2, a_3, a_4)$  which are related in this way and easily visible in the picture are  $(13, 21, 24, 124)$  and  $(13, 24, 37, 156)$ .

Take a quadruple  $(C_1, C_2, C_3, C_4)$  of touching circles with curvatures

$$(a_1, a_2, a_3, a_4) = (13, 21, 24, 124).$$

There is another circle  $C'_1$  tangent to  $C_2, C_3$  and  $C_4$ , and it has curvature  $a'_1 = 325$ . To find a general relation between  $a_1$  and  $a'_1$  we may note that  $a_1, a'_1$  are roots of (6.1) regarded as a quadratic in  $a_1$  and thereby obtain the relation

$$a'_1 = -a_1 + 2a_2 + 2a_3 + 2a_4.$$

This may of course be written as

$$\begin{pmatrix} a'_1 \\ a_2 \\ a_3 \\ a_4 \end{pmatrix} = \begin{pmatrix} -1 & 2 & 2 & 2 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} a_1 \\ a_2 \\ a_3 \\ a_4 \end{pmatrix}.$$

That is, if one starts with some fixed vector such as  $x_0 = (13, 21, 24, 124)$  then one may obtain another quadruple of curvatures of circles in the Apollonian packing by applying the matrix

$$S_1 := \begin{pmatrix} -1 & 2 & 2 & 2 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}.$$

By playing the same game with  $C'_2, C'_3$  and  $C'_4$  we can make the same assertion with the matrices

$$S_2 := \begin{pmatrix} 1 & 0 & 0 & 0 \\ 2 & -1 & 2 & 2 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix},$$

$$S_3 := \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 2 & 2 & -1 & 2 \\ 0 & 0 & 0 & 1 \end{pmatrix}$$



and

$$S_4 := \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 2 & 2 & 2 & -1 \end{pmatrix}.$$

This leads naturally to consideration of the orbit  $\langle \tilde{S} \rangle x_0 \subseteq \mathbb{Z}^4$ , where

$$\tilde{S} := \{S_1, S_2, S_3, S_4\},$$

every vector in which consists entirely of curvatures of circles in the Apollonian packing. This puts us in a situation very similar to that studied in Theorem 6.2, except now we appear to be dealing with a subgroup of  $\mathrm{GL}_4(\mathbb{Z})$  rather than of  $\mathrm{SL}_2(\mathbb{Z})$ .

It turns out, however, that this situation is essentially a two-dimensional one in disguise, and for this we need to add to the list of areas of mathematics we touch upon by hinting at Lie theory and special relativity! The matrices  $S_1, S_2, S_3, S_4$  belong to  $\mathrm{SO}_F(\mathbb{Z})$ , the subgroup of  $\mathrm{GL}_4(\mathbb{Z})$  consisting of  $4 \times 4$  matrices with determinant one which preserve the quadratic form  $F(\vec{x}) = 2(x_1^2 + x_2^2 + x_3^2 + x_4^2) - (x_1 + x_2 + x_3 + x_4)^2$  (cf. (6.1)). By the standard theory of quadratic forms (over  $\mathbb{R}$ ) this is equivalent to the Lorentz form  $L(\vec{y}) = y_1^2 + y_2^2 + y_3^2 - y_4^2$ , and so we may identify  $\mathrm{SO}_F(\mathbb{R})$  with the orthogonal group  $\mathrm{SO}(3, 1)$  preserving this latter form. But it is very well-known that this group admits  $\mathrm{SU}(2)$  as a double cover: this is because the set  $\{\vec{y} : L(\vec{y}) = -1\}$  may be identified with the set of  $2 \times 2$  hermitian matrices  $M$  with determinant 1 via

$$(y_1, y_2, y_3, y_4) \mapsto \begin{pmatrix} y_4 + y_3 & y_1 - iy_2 \\ y_1 + iy_2 & y_4 - y_3 \end{pmatrix},$$

and so any  $P \in \mathrm{SU}(2)$  gives rise to an element of  $\mathrm{SO}(3, 1)$  via the transformation  $M \mapsto PMP^*$ .

By lifting to this double cover the group  $\langle \tilde{S} \rangle$  can be lifted to a subgroup of  $\mathrm{SL}_2(\mathbb{Z}[i])$ . The proof of Theorem 6.2 goes through with relatively minimal changes, although once again the group generated by  $S_1, S_2, S_3$  and  $S_4$  contains unipotents and so, if the aim is simply to find infinitely many circles or pairs/quadruples of touching circles with almost-prime curvatures, more elementary approaches work just as well. Those elementary approaches do not, however, give sharp quantitative results, whereas the techniques we have sketched do. To explain one such result, imagine Figure 1 being generated as follows. Start with the outer circle (which has curvature  $-6$ ) and the three largest inner circles, with curvatures 13, 21 and 24. This is the *first generation*. The second generation consists of those circles touching three from the first generation: they have curvatures 28, 37, 61 and 124. The third generation contains those new circles touching three circles from either the first or the second generations: these have curvatures 45, 60, 69, 93, 124, 132, 133, 156, 220, 292, 301 and 325. Carry on in this vein: the  $n$ th generation will contain  $4 \cdot 3^{n-2}$  circles.

**Theorem 6.3** (Bourgain, Gamburd, Sarnak). *The number of circles at generation  $n$  which have prime curvature is bounded by  $C3^n/n$ , for some absolute constant  $C$ .*

We conclude by remarking that there are some very interesting unsolved questions connected with Apollonian packings [24]. In that paper the very interesting

question is raised of whether, in Figure 1, a positive proportion of all positive integers appear as curvatures. J. Bourgain has recently indicated to me that he and Elena Fuchs have obtained new information on this question. See also [39] for an asymptotic formula for the number of circles in the packing of curvature at most  $X$ . It seems that the question of describing this set of integers more precisely remains open: are they given, from some point on, by finitely many congruence conditions?

## 7. ACKNOWLEDGEMENTS

I am very grateful to Cliff Reiter for his permission to use the image of an Apollonian packing. In addition to the cited references, the notes from the fourth Marker Lecture of Terence Tao were very useful in preparing what is written here. I thank Jean Bourgain, Emmanuel Breuillard, Alex Gamburd, Harald Helfgott, Alex Kontorovich, Peter Sarnak and Terry Tao for their comments on an earlier draft of this survey. Finally, I am very grateful to Lilian Matthiesen and Vicky Neale for correcting a number of typographical errors.

## REFERENCES

1. D. Austin, *When kissing involves trigonometry*, AMS monthly feature article, available at <http://www.ams.org/featurecolumn/archive/kissing.html>.
2. J. Bourgain, *The sum-product theorem in  $\mathbb{Z}_q$  with  $q$  arbitrary*, J. Anal. Math. **106** (2008), 1–93.
3. J. Bourgain and A. Gamburd, *Uniform expansion bounds for Cayley graphs of  $\mathrm{SL}_2(\mathbb{F}_p)$* , Ann. Math. **167** (2008), 625–642.
4. ———, *Expansion and random walks in  $\mathrm{SL}_d(\mathbb{Z}/p^n\mathbb{Z})$ . I*, J. Eur. Math. Soc. (JEMS) **10** (2008), no. 4, 987–1011.
5. ———, *Expansion and random walks in  $\mathrm{SL}_d(\mathbb{Z}/p^n\mathbb{Z})$ , II*, with an appendix by J. Bourgain, J. Eur. Math. Soc. (JEMS) **11** (2009), 1057–1103.
6. J. Bourgain, A. Gamburd and P. Sarnak, *Affine linear sieve, expanders, and sum-product*, preprint.
7. J. Bourgain, A. A. Glibichuk and S. V. Konyagin, *Estimates for the number of sums and products and for exponential sums in fields of prime order*, J. London Math. Soc. (2) **73** (2006), no. 2, 380–398.
8. J. Bourgain, N. H. Katz and T. C. Tao, *A sum-product estimate in finite fields and applications*, Geom. Funct. Anal. (GAFA) **14** (2004), 27–57.
9. J. Bourgain and S. Konyagin, *Estimates for the number of sums and products and for exponential sums over subgroups in fields of prime order*, C. R. Acad. Sci. Paris, Ser. I, **337** (2003), 75–80.
10. E. Breuillard and B. J. Green, *Approximate groups, I: the torsion-free nilpotent case*, preprint available at <http://arxiv.org/abs/0906.3598>.
11. ———, *Approximate groups, II: the solvable linear case*, preprint available at <http://arxiv.org/abs/0907.0927>.
12. ———, *Approximate groups, III: the bounded linear case*, in preparation.
13. M.-C. Chang, *Product theorems in  $\mathrm{SL}_2$  and  $\mathrm{SL}_3$* , J. Inst. Math. Jussieu **7** (2008), no. 1, 1–25.
14. J. R. Chen, *On the representation of a larger even integer as the sum of a prime and the product of at most two primes*, Sci. Sinica **16** (1973), 157–176.
15. G. A. Edgar and C. Miller, *Borel subsets of the reals*, Proc. Amer. Math. Soc. **131** (2003), 1121–1129.
16. G. Elekes and Z. Király, *On the combinatorics of projective mappings*, J. Algebraic Combin. **14** (2001), no. 3, 183–197.
17. P. Erdős and E. Szemerédi, *On sums and products of integers*, Studies in pure mathematics, 213–218, Birkhäuser, Basel, 1983.
18. A. Eskin, S. Mozes and H. Oh, *On uniform exponential growth for linear groups*, Invent. Math. **160** (2005), no. 1, 1–30.

19. D. Fisher, N. H. Katz and I. Peng, *On Freiman's theorem in nilpotent groups*, preprint available at <http://arxiv.org/abs/0901.1409>.
20. G. A. Freiman, *Foundations of a structural theory of set addition*, Translations of Mathematical Monographs, **37**. American Mathematical Society, Providence, R. I., 1973. vii+108 pp.
21. A. Gamburd, *On the spectral gap for infinite index "congruence" subgroups of  $SL_2(\mathbb{Z})$* , Israel J. Math. **127** (2002), 157–200.
22. W. T. Gowers, *A new proof of Szemerédi's theorem for arithmetic progressions of length four*, Geom. Funct. Anal. **8** (1998), no. 3, 529–551.
23. ———, *Rough structure and classification*, GAFA 2000 (Tel Aviv, 1999). Geom. Funct. Anal. 2000, Special Volume, Part I, 79–117.
24. R. Graham, J. Lagarias, C. Mallows, A. Wilks and C. Yan, *Apollonian circle packings: number theory*, J. Number Theory **100** (2003), no. 1, 1–45.
25. B. J. Green, *Finite field models in additive combinatorics*, Surveys in Combinatorics 2005, London Math. Soc. Lecture Notes **327**, 1–27.
26. ———, *Additive combinatorics: review of the book of Tao and Vu*, Bull. Amer. Math. Soc. (N.S.) **46** (2009), no. 3, 489–497.
27. ———, *The Polynomial Freiman-Ruzsa conjecture*, guest blog at Terry Tao's weblog.
28. ———, notes from a 2009 Cambridge Part III course on Additive Combinatorics, Chapter 2, available on the author's webpage.
29. ———, *Approximate structure in additive combinatorics: Barbados-Radcliffe Lectures*, book in preparation.
30. B. J. Green and I. Z. Ruzsa, *Freiman's theorem in an arbitrary abelian group*, J. Lond. Math. Soc. (2) **75** (2007), no. 1, 163–175.
31. B. J. Green and T. C. Tao, *An equivalence between inverse sumset theorems and inverse conjectures for the  $U^3$ -norm*, preprint available at <http://arxiv.org/abs/0906.3100>.
32. M. Gromov, *Groups of polynomial growth and expanding maps*, Inst. Hautes Études Sci. Publ. Math. No. **53** (1981), 53–73.
33. H. A. Helfgott, *Growth and generation in  $SL_2(\mathbb{Z}/p\mathbb{Z})$* , Ann. of Math. (2) **167** (2008), no. 2, 601–623.
34. ———, *Growth in  $SL_3(\mathbb{Z}/p\mathbb{Z})$* , to appear in J. European Math. Soc.
35. S. Hoory, N. Linial and A. Wigderson, *Expander graphs and their applications*, Bull. Amer. Math. Soc. **43**, no. 4 (2006), 439–561.
36. E. Hrushovski, *Stable group theory and approximate subgroups*, preprint available at <http://arxiv.org/abs/0909.2190>.
37. H. Iwaniec, *Unpublished notes on sieve theory*.
38. H. Iwaniec and E. Kowalski, *Analytic number theory*, American Mathematical Society Colloquium Publications, **53**. American Mathematical Society, Providence, RI, 2004. xii+615 pp.
39. A. Kontorovich and H. Oh, *Apollonian circle packings and closed horospheres on hyperbolic 3-manifolds*, preprint available at <http://arxiv.org/abs/0811.2236>.
40. A. Lubotzky, *Cayley graphs: eigenvalues, expanders and random walks*, Surveys in combinatorics, 1995 (Stirling), 155–189, London Math. Soc. Lecture Note Ser., **218**, Cambridge Univ. Press, Cambridge, 1995.
41. A. Lubotzky, R. Phillips and P. Sarnak, *Ramanujan graphs*, Combinatorica **8** (1988), no. 3, 261–277.
42. G. A. Margulis, *Explicit constructions of expanders*, Problemy Peredači Informacii **9** (1973), no. 4, 71–80.
43. ———, *Explicit constructions of graphs without short cycles and low density codes*, Combinatorica **2** (1982), no. 1, 71–78.
44. C. Matthews, L. Vaserstein and B. Weisfeiler, *Congruence properties of Zariski-dense subgroups, I*, Proc. London Math. Soc. **45** (1984), 514–532.
45. M. V. Nori, *On subgroups of  $GL_n(\mathbb{F}_p)$* , Invent. Math. **88** (1987), no. 2, 257–275.
46. A. A. Razborov, *A product theorem in free groups*, preprint available on the author's webpage.
47. I. Z. Ruzsa, *On the cardinality of  $A + A$  and  $A - A$* , Combinatorics (Proc. Fifth Hungarian Colloq., Keszthely, 1976), Vol. II, pp. 933–938, Colloq. Math. Soc. János Bolyai, 18, North-Holland, Amsterdam-New York, 1978.

48. ———, *Generalized arithmetical progressions and sumsets*, Acta Math. Hungar. **65** (1994), no. 4, 379–388.
49. ———, *An analogue of Freiman's theorem in groups*, Astérisque **258** (1999), 323–326.
50. T. Sanders, *From polynomial growth to metric balls in polynomial groups*, preprint.
51. P. Sarnak, *What is... an expander?* Notices Amer. Math. Soc. **51** (2004), no. 7, 762–763.
52. ———, *Equidistribution and primes*, available on the author's website.
53. ———, *Letter to J. Lagarias*, available on the author's website.
54. P. Sarnak and X. X. Xue, *Bounds for multiplicities of automorphic representations*, Duke Math. J. **64** (1991), no. 1, 207–227.
55. Y. Shalom and T. Tao, *A finitary version of Gromov's polynomial growth theorem*, preprint available at <http://arxiv.org/abs/0910.4148>.
56. J. Solymosi, *On sum-sets and product-sets of complex numbers*, J. Théor. Nombres Bordeaux **17** (2005), no. 3, 921–924.
57. ———, *An upper bound on the multiplicative energy*, preprint available at <http://arxiv.org/abs/0806.1040>.
58. M. Suzuki, *Group theory I*, Springer-Verlag, New York, 1982.
59. T. C. Tao, *Product set estimates for non-commutative groups*, Combinatorica **28** (2008), no. 5, 547–594.
60. ———, *The sum-product phenomenon in arbitrary rings*, to appear in Contrib. Discrete Math.
61. ———, *Freiman's theorem for solvable groups*, preprint available at <http://arxiv.org/abs/0906.3535>.
62. T. C. Tao and others, discussion at <http://terrytao.wordpress.com/2009/06/21/freimans-theorem-for-solvable-groups/>.
63. T. C. Tao and V. H. Vu, *Additive Combinatorics*, Cambridge Studies in Advanced Mathematics **105**, Cambridge University Press 2006.

CENTRE FOR MATHEMATICAL SCIENCES, WILBERFORCE ROAD, CAMBRIDGE CB3 0WA

*Current address:* Radcliffe Institute for Advanced Study, 8 Garden Street, Cambridge MA 02138

*E-mail address:* [b.j.green@dpms.cam.ac.uk](mailto:b.j.green@dpms.cam.ac.uk)



# MULTIVARIATE STABLE POLYNOMIALS: THEORY AND APPLICATIONS

DAVID G. WAGNER

*In memoriam Julius Borcea.*

ABSTRACT. Univariate polynomials with only real roots – while special – do occur often enough that their properties can lead to interesting conclusions in diverse areas. Due mainly to the recent work of two young mathematicians, Julius Borcea and Petter Brändén, a very successful multivariate generalization of this method has been developed. The first part of this paper surveys some of the main results of this theory of “multivariate stable” polynomials – the most central of these results is the characterization of linear transformations preserving stability of polynomials. The second part presents various applications of this theory in complex analysis, matrix theory, probability and statistical mechanics, and combinatorics.

## 1. INTRODUCTION.

I have been asked by the AMS to survey the recent work of Julius Borcea and Petter Brändén on their multivariate generalization of the theory of univariate polynomials with only real roots, and its applications. It is exciting work – elementary but subtle, and with spectacular consequences. Borcea and Brändén take center stage but there are many other actors, many of whom I am unable to mention in this brief treatment. Notably, Leonid Gurvits provides a transparent proof of a vast generalization of the famous van der Waerden Conjecture.

Space is limited and I have been advised to use “Bourbaki style”, and so this is an account of the essentials of the theory and a few of its applications, with complete proofs as far as possible. Some relatively straightforward arguments have been left as exercises to engage the reader, and some more specialized topics are merely sketched or even omitted. For the full story and the history and context of the subject one must go to the references cited, the references they cite, and so on. The introduction of [4], in particular, gives a good account of the genesis of the theory.

Here is a brief summary of the contents. Section 2 introduces stable polynomials, gives some examples, presents their elementary properties, and develops multivariate generalizations of two classical univariate results: the Hermite-Kekeya-Obreschkoff and Hermite-Biehler Theorems. We also state the Pólya-Schur Theorem characterizing “multiplier sequences”, as this provides an inspiration for much of the multivariate theory. Section 3 restricts attention to multiaffine stable polynomials: we present a characterization of multiaffine real stable polynomials by means

---

2000 *Mathematics Subject Classification.* Primary: 32A60; Secondary: 05A20, 05B35, 15A45, 15A48, 60G55, 60K35.

Research supported by NSERC Discovery Grant OGP0105392.

of parameterized quadratic inequalities, and characterize those linear transformations which take multiaffine stable polynomials to stable polynomials. In Section 4 we use parts of the forgoing for Borcea and Brändén's splendid proof of the Grace-Walsh-Szegő Coincidence Theorem. In Section 5, the Grace-Walsh-Szegő Theorem is used to extend the results of Section 3 from multiaffine to arbitrary stable polynomials. This culminates in an amazing multivariate generalization of the Pólya-Schur Theorem, the proof of which requires the development of a multivariate extension of the Szasz Principle (which is omitted, regretfully, for lack of space). Section 6 presents Borcea and Brändén's resolution of some matrix-theoretic conjectures of Johnson. Section 7 presents the derivation by Borcea, Brändén, and Liggett of negative association inequalities for the symmetric exclusion process, a fundamental model in probability and statistical mechanics. Section 8 presents Gurvits's sweeping generalization of the van der Waerden Conjecture. Finally, Section 9 briefly mentions a few further topics that could not be included fully for lack of space.

I thank Petter Brändén kindly for his helpful comments on preliminary drafts of this paper.

## 2. STABLE POLYNOMIALS.

We use the following shorthand notation for multivariate polynomials. Let  $[m] = \{1, 2, \dots, m\}$ , let  $\mathbf{x} = (x_1, \dots, x_m)$  be a sequence of indeterminates, and let  $\mathbb{C}[\mathbf{x}]$  be the ring of complex polynomials in the indeterminates  $\mathbf{x}$ . For a function  $\alpha : [m] \rightarrow \mathbb{N}$ , let  $\mathbf{x}^\alpha = x_1^{\alpha(1)} \cdots x_m^{\alpha(m)}$  be the corresponding monomial. For  $S \subseteq [m]$  we also let  $\mathbf{x}^S = \prod_{i \in S} x_i$ . Similarly, for  $i \in [m]$  let  $\partial_i = \partial/\partial x_i$ , let  $\boldsymbol{\partial} = (\partial_1, \dots, \partial_m)$ , let  $\boldsymbol{\partial}^\alpha = \partial_1^{\alpha(1)} \cdots \partial_m^{\alpha(m)}$  and let  $\boldsymbol{\partial}^S = \prod_{i \in S} \partial_i$ . The constant functions on  $[m]$  with images 0 or 1 are denoted by  $\mathbf{0}$  and  $\mathbf{1}$ , respectively. The  $\mathbf{x}$  indeterminates are always indexed by  $[m]$ .

Let  $\mathcal{H} = \{z \in \mathbb{C} : \text{Im}(z) > 0\}$  denote the open upper half of the complex plane, and  $\overline{\mathcal{H}}$  the closure of  $\mathcal{H}$  in  $\mathbb{C}$ . A polynomial  $f \in \mathbb{C}[\mathbf{x}]$  is *stable* provided that either  $f \equiv 0$  identically, or whenever  $\mathbf{z} = (z_1, \dots, z_m) \in \mathcal{H}^m$  then  $f(\mathbf{z}) \neq 0$ . We use  $\mathfrak{S}[\mathbf{x}]$  to denote the set of stable polynomials in  $\mathbb{C}[\mathbf{x}]$ , and  $\mathfrak{S}_{\mathbb{R}}[\mathbf{x}] = \mathfrak{S}[\mathbf{x}] \cap \mathbb{R}[\mathbf{x}]$  for the set of *real stable* polynomials in  $\mathbb{R}[\mathbf{x}]$ . (Borcea and Brändén do not consider the zero polynomial to be stable, but I find the above convention more convenient.)

We rely on the following essential fact at several points.

**Hurwitz's Theorem** (Theorem 1.3.8 of [14]). *Let  $\Omega \subseteq \mathbb{C}^m$  be a connected open set, and let  $(f_n : n \in \mathbb{N})$  be a sequence of functions, each analytic and nonvanishing on  $\Omega$ , which converges to a limit  $f$  uniformly on compact subsets of  $\Omega$ . Then  $f$  is either nonvanishing on  $\Omega$  or identically zero.*

Consequently, a polynomial obtained as the limit of a convergent sequence of stable polynomials is itself stable.

### 2.1. Examples.

**Proposition 2.1** (Proposition 2.4 of [1]). *For  $i \in [m]$ , let  $A_i$  be an  $n$ -by- $n$  matrix and let  $x_i$  be an indeterminate, and let  $B$  be an  $n$ -by- $n$  matrix. If  $A_i$  is positive semidefinite for all  $i \in [m]$  and  $B$  is Hermitian then*

$$f(\mathbf{x}) = \det(x_1 A_1 + x_2 A_2 + \cdots + x_m A_m + B)$$

is real stable.

*Proof.* Let  $\bar{f}$  denote the coefficientwise complex conjugate of  $f$ . Since  $\bar{A}_i = A_i^\top$  for all  $i \in [m]$ , and  $\bar{B} = B^\top$ , it follows that  $\bar{f} = f$ , so that  $f \in \mathbb{R}[\mathbf{x}]$ . By Hurwitz's Theorem and a routine perturbation argument, it suffices to prove that  $f$  is stable when each  $A_i$  is positive definite. Consider any  $\mathbf{z} = \mathbf{a} + i\mathbf{b} \in \mathcal{H}^m$ , with  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^m$  and  $b_i > 0$  for all  $i \in [m]$  (abbreviated to  $\mathbf{b} > \mathbf{0}$ ). Now  $Q = \sum_{i=1}^m b_i A_i$  is positive definite, and hence has a positive definite square-root  $Q^{1/2}$ . Also note that  $H = \sum_{i=1}^m a_i A_i + B$  is Hermitian, and that

$$f(\mathbf{z}) = \det(Q) \det(iI + Q^{-1/2} H Q^{-1/2}).$$

Since  $\det(Q) \neq 0$ , if  $f(\mathbf{z}) = 0$  then  $-i$  is an eigenvalue of  $Q^{-1/2} H Q^{-1/2}$ , contradicting the fact that this matrix is Hermitian. Thus,  $f(\mathbf{z}) \neq 0$  for all  $\mathbf{z} \in \mathcal{H}^m$ . That is,  $f$  is stable.  $\square$

**Corollary 2.2.** *Let  $Q$  be an  $n$ -by- $m$  complex matrix, and let  $X = \text{diag}(x_1, \dots, x_m)$  be a diagonal matrix of indeterminates. Then  $f(\mathbf{x}) = \det(QXQ^\dagger)$  is real stable.*

*Proof.* Let  $Q = (q_{ij})$ , and for  $j \in [m]$  let  $A_j$  denote the  $n$ -by- $n$  matrix with  $hi$ -th entry  $q_{hj}\bar{q}_{ij}$ . That is,  $A_j = Q_j Q_j^\dagger$  in which  $Q_j$  denotes the  $j$ -th column of  $Q$ . Since each  $A_j$  is positive semidefinite and  $QXQ^\dagger = x_1 A_1 + \dots + x_m A_m$ , the conclusion follows directly from Proposition 2.1.  $\square$

**2.2. Elementary properties.** The following simple observation often allows multivariate problems to be reduced to univariate ones, as will be seen.

**Lemma 2.3.** *A polynomial  $f \in \mathbb{C}[\mathbf{x}]$  is stable if and only if for all  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^m$  with  $\mathbf{b} > \mathbf{0}$ ,  $f(\mathbf{a} + \mathbf{b}t)$  is stable in  $\mathfrak{S}[t]$ .*

*Proof.* Since  $\mathcal{H}^m = \{\mathbf{a} + \mathbf{b}t : \mathbf{a}, \mathbf{b} \in \mathbb{R}^m, \mathbf{b} > \mathbf{0}, \text{ and } t \in \mathcal{H}\}$ , the result follows.  $\square$

For  $f \in \mathbb{C}[\mathbf{x}]$  and  $i \in [m]$ , let  $\deg_i(f)$  denote the degree of  $x_i$  in  $f$ .

**Lemma 2.4.** *These operations preserve stability of polynomials in  $\mathbb{C}[\mathbf{x}]$ .*

- (a) **Permutation:** for any permutation  $\sigma : [m] \rightarrow [m]$ ,  $f \mapsto f(x_{\sigma(1)}, \dots, x_{\sigma(m)})$ .
- (b) **Scaling:** for  $c \in \mathbb{C}$  and  $\mathbf{a} \in \mathbb{R}^m$  with  $\mathbf{a} > \mathbf{0}$ ,  $f \mapsto cf(a_1 x_1, \dots, a_m x_m)$ .
- (c) **Diagonalization:** for  $\{i, j\} \subseteq [m]$ ,  $f \mapsto f(\mathbf{x})|_{x_i=x_j}$ .
- (d) **Specialization:** for  $a \in \mathcal{H}$ ,  $f \mapsto f(a, x_2, \dots, x_m)$ .
- (e) **Inversion:** if  $\deg_1(f) = d$ ,  $f \mapsto x_1^d f(-x_1^{-1}, x_2, \dots, x_m)$ .
- (f) **Differentiation (or "Contraction"):**  $f \mapsto \partial_1 f(\mathbf{x})$ .

*Proof.* Parts (a,b,c) are clear. Part (d) is also clear in the case that  $\text{Im}(a) > 0$ . For  $a \in \mathbb{R}$  apply part (d) with values in the sequence  $(a + i2^{-n} : n \in \mathbb{N})$ , and then apply Hurwitz's Theorem to the limit as  $n \rightarrow \infty$ . Part (e) follows from the fact that  $\mathcal{H}$  is invariant under the operation  $z \mapsto -z^{-1}$ . For part (f), let  $d = \deg_1(f)$ , and consider the sequence  $f_n = n^{-d} f(nx_1, x_2, \dots, x_m)$  for all  $n \geq 1$ . Each  $f_n$  is stable and the sequence converges to a polynomial, so the limit is stable. Since  $\deg_1(f) = d$ , this limit is not identically zero. This implies that for all  $z_2, \dots, z_m \in \mathcal{H}$ , the polynomial  $g(x) = f(x, z_2, \dots, z_m) \in \mathbb{C}[x]$  has degree  $d$ . Clearly  $g'(x) = \partial_1 f(x, z_2, \dots, z_m)$ . Let  $\xi_1, \dots, \xi_d$  be the roots of  $g(x)$ , so that  $g(x) = c \prod_{h=1}^d (x - \xi_h)$  for some  $c \in \mathbb{C}$ . Since  $f$  is stable,  $\text{Im}(\xi_h) \leq 0$  for all  $h \in [d]$ . Now

$$\frac{g'(x)}{g(x)} = \frac{d}{dx} \log g(x) = \sum_{h=1}^d \frac{1}{x - \xi_h}.$$



If  $\text{Im}(z) > 0$  then  $\text{Im}(1/(z - \xi_h)) < 0$  for all  $h \in [d]$ , so that  $g'(z) \neq 0$ . Thus, if  $\mathbf{z} \in \mathcal{H}^m$  then  $\partial_1 f(\mathbf{z}) \neq 0$ . That is,  $\partial_1 f$  is stable.  $\square$

Of course, by permutation, parts (d,e,f) of Lemma 2.4 apply for any index  $i \in [m]$  as well (not just  $i = 1$ ). Part (f) is essentially the Gauss-Lucas Theorem: the roots of  $g'(x)$  lie in the convex hull of the roots of  $g(x)$ .

**2.3. Univariate stable polynomials.** A nonzero univariate polynomial is real stable if and only if it has only real roots. Let  $f$  and  $g$  be two such polynomials, let  $\xi_1 \leq \xi_2 \leq \dots \leq \xi_k$  be the roots of  $f$ , and let  $\theta_1 \leq \theta_2 \leq \dots \leq \theta_\ell$  be the roots of  $g$ . These roots are *interlaced* if they are ordered so that  $\xi_1 \leq \theta_1 \leq \xi_2 \leq \theta_2 \leq \dots$  or  $\theta_1 \leq \xi_1 \leq \theta_2 \leq \xi_2 \leq \dots$ . For each  $i \in [\ell]$ , let  $\widehat{g}_i = g/(x - \theta_i)$ . If  $\deg f \leq \deg g$  and the roots of  $g$  are simple, then there is a unique  $(a, b_1, \dots, b_\ell) \in \mathbb{R}^{\ell+1}$  such that

$$f = ag + b_1 \widehat{g}_1 + \dots + b_\ell \widehat{g}_\ell.$$

**Exercise 2.5.** Let  $f, g \in \mathfrak{S}_{\mathbb{R}}[x]$  be nonzero and such that  $fg$  has only simple roots, let  $\deg f \leq \deg g$ , and let  $\theta_1 < \dots < \theta_\ell$  be the roots of  $g$ . The following are equivalent:

- (a) The roots of  $f$  and  $g$  are interlaced.
- (b) The sequence  $f(\theta_1), f(\theta_2), \dots, f(\theta_\ell)$  alternates in sign (strictly).
- (c) In  $f = ag + \sum_{i=1}^{\ell} b_i \widehat{g}_i$ , all of  $b_1, \dots, b_\ell$  have the same sign (and are nonzero).

The *Wronskian* of  $f, g \in \mathbb{C}[x]$  is  $W[f, g] = f' \cdot g - f \cdot g'$ . If  $f = ag + \sum_{i=1}^{\ell} b_i \widehat{g}_i$  as in Exercise 2.5 then

$$\frac{W[f, g]}{g^2} = \frac{d}{dx} \left( \frac{f}{g} \right) = \sum_{i=1}^{\ell} \frac{-b_i}{(x - \theta_i)^2}.$$

It follows that if  $f$  and  $g$  are as in Exercise 2.5(a) then  $W[f, g]$  is either positive for all real  $x$ , or negative for all real  $x$ . Since  $W[g, f] = -W[f, g]$  the condition that  $\deg f \leq \deg g$  is immaterial. Any pair  $f, g$  with interlacing roots can be approximated arbitrarily closely by such a pair with all roots of  $fg$  simple. It follows that for any pair  $f, g$  with interlacing roots, the Wronskian  $W[f, g]$  is either nonnegative on all of  $\mathbb{R}$  or nonpositive on all of  $\mathbb{R}$ .

Nonzero univariate polynomials  $f, g \in \mathfrak{S}_{\mathbb{R}}[x]$  are in *proper position*, denoted by  $f \ll g$ , if  $W[f, g] \leq 0$  on all of  $\mathbb{R}$ . For convenience we also let  $0 \ll f$  and  $f \ll 0$  for any  $f \in \mathfrak{S}_{\mathbb{R}}[x]$ ; in particular  $0 \ll 0$ .

**Exercise 2.6.** Let  $f, g \in \mathfrak{S}_{\mathbb{R}}[x]$  be real stable. Then  $f \ll g$  and  $g \ll f$  if and only if  $cf = dg$  for some  $c, d \in \mathbb{R}$  not both zero.

**Hermite-Kakeya-Obreschkoff (HKO) Theorem** (Theorem 6.3.8 of [14]). *Let  $f, g \in \mathbb{R}[x]$ . Then  $af + bg \in \mathfrak{S}_{\mathbb{R}}[x]$  for all  $a, b \in \mathbb{R}$  if and only if  $f, g \in \mathfrak{S}_{\mathbb{R}}[x]$  and either  $f \ll g$  or  $g \ll f$ .*

**Hermite-Biehler (HB) Theorem** (Theorem 6.3.4 of [14]). *Let  $f, g \in \mathbb{R}[x]$ . Then  $g + if \in \mathfrak{S}[x]$  if and only if  $f, g \in \mathfrak{S}_{\mathbb{R}}[x]$  and  $f \ll g$ .*

*Proofs of HKO and HB.* It suffices to prove these when  $fg$  has only simple roots.

For HKO we can assume that  $\deg(f) \leq \deg(g)$ . Exercise 2.5 shows that if the roots of  $f$  and  $g$  are interlaced then for all  $a, b \in \mathbb{R}$ , the roots of  $g$  and  $af + bg$

are interlaced, so that  $af + bg$  is real stable. The converse is trivial if  $cf = dg$  for some  $c, d \in \mathbb{R}$  not both zero, so assume otherwise. From the hypothesis, both  $f$  and  $g$  are real stable. If there are  $z_0, z_1 \in \mathcal{H}$  for which  $\operatorname{Im}(f(z_0)/g(z_0)) < 0$  and  $\operatorname{Im}(f(z_1)/g(z_1)) > 0$ , then for some  $\lambda \in [0, 1]$  the number  $z_\lambda = (1-\lambda)z_0 + \lambda z_1$  is such that  $\operatorname{Im}(f(z_\lambda)/g(z_\lambda)) = 0$ . Thus  $f(z_\lambda) - ag(z_\lambda) = 0$  for some real number  $a \in \mathbb{R}$ . Since  $f - ag$  is stable (by hypothesis) and  $z_\lambda \in \mathcal{H}$ , this implies that  $f - ag \equiv 0$ , a contradiction. Thus  $\operatorname{Im}(f(z)/g(z))$  does not change sign for  $z \in \mathcal{H}$ . This implies Exercise 2.5(c): all the  $b_i$  have the same sign (consider  $f/g$  at the points  $\theta_i + i\epsilon$  for  $\epsilon > 0$  approaching 0). Thus, the roots of  $f$  and  $g$  are interlaced.

For HB, let  $p = g + if$ . Considering  $ip = -f + ig$  if necessary, we can assume that  $\deg f \leq \deg g$ . If  $f \ll g$  then Exercise 2.5(c) implies that  $\operatorname{Im}(f(z)/g(z)) \leq 0$  for all  $z \in \mathcal{H}$ , so that  $g + if$  is stable. For the converse, let  $p(x) = c \prod_{i=1}^d (x - \xi_i)$ , so that  $\operatorname{Im}(\xi_i) \leq 0$  for all  $i \in [d]$ . Now  $|z - \xi_i| \geq |\bar{z} - \xi_i|$  for all  $z \in \mathcal{H}$  and  $i \in [d]$ , so that  $|p(z)| \geq |p(\bar{z})|$  for all  $z \in \mathcal{H}$ . For any  $z \in \mathcal{H}$  with  $f(z) \neq 0$  we have

$$\left| \frac{g(z)}{f(z)} + i \right| \geq \left| \frac{g(\bar{z})}{f(\bar{z})} + i \right| = \left| \frac{g(z)}{f(z)} - i \right|,$$

and it follows that  $\operatorname{Im}(g(z)/f(z)) \geq 0$  for all  $z \in \mathcal{H}$  with  $f(z) \neq 0$ . Since  $fg$  has simple roots it follows that  $g(x) + yf(x)$  is stable in  $\mathfrak{S}[x, y]$ . By contraction and specialization, both  $f$  and  $g$  are real stable. By scaling and specialization,  $af + bg$  is stable for all  $a, b \in \mathbb{R}$ . By HKO, the roots of  $f$  and  $g$  are interlaced. Since  $\operatorname{Im}(f(z)/g(z)) \leq 0$  for all  $z \in \mathcal{H}$ , all the  $b_i$  in Exercise 2.5(c) are positive, so that  $W[f, g]$  is negative on all of  $\mathbb{R}$ : that is  $f \ll g$ .  $\square$

For  $\lambda : \mathbb{N} \rightarrow \mathbb{R}$ , let  $T_\lambda : \mathbb{R}[x] \rightarrow \mathbb{R}[x]$  be the linear transformation defined by  $T_\lambda(x^n) = \lambda(n)x^n$  and linear extension. A *multiplier sequence (of the first kind)* is such a  $\lambda$  for which  $T_\lambda(f)$  is real stable whenever  $f$  is real stable. Pólya and Schur characterized multiplier sequences as follows.

**Pólya-Schur Theorem** (Theorem 1.7 of [4]). *Let  $\lambda : \mathbb{N} \rightarrow \mathbb{R}$ . The following are equivalent:*

- (a)  $\lambda$  is a multiplier sequence.
- (b)  $F_\lambda(x) = \sum_{n=0}^{\infty} \lambda(n)x^n/n!$  is an entire function which is the limit, uniformly on compact sets, of real stable polynomials with all roots of the same sign.
- (c) Either  $F_\lambda(x)$  or  $F_\lambda(-x)$  has the form

$$Cx^n e^{ax} \prod_{j=1}^{\infty} (1 + \alpha_j x),$$

in which  $C \in \mathbb{R}$ ,  $n \in \mathbb{N}$ ,  $a \geq 0$ , all  $\alpha_j \geq 0$ , and  $\sum_{j=1}^{\infty} \alpha_j$  is finite.

- (d) For all  $n \in \mathbb{N}$ , the polynomial  $T_\lambda((1+x)^n)$  is real stable with all roots of the same sign.

One of the main results of Borcea and Brändén's theory is a great generalization of the Pólya-Schur Theorem – a characterization of all stability preservers: linear transformations  $T : \mathbb{C}[\mathbf{x}] \rightarrow \mathbb{C}[\mathbf{x}]$  such that  $T(f)$  is stable whenever  $f$  is stable. (Also the analogous characterization of real stability preservers.) This is discussed in some detail in Section 5.3.

**2.4. Multivariate analogues of the HKO and HB Theorems.** By analogy with the univariate HB Theorem, polynomials  $f, g \in \mathbb{R}[\mathbf{x}]$  are said to be in *proper position*, denoted by  $f \ll g$ , when  $g + if \in \mathfrak{S}[\mathbf{x}]$ . (As will be seen, this implies that  $f, g \in \mathfrak{S}_{\mathbb{R}}[\mathbf{x}]$ .) Thus, the multivariate analogue of Hermite-Biehler is a definition, not a theorem.

**Proposition 2.7** (Lemma 1.8 and Remark 1.3 of [5]). *Let  $f, g \in \mathbb{C}[\mathbf{x}]$ .*

(a) *If  $f, g \in \mathbb{R}[\mathbf{x}]$  then  $f \ll g$  if and only if  $g + yf \in \mathfrak{S}_{\mathbb{R}}[\mathbf{x}, y]$ .*

(b) *If  $0 \neq f \in \mathfrak{S}[\mathbf{x}]$  then  $g + yf \in \mathfrak{S}[\mathbf{x}, y]$  if and only if for all  $\mathbf{z} \in \mathcal{H}^m$ ,*

$$\operatorname{Im} \left( \frac{g(\mathbf{z})}{f(\mathbf{z})} \right) \geq 0.$$

*Proof.* If  $g + yf \in \mathfrak{S}_{\mathbb{R}}[\mathbf{x}, y]$  then  $g + if \in \mathfrak{S}[\mathbf{x}]$ , by specialization. Conversely, assume that  $h = g + if \in \mathfrak{S}[\mathbf{x}]$  with  $f, g \in \mathbb{R}[\mathbf{x}]$ , and let  $z = a + ib$  with  $a, b \in \mathbb{R}$  and  $b > 0$ . By Lemma 2.3, for all  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^m$  with  $\mathbf{b} > \mathbf{0}$  we have  $h(\mathbf{a} + \mathbf{b}t) \in \mathfrak{S}[t]$ . By HB,  $\widehat{f}(t) = f(\mathbf{a} + \mathbf{b}t)$  and  $\widehat{g}(t) = g(\mathbf{a} + \mathbf{b}t)$  are such that  $\widehat{f} \ll \widehat{g}$ . By HKO,  $c\widehat{f} + d\widehat{g} \in \mathfrak{S}_{\mathbb{R}}[t]$  for all  $c, d \in \mathbb{R}$ . By HKO again, the roots of  $b\widehat{f}$  and of  $\widehat{g} + a\widehat{f}$  are interlaced. Since  $W[b\widehat{f}, \widehat{g} + a\widehat{f}] = bW[\widehat{f}, \widehat{g}] \leq 0$  on  $\mathbb{R}$ , it follows that  $b\widehat{f} \ll \widehat{g} + a\widehat{f}$ . Finally, by HB again,  $\widehat{g} + (a + ib)\widehat{f} \in \mathfrak{S}[t]$ . Since this holds for all  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^m$  with  $\mathbf{b} > \mathbf{0}$ , Lemma 2.3 implies that  $g + (a + ib)f \in \mathfrak{S}[\mathbf{x}]$ . Since this holds for all  $a, b \in \mathbb{R}$  with  $b > 0$ ,  $g + yf \in \mathfrak{S}[\mathbf{x}, y]$ . This proves part (a).

For part (b), first let  $g + yf$  be stable. By specialization,  $g$  is also stable. If  $g \equiv 0$  then there is nothing to prove. Otherwise, consider any  $\mathbf{z} \in \mathcal{H}^m$ , so that  $f(\mathbf{z}) \neq 0$  and  $g(\mathbf{z}) \neq 0$ . There is a unique solution  $z \in \mathbb{C}$  to  $g(\mathbf{z}) + zf(\mathbf{z}) = 0$ , and since  $g + yf$  is stable,  $\operatorname{Im}(z) \leq 0$ . Hence,  $\operatorname{Im}(g(\mathbf{z})/f(\mathbf{z})) = \operatorname{Im}(-z) \geq 0$ . This argument can be reversed to prove the converse implication.  $\square$

**Exercise 2.8** (Corollary 2.4 of [4]).  $\mathfrak{S}[\mathbf{x}] = \{g + if : f, g \in \mathfrak{S}_{\mathbb{R}}[\mathbf{x}] \text{ and } f \ll g\}$ .

Here is the multivariate HKO Theorem of Borcea and Brändén.

**Theorem 2.9** (Theorem 1.6 of [4]). *Let  $f, g \in \mathbb{R}[\mathbf{x}]$ . Then  $af + bg \in \mathfrak{S}_{\mathbb{R}}[\mathbf{x}]$  for all  $a, b \in \mathbb{R}$  if and only if  $f, g \in \mathfrak{S}_{\mathbb{R}}[\mathbf{x}]$  and either  $f \ll g$  or  $g \ll f$ .*

*Proof.* First assume that  $f \ll g$ , and let  $a, b \in \mathbb{R}$  with  $b > 0$ . By Proposition 2.7(a),  $g + yf \in \mathfrak{S}_{\mathbb{R}}[\mathbf{x}, y]$ . By scaling and specialization,  $bg + (a + i)f \in \mathfrak{S}[\mathbf{x}]$ . By Proposition 2.7(a) again,  $f \ll (af + bg)$ . Thus  $af + bg \in \mathfrak{S}_{\mathbb{R}}[\mathbf{x}]$  for all  $a, b \in \mathbb{R}$ . The case that  $g \ll f$  is similar.

Conversely, assume that  $af + bg \in \mathfrak{S}_{\mathbb{R}}[\mathbf{x}]$  for all  $a, b \in \mathbb{R}$ . Let  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^m$  with  $\mathbf{b} > \mathbf{0}$ , and let  $\widehat{f}(t) = f(\mathbf{a} + \mathbf{b}t)$  and  $\widehat{g}(t) = g(\mathbf{a} + \mathbf{b}t)$ . By Lemma 2.3,  $a\widehat{f} + b\widehat{g} \in \mathfrak{S}_{\mathbb{R}}[t]$  for all  $a, b \in \mathbb{R}$ . By HKO, for each  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^m$  with  $\mathbf{b} > \mathbf{0}$ , either  $\widehat{f} \ll \widehat{g}$  or  $\widehat{g} \ll \widehat{f}$ .

If  $\widehat{f} \ll \widehat{g}$  for all  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^m$  with  $\mathbf{b} > \mathbf{0}$ , then by HB,  $\widehat{g} + i\widehat{f} \in \mathfrak{S}[t]$  for all  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^m$  with  $\mathbf{b} > \mathbf{0}$ . Thus  $g + if \in \mathfrak{S}[\mathbf{x}]$  by Lemma 2.3, which is to say that  $f \ll g$  (by definition). Similarly, if  $\widehat{g} \ll \widehat{f}$  for all  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^m$  with  $\mathbf{b} > \mathbf{0}$  then  $g \ll f$ .

It remains to consider the case that  $f(\mathbf{a}_0 + \mathbf{b}_0t) \ll g(\mathbf{a}_0 + \mathbf{b}_0t)$  for some  $\mathbf{a}_0, \mathbf{b}_0 \in \mathbb{R}^m$  with  $\mathbf{b}_0 > \mathbf{0}$ , and  $g(\mathbf{a}_1 + \mathbf{b}_1t) \ll f(\mathbf{a}_1 + \mathbf{b}_1t)$  for another  $\mathbf{a}_1, \mathbf{b}_1 \in \mathbb{R}^m$  with  $\mathbf{b}_1 > \mathbf{0}$ . For  $0 \leq \lambda \leq 1$ , let  $\mathbf{a}_\lambda = (1 - \lambda)\mathbf{a}_0 + \lambda\mathbf{a}_1$  and  $\mathbf{b}_\lambda = (1 - \lambda)\mathbf{b}_0 + \lambda\mathbf{b}_1$ . Since roots of polynomials move continuously as the coefficients are varied continuously, there is a value  $0 \leq \lambda \leq 1$  for which both  $f(\mathbf{a}_\lambda + \mathbf{b}_\lambda t) \ll g(\mathbf{a}_\lambda + \mathbf{b}_\lambda t)$  and  $g(\mathbf{a}_\lambda + \mathbf{b}_\lambda t) \ll f(\mathbf{a}_\lambda + \mathbf{b}_\lambda t)$ . From Exercise 2.6, it follows that  $cf(\mathbf{a}_\lambda + \mathbf{b}_\lambda t) = dg(\mathbf{a}_\lambda + \mathbf{b}_\lambda t)$

for some  $c, d \in \mathbb{R}$  not both zero. Now  $h = cf - dg \in \mathfrak{S}[\mathbf{x}]$  by hypothesis, and since  $h(\mathbf{a}_\lambda + \mathbf{b}_\lambda t) \equiv 0$  identically, it follows that  $h(\mathbf{a}_\lambda + \mathbf{b}_\lambda) = 0$ . Since  $\mathbf{b}_\lambda > \mathbf{0}$  and  $h$  is stable, this implies that  $h \equiv 0$ , so that  $cf = dg$  in  $\mathfrak{S}[\mathbf{x}]$ . In this case, both  $f \ll g$  and  $g \ll f$  hold.  $\square$

For  $f, g \in \mathbb{C}[\mathbf{x}]$  and  $i \in [m]$ , let  $W_i[f, g] = \partial_i f \cdot g - f \cdot \partial_i g$  be the  $i$ -th Wronskian of the pair  $(f, g)$ .

**Corollary 2.10** (Theorem 1.9 of [5]). *Let  $f, g \in \mathbb{R}[\mathbf{x}]$ . The following are equivalent:*

- (a)  $g + if$  is stable in  $\mathfrak{S}[\mathbf{x}]$ , that is  $f \ll g$ ;
- (b)  $g + yf$  is real stable in  $\mathfrak{S}_{\mathbb{R}}[\mathbf{x}, y]$ ;
- (c)  $af + bg \in \mathfrak{S}_{\mathbb{R}}[\mathbf{x}]$  for all  $a, b \in \mathbb{R}$ , and  $W_i[f, g](\mathbf{a}) \leq 0$  for all  $i \in [m]$  and  $\mathbf{a} \in \mathbb{R}^m$ .

*Proof.* Proposition 2.7(a) shows that (a) and (b) are equivalent.

If (a) holds then Theorem 2.9 implies that  $af + bg \in \mathfrak{S}_{\mathbb{R}}[\mathbf{x}]$  for all  $a, b \in \mathbb{R}$ . To prove the rest of (c), let  $i \in [m]$  and  $\mathbf{a} \in \mathbb{R}^m$ , and let  $\delta_i \in \mathbb{R}^m$  be the unit vector with a one in the  $i$ -th position. Since  $f \ll g$ , for any  $\mathbf{b} \in \mathbb{R}^m$  with  $\mathbf{b} > \mathbf{0}$  we have  $f(\mathbf{a} + (\mathbf{b} + \delta_i)t) \ll g(\mathbf{a} + (\mathbf{b} + \delta_i)t)$  in  $\mathfrak{S}_{\mathbb{R}}[t]$ , from Proposition 2.7(a) and Lemma 2.3. By the Wronskian condition for univariate polynomials in proper position,

$$W[f(\mathbf{a} + (\mathbf{b} + \delta_i)t), g(\mathbf{a} + (\mathbf{b} + \delta_i)t)] \leq 0$$

for all  $t \in \mathbb{R}$ . Taking the limit as  $\mathbf{b} \rightarrow \mathbf{0}$  and evaluating at  $t = 0$  yields

$$W_i[f, g](\mathbf{a}) = W[f(\mathbf{a} + \delta_i t), g(\mathbf{a} + \delta_i t)]|_{t=0} \leq 0,$$

by continuity. Thus (a) implies (c).

To prove that (c) implies (b), let  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^m$  with  $\mathbf{b} = (b_1, \dots, b_m) > \mathbf{0}$ , and let  $a, b \in \mathbb{R}$  with  $b > 0$ . By Lemma 2.3, to show that  $g + yf \in \mathfrak{S}_{\mathbb{R}}[\mathbf{x}, y]$  it suffices to show that  $g(\mathbf{a} + \mathbf{b}t) + (a + ib)f(\mathbf{a} + \mathbf{b}t) \in \mathfrak{S}[t]$ . From (c) it follows that  $p = g + af$  and  $q = bf$  are such that  $\alpha p + \beta q \in \mathfrak{S}_{\mathbb{R}}[\mathbf{x}]$  for all  $\alpha, \beta \in \mathbb{R}$ . By Theorem 2.9, either  $p \ll q$  or  $q \ll p$ . Now

$$\begin{aligned} W[q(\mathbf{a} + \mathbf{b}t), p(\mathbf{a} + \mathbf{b}t)] &= bW[f(\mathbf{a} + \mathbf{b}t), g(\mathbf{a} + \mathbf{b}t)] \\ &= b \sum_{i=1}^m b_i W_i[f, g](\mathbf{a} + \mathbf{b}t) \leq 0, \end{aligned}$$

by the Wronskian condition in part (c). Thus  $q(\mathbf{a} + \mathbf{b}t) \ll p(\mathbf{a} + \mathbf{b}t)$ , so that  $p(\mathbf{a} + \mathbf{b}t) + iq(\mathbf{a} + \mathbf{b}t) \in \mathfrak{S}[t]$ . Since  $p + iq = g + (a + ib)f$ , this shows that (c) implies (b).  $\square$

**Exercise 2.11** (Corollary 1.10 of [5]). Let  $f, g \in \mathfrak{S}_{\mathbb{R}}[\mathbf{x}]$  be real stable. Then  $f \ll g$  and  $g \ll f$  if and only if  $cf = dg$  for some  $c, d \in \mathbb{R}$  not both zero.

**Proposition 2.12** (Lemma 3.2 of [5]). *Let  $V$  be a  $\mathbb{K}$ -vector subspace of  $\mathbb{K}[\mathbf{x}]$ , with either  $\mathbb{K} = \mathbb{R}$  or  $\mathbb{K} = \mathbb{C}$ .*

- (a) *If  $\mathbb{K} = \mathbb{R}$  and  $V \subseteq \mathfrak{S}_{\mathbb{R}}[\mathbf{x}]$  then  $\dim_{\mathbb{R}} V \leq 2$ .*
- (b) *If  $\mathbb{K} = \mathbb{C}$  and  $V \subseteq \mathfrak{S}[\mathbf{x}]$  then  $\dim_{\mathbb{C}} V \leq 1$ .*

*Proof.* For part (a), suppose to the contrary that  $f, g, h \in V$  are linearly independent over  $\mathbb{R}$  (and hence not identically zero). By Theorem 2.9, either  $f \ll g$  or  $g \ll f$ , and similarly for the other pairs  $\{f, h\}$  and  $\{g, h\}$ . Renaming these polynomials as necessary, we may assume that  $f \ll h$  and  $h \ll g$ . Now, for all  $\lambda \in [0, 1]$  let  $p_\lambda = (1 - \lambda)f + \lambda g$ , and note that each  $p_\lambda \not\equiv 0$ . By Theorem 2.9, for each  $\lambda \in [0, 1]$  either  $h \ll p_\lambda$  or  $p_\lambda \ll h$ . Since  $p_0 = f \ll h$  and  $h \ll g = p_1$ , by continuity of

the roots of  $\{p_\lambda : \lambda \in [0, 1]\}$  there is a  $\lambda \in [0, 1]$  such that  $h \ll p_\lambda$  and  $p_\lambda \ll h$ . But then, by Exercise 2.11, either  $\{f, g\}$  is linearly dependent or  $h$  is in the span of  $\{f, g\}$ , contradicting the supposition.

For part (b), let  $\text{Re}(V) = \{\text{Re}(h) : h \in V\}$ . Then  $\text{Re}(V)$  is a real subspace of  $\mathfrak{S}_{\mathbb{R}}[\mathbf{x}]$ , so that  $\dim_{\mathbb{R}} \text{Re}(V) \leq 2$  by part (a). If  $\dim_{\mathbb{R}} \text{Re}(V) \leq 1$  then  $\dim_{\mathbb{C}} V \leq 1$ . In the remaining case let  $\{p, q\}$  be a basis of  $\text{Re}(V)$  with  $f = p + iq \in V$ . By Corollary 2.10,  $W_i[q, p](\mathbf{a}) \leq 0$  for all  $i \in [m]$  and  $\mathbf{a} \in \mathbb{R}^m$ . Since  $p$  and  $q$  are not linearly dependent, there is an index  $k \in [m]$  such that  $W_k[q, p] \neq 0$ .

Consider any  $g \in V$ . There are reals  $a, b, c, d \in \mathbb{R}$  such that

$$g = (ap + bq) + i(cp + dq).$$

Since  $g$  is stable,  $W_k[cp + dq, ap + bq](\mathbf{a}) = (ad - bc)W_k[q, p](\mathbf{a}) \leq 0$  for all  $\mathbf{a} \in \mathbb{R}^m$ . Since  $W_k[q, p] \neq 0$ , it follows that  $ad - bc \geq 0$ . Now, for any  $v, w \in \mathbb{R}$ ,  $g + (v + iw)f$  is in  $V$ . Since

$$g + (v + iw)f = (a + v)p + (b - w)q + i((c + w)p + (d + v)q),$$

this argument shows that  $H = (a + v)(d + v) - (b - w)(c + w) \geq 0$  for all  $v, w \in \mathbb{R}$ . But

$$4H = (2v + a + d)^2 + (2w + c - b)^2 - (a - d)^2 - (b + c)^2,$$

so that  $H \geq 0$  for all  $v, w \in \mathbb{R}$  if and only if  $a = d$  and  $b = -c$ . This implies that  $g = (a + ic)f$ , so that  $\dim_{\mathbb{C}} V = 1$ .  $\square$

### 3. MULTIAFFINE STABLE POLYNOMIALS.

A polynomial  $f$  is *multiaffine* if each indeterminate occurs at most to the first power in  $f$ . For a set  $\mathcal{S}$  of polynomials, let  $\mathcal{S}^{\text{MA}}$  denote the set of multiaffine polynomials in  $\mathcal{S}$ . For multiaffine  $f \in \mathbb{C}[\mathbf{x}]^{\text{MA}}$  and  $i \in [m]$  we use the ‘‘ultra-shorthand’’ notation  $f = f^i + x_i f_i$  in which  $f^i = f|_{x_i=0}$  and  $f_i = \partial_i f$ . This notation is extended to multiple distinct indices in the obvious way – in particular,

$$f = f^{ij} + x_i f_i^j + x_j f_j^i + x_i x_j f_{ij}.$$

**3.1. A criterion for real stability.** For  $f \in \mathbb{C}[\mathbf{x}]$  and  $\{i, j\} \subseteq [m]$ , let

$$\Delta_{ij} f = \partial_i f \cdot \partial_j f - f \cdot \partial_i \partial_j f.$$

Notice that for  $f \in \mathbb{C}[\mathbf{x}]^{\text{MA}}$ ,

$$\Delta_{ij} f = f_i^j f_j - f^j f_{ij} = W_i[f^j, f_j] = -W_i[f_j, f^j],$$

and

$$\Delta_{ij} f = f_i^j f_j^i - f^{ij} f_{ij}.$$

**Theorem 3.1** (Theorem 5.6 of [8] and Theorem 3 of [16]). *Let  $f \in \mathbb{R}[\mathbf{x}]^{\text{MA}}$  be multiaffine. The following are equivalent:*

- (a)  $f$  is real stable.
- (b) For all  $\{i, j\} \subseteq [m]$  and all  $\mathbf{a} \in \mathbb{R}^m$ ,  $\Delta_{ij} f(\mathbf{a}) \geq 0$ .
- (c) Either  $m = 1$ , or there exists  $\{i, j\} \subseteq [m]$  such that  $f_i, f^i, f_j$  and  $f^j$  are real stable, and  $\Delta_{ij} f(\mathbf{a}) \geq 0$  for all  $\mathbf{a} \in \mathbb{R}^m$ .

*Proof.* To see that (a) implies (b), fix  $\{i, j\} \subseteq [m]$ . Proposition 2.7(a) shows that  $f_j \ll f^j$ , and from the calculation above and Corollary 2.10, it follows that  $\Delta_{ij}f(\mathbf{a}) = -W_i[f_j, f^j](\mathbf{a}) \geq 0$  for all  $\mathbf{a} \in \mathbb{R}^m$ .

We show that (b) implies (a) by induction on  $m$ , the base case  $m = 1$  being trivial. For the induction step let  $f$  be as in part (b), let  $a \in \mathbb{R}$ , and let  $g = f|_{x_m=a}$ . For all  $\{i, j\} \subseteq [m-1]$  and  $\mathbf{a} \in \mathbb{R}^{m-1}$ ,  $\Delta_{ij}g(\mathbf{a}) = \Delta_{ij}f(\mathbf{a}, a) \geq 0$ . By induction,  $g = f^m + af_m$  is real stable for all  $a \in \mathbb{R}$ ; it follows that  $af_m + bf^m \in \mathfrak{S}_{\mathbb{R}}[x_1, \dots, x_{m-1}]$  for all  $a, b \in \mathbb{R}$ . Furthermore, for all  $j \in [m-1]$  and  $\mathbf{a} \in \mathbb{R}^{m-1}$ ,  $W_j[f_m, f^m](\mathbf{a}) = -\Delta_{jm}f(\mathbf{a}, 1) \leq 0$ . This verifies condition (c) of Corollary 2.10 for the pair  $(f_m, f^m)$ , and it follows that  $f = f^m + x_m f_m \in \mathfrak{S}_{\mathbb{R}}[\mathbf{x}]$ , completing the induction.

It is clear that (a) and (b) imply (c) – we show that (c) implies (b) below. This is clear if  $m \leq 2$ , so assume that  $m \geq 3$ . To begin with, let  $\{h, i, j\} \subseteq [m]$  be three distinct indices, and consider  $\Delta_{ij}f$  as a polynomial in  $x_h$ . That is,  $\Delta_{ij}f = A_{hij}x_h^2 + B_{hij}x_h + C_{hij}$  in which

$$\begin{aligned} A_{hij} &= f_{hi}^j f_{hj}^i - f_h^{ij} f_{hij} = \Delta_{ij}f_h, \\ B_{hij} &= f_{hi}^j f_j^{hi} - f_h^{ij} f_{ij}^h + f_i^{hj} f_{hj}^i - f^{hij} f_{hij}, \text{ and} \\ C_{hij} &= f_i^{hj} f_j^{hi} - f^{hij} f_{ij}^h = \Delta_{ij}f^h. \end{aligned}$$

If  $\Delta_{ij}f(\mathbf{a}) \geq 0$  for all  $\mathbf{a} \in \mathbb{R}^m$  then this quadratic polynomial in  $x_h$ :

$$\Delta_{ij}f(a_1, \dots, a_{h-1}, x_h, a_{h+1}, \dots, a_m)$$

has a nonpositive discriminant for all  $\mathbf{a} \in \mathbb{R}^m$ . That is,  $D_{hij} = B_{hij}^2 - 4A_{hij}C_{hij}$  is such that  $D_{hij}(\mathbf{a}) \leq 0$  for all  $\mathbf{a} \in \mathbb{R}^m$ .

It is a surprising fact that as a polynomial in  $\{x_k : k \in [m] \setminus \{h, i, j\}\}$ ,  $D_{hij}$  is invariant under all six permutations of its indices, as is seen by direct calculation:

$$\begin{aligned} D_{hij} &= (f_h^{ij} f_{ij}^h)^2 + (f_i^{hj} f_{hj}^i)^2 + (f_j^{hi} f_{hi}^j)^2 + (f_{hij} f^{hij})^2 \\ &\quad - 2(f_h^{ij} f_{ij}^h f_i^{hj} f_{hj}^i + f_i^{hj} f_{hj}^i f_j^{hi} f_{hi}^j + f_j^{hi} f_{hi}^j f_h^{ij} f_{ij}^h) \\ &\quad - 2(f_h^{ij} f_{ij}^h + f_i^{hj} f_{hj}^i + f_j^{hi} f_{hi}^j) f^{hij} f_{hij} \\ &\quad + 4f_{ij}^h f_{hj}^i f_{hi}^j f^{hij} + 4f_h^{ij} f_i^{hj} f_j^{hi} f_{hij}. \end{aligned}$$

Now for the proof that (c) implies (b) when  $m \geq 3$ . Consider any  $h \in [m] \setminus \{i, j\}$ . Then

$$\Delta_{hi}f = A_{jhi}x_j^2 + B_{jhi}x_j + C_{jhi}$$

has discriminant  $D_{jhi} = D_{hij}$ . Since  $f_j$  and  $f^j$  are real stable, we have  $A_{jhi}(\mathbf{a}) = \Delta_{hi}f_j(\mathbf{a}) \geq 0$  and  $C_{jhi}(\mathbf{a}) = \Delta_{hi}f^j(\mathbf{a}) \geq 0$  for all  $\mathbf{a} \in \mathbb{R}^m$ . Since  $\Delta_{ij}f(\mathbf{a}) \geq 0$  for all  $\mathbf{a} \in \mathbb{R}^m$  it follows that  $D_{jhi}(\mathbf{a}) = D_{hij}(\mathbf{a}) \leq 0$  for all  $\mathbf{a} \in \mathbb{R}^m$ . It follows that  $\Delta_{hi}f(\mathbf{a}) \geq 0$  for all  $\mathbf{a} \in \mathbb{R}^m$ . (Note that if  $B^2 - 4AC \leq 0$  and either  $A = 0$  or  $C = 0$ , then  $B = 0$ .) A similar argument using the fact that  $f_i$  and  $f^i$  are real stable shows that  $\Delta_{hj}f(\mathbf{a}) \geq 0$  for all  $\mathbf{a} \in \mathbb{R}^m$ .

It remains to show that  $\Delta_{hk}f(\mathbf{a}) \geq 0$  for all  $\mathbf{a} \in \mathbb{R}^m$  when  $\{h, k\}$  is disjoint from  $\{i, j\}$ . We have seen that  $\Delta_{hi}f(\mathbf{a}) \geq 0$  for all  $\mathbf{a} \in \mathbb{R}^m$ , and we know that both  $f_i$  and  $f^i$  are real stable. The argument above applies once more:  $\Delta_{hi}f(\mathbf{a}) \geq 0$  for all  $\mathbf{a} \in \mathbb{R}^m$ , so that  $D_{ihk}(\mathbf{a}) = D_{khi}(\mathbf{a}) \leq 0$  for all  $\mathbf{a} \in \mathbb{R}^m$ , and then since  $A_{ihk}(\mathbf{a}) \geq 0$  and  $C_{ihk}(\mathbf{a}) \geq 0$  for all  $\mathbf{a} \in \mathbb{R}^m$  it follows that  $\Delta_{hk}f(\mathbf{a}) \geq 0$  for all  $\mathbf{a} \in \mathbb{R}^m$ . Thus (c) implies (b).  $\square$

### 3.2. Linear transformations preserving stability – multiaffine case.

**Lemma 3.2** (Lieb-Sokal Lemma, Lemma 2.1 of [5]). *Let  $g(\mathbf{x}) + yf(\mathbf{x}) \in \mathfrak{S}[\mathbf{x}, y]$  be stable and such that  $\deg_z(f) \leq 1$ . Then  $g - \partial_1 f \in \mathfrak{S}[\mathbf{x}]$  is stable.*

*Proof.* Since  $g$  is stable (by specialization to  $y = 0$ ), there is nothing to prove if  $\partial_1 f \equiv 0$  identically, so assume otherwise (and hence that  $f \not\equiv 0$ ). By permutation we can assume that  $i = 1$ . Since  $f$  is stable and  $z_1, z \in \mathcal{H}$  imply that  $z_1 - z^{-1} \in \mathcal{H}$ , it follows that

$$yf(x_1 - y^{-1}, x_2, \dots, x_m) = -\partial_1 f(\mathbf{x}) + yf(\mathbf{x})$$

is stable. Proposition 2.7(b) implies that for all  $\mathbf{z} \in \mathcal{H}^m$ ,

$$\operatorname{Im} \left( \frac{g(\mathbf{z}) - \partial_1 f(\mathbf{z})}{f(\mathbf{z})} \right) = \operatorname{Im} \left( \frac{g(\mathbf{z})}{f(\mathbf{z})} \right) + \operatorname{Im} \left( \frac{-\partial_1 f(\mathbf{z})}{f(\mathbf{z})} \right) \geq 0.$$

Thus, by Proposition 2.7(b) again,  $g - \partial_1 f + yf$  is stable. Specializing to  $y = 0$  shows that  $g - \partial_1 f$  is stable.  $\square$

**Exercise 3.3** (Lemma 3.1 of [5]). Let  $f \in \mathbb{C}[\mathbf{x}]^{\text{MA}}$  and  $\mathbf{w} \in \mathcal{H}^m$ . Then for all  $\epsilon > 0$  sufficiently small,  $(\mathbf{x} + \mathbf{w})^{[m]} + \epsilon f(\mathbf{x})$  is stable. (Here  $(\mathbf{x} + \mathbf{w})^{[m]} = \prod_{i=1}^m (x_i + w_i)$ .)

For a linear transformation  $T : \mathbb{C}[\mathbf{x}]^{\text{MA}} \rightarrow \mathbb{C}[\mathbf{x}]$  of multiaffine polynomials, define the *algebraic symbol* of  $T$  to be the polynomial

$$T((\mathbf{x} + \mathbf{y})^{[m]}) = T \left( \prod_{i=1}^m (x_i + y_i) \right) = \sum_{S \subseteq [m]} T(\mathbf{x}^S) \mathbf{y}^{[m] \setminus S}$$

in  $\mathbb{C}[x_1, \dots, x_m, y_1, \dots, y_m] = \mathbb{C}[\mathbf{x}, \mathbf{y}]$ .

**Theorem 3.4** (Theorem 1.1 of [5]). *Let  $T : \mathbb{C}[\mathbf{x}]^{\text{MA}} \rightarrow \mathbb{C}[\mathbf{x}]$  be a linear transformation. Then  $T$  maps  $\mathfrak{S}[\mathbf{x}]^{\text{MA}}$  into  $\mathfrak{S}[\mathbf{x}]$  if and only if either*

- (a)  $T(f) = \eta(f) \cdot p$  for some linear functional  $\eta : \mathbb{C}[\mathbf{x}]^{\text{MA}} \rightarrow \mathbb{C}$  and  $p \in \mathfrak{S}[\mathbf{x}]$ , or
- (b) the polynomial  $T((\mathbf{x} + \mathbf{y})^{[m]})$  is stable in  $\mathfrak{S}[\mathbf{x}, \mathbf{y}]$ .

*Proof.* First, assume (b) that  $T((\mathbf{x} + \mathbf{y})^{[m]}) \in \mathfrak{S}[\mathbf{x}, \mathbf{y}]$  is stable. By inversion, it follows that  $\mathbf{y}^{[m]} T((\mathbf{x} - \mathbf{y}^{-1})^{[m]})$  is also stable. Thus, if  $f \in \mathfrak{S}[w_1, \dots, w_m]$  is stable then

$$\mathbf{y}^{[m]} T((\mathbf{x} - \mathbf{y}^{-1})^{[m]}) f(\mathbf{w}) = \sum_{S \subseteq [m]} T(\mathbf{x}^S) (-\mathbf{y})^S f(\mathbf{w})$$

is stable. If  $f$  is also multiaffine then repeated application of the Lieb-Sokal Lemma 3.2 (replacing  $y_i$  by  $-\partial/\partial w_i$  for  $i \in [m]$ ) shows that

$$\sum_{S \subseteq [m]} T(\mathbf{x}^S) \frac{\partial^S}{\partial \mathbf{w}^S} f(\mathbf{w})$$

is stable. Finally, specializing to  $\mathbf{w} = \mathbf{0}$  shows that  $T(f(\mathbf{x}))$  is stable. Thus, the linear transformation  $T$  maps  $\mathfrak{S}[\mathbf{x}]^{\text{MA}}$  into  $\mathfrak{S}[\mathbf{x}]$ . This is clearly also the case if (a) holds.

Conversely, assume that  $T$  maps  $\mathfrak{S}[\mathbf{x}]^{\text{MA}}$  into  $\mathfrak{S}[\mathbf{x}]$ . Then for any  $\mathbf{w} \in \mathcal{H}^m$ ,  $(\mathbf{x} + \mathbf{w})^{[m]} \in \mathfrak{S}[\mathbf{x}]^{\text{MA}}$ , so that  $T((\mathbf{x} + \mathbf{w})^{[m]}) \in \mathfrak{S}[\mathbf{x}]$ .

First, assume that there is a  $\mathbf{w} \in \mathcal{H}^m$  for which  $T((\mathbf{x} + \mathbf{w})^{[m]}) \equiv 0$  identically. For any  $f \in \mathbb{C}[\mathbf{x}]^{\text{MA}}$  let  $\epsilon > 0$  be as in Exercise 3.3. Then  $\epsilon T(f) = T((\mathbf{x} + \mathbf{w})^{[m]} + \epsilon f)$  is stable, so that  $T(f)$  is stable. Thus, the image of  $\mathbb{C}[\mathbf{x}]^{\text{MA}}$  under  $T$  is a  $\mathbb{C}$ -subspace of  $\mathfrak{S}[\mathbf{x}]$ . By Proposition 2.12(b),  $T$  has the form of case (a).

Secondly, if  $T((\mathbf{x} + \mathbf{w})^{[m]}) \neq 0$  for all  $\mathbf{w} \in \mathcal{H}^m$  then, since each of these polynomials is in  $\mathfrak{S}[\mathbf{x}]$ , we have  $T((\mathbf{x} + \mathbf{w})^{[m]})|_{\mathbf{x}=\mathbf{z}} \neq 0$  for all  $\mathbf{z} \in \mathcal{H}^m$  and  $\mathbf{w} \in \mathcal{H}^m$ . This shows that  $T((\mathbf{x} + \mathbf{y})^{[m]})$  is stable in  $\mathfrak{S}[\mathbf{x}, \mathbf{y}]$ , which is the form of case (b).  $\square$

Theorem 3.4 has a corresponding real form – the proof is completely analogous.

**Theorem 3.5** (Theorem 1.2 of [5]). *Let  $T : \mathbb{R}[\mathbf{x}]^{\text{MA}} \rightarrow \mathbb{R}[\mathbf{x}]$  be a linear transformation. Then  $T$  maps  $\mathfrak{S}_{\mathbb{R}}[\mathbf{x}]^{\text{MA}}$  into  $\mathfrak{S}_{\mathbb{R}}[\mathbf{x}]$  if and only if either*

- (a)  $T(f) = \eta(f) \cdot p + \xi(f) \cdot q$  for some linear functionals  $\eta, \xi : \mathbb{R}[\mathbf{x}]^{\text{MA}} \rightarrow \mathbb{R}$  and  $p, q \in \mathfrak{S}_{\mathbb{R}}[\mathbf{x}]$  such that  $p \ll q$ , or
- (b) the polynomial  $T((\mathbf{x} + \mathbf{y})^{[m]})$  is real stable in  $\mathfrak{S}_{\mathbb{R}}[\mathbf{x}, \mathbf{y}]$ , or
- (c) the polynomial  $T((\mathbf{x} - \mathbf{y})^{[m]})$  is real stable in  $\mathfrak{S}_{\mathbb{R}}[\mathbf{x}, \mathbf{y}]$ .

*Proof.* **Exercise 3.6.**  $\square$

#### 4. THE GRACE-WALSH-SZEGŐ COINCIDENCE THEOREM.

Let  $f \in \mathbb{C}[x]$  be a univariate polynomial of degree at most  $m$ , and let  $\mathbf{x} = (x_1, \dots, x_m)$  as usual. For  $0 \leq j \leq m$ , the  $j$ -th elementary symmetric function of  $\mathbf{x}$  is

$$e_j(\mathbf{x}) = \sum_{1 \leq i_1 < \dots < i_j \leq m} x_{i_1} \cdots x_{i_j} = \sum_{S \subseteq [m]: |S|=j} \mathbf{x}^S.$$

The  $m$ -th polarization of  $f$  is the polynomial obtained as the image of  $f$  under the linear transformation  $\text{Pol}_m$  defined by  $x^j \mapsto \binom{m}{j}^{-1} e_j(\mathbf{x})$  for all  $0 \leq j \leq m$ , and linear extension. In other words,  $\text{Pol}_m f$  is the unique multiaffine polynomial in  $\mathbb{C}[\mathbf{x}]^{\text{MA}}$  that is invariant under all permutations of  $[m]$  and such that  $\text{Pol}_m f(x, \dots, x) = f(x)$ . A circular region is a nonempty subset  $\mathcal{A}$  of  $\mathbb{C}$  that is either open or closed, and which is bounded by either a circle or a straight line.

**Theorem 4.1** (Grace-Walsh-Szegő, Theorem 3.4.1b of [14]). *Let  $f \in \mathbb{C}[x]$  have degree at most  $m$  and let  $\mathcal{A}$  be a circular region. If either  $\deg(f) = m$  or  $\mathcal{A}$  is convex, then for every  $\mathbf{z} \in \mathcal{A}^m$  there exists  $z \in \mathcal{A}$  such that  $\text{Pol}_m f(\mathbf{z}) = f(z)$ .*

Figure 1 illustrates the Grace-Walsh-Szegő (GWS) Theorem for the polynomial  $f(x) = x^5 + 10x^2 + 1$ . The black dots mark the solutions to  $f(x) = 0$ . Any permutation of the red (grey) dots is a solution to  $\text{Pol}_5 f(x_1, \dots, x_5) = 0$ . By GWS, any circular region containing all the red dots must contain at least one of the black dots. The figure indicates the boundaries of several circular regions for which this condition is met.

The proof of GWS in this section is adapted from Borcea and Brändén [6].

**4.1. Reduction to the case of stable polynomials.** First of all, it suffices to prove GWS for open circular regions, since a closed circular region is the intersection of all the open circular regions which contain it. Second, it suffices to show that for any  $g \in \mathbb{C}[x]$  of degree at most  $m$ , if  $\deg(g) = m$  or  $\mathcal{A}$  is convex, and  $\mathbf{z} \in \mathcal{A}^m$  is such that  $\text{Pol}_m g(\mathbf{z}) = 0$ , then there exists  $z \in \mathcal{A}$  such that  $g(z) = 0$ . This implies the stated form of GWS by applying this special case to  $g(x) = f(x) - c$ , where  $c = \text{Pol}_m f(\mathbf{z})$ . Stated otherwise, it suffices to show that if  $f(z) \neq 0$  for all  $z \in \mathcal{A}$  then  $\text{Pol}_m f(\mathbf{z}) \neq 0$  for all  $\mathbf{z} \in \mathcal{A}^m$  (provided that either  $\deg(f) = m$  or  $\mathcal{A}$  is convex).

Let  $\mathcal{M}$  be the set of Möbius transformations  $z \mapsto \phi(z) = (az + b)/(cz + d)$  with  $a, b, c, d \in \mathbb{C}$  and  $ab - cd = \pm 1$ . Then  $\mathcal{M}$  with the operation of functional composition is a group of conformal transformations of the Riemann sphere  $\widehat{\mathbb{C}} = \mathbb{C} \cup \{\infty\}$ ,



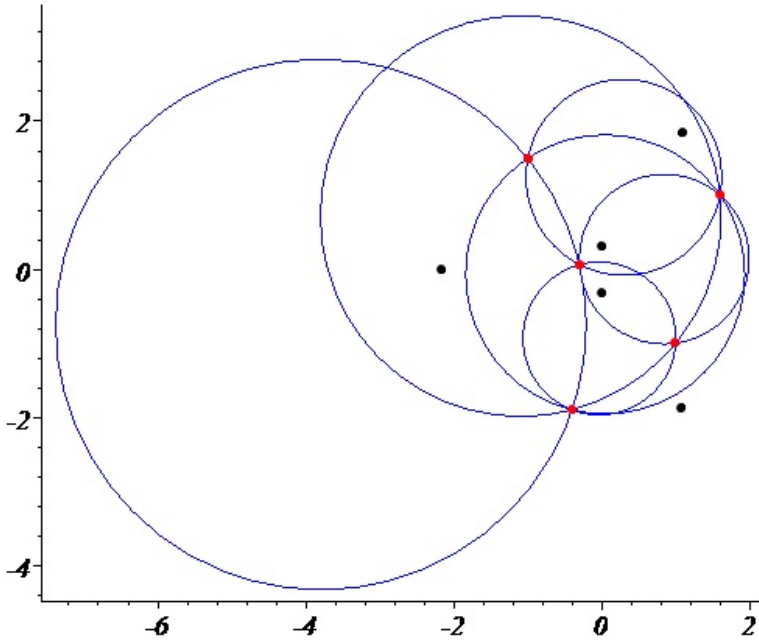


FIGURE 1. Illustration of the Grace-Walsh-Szegő Theorem.

and it acts simply transitively on the set of all ordered triples of distinct points of  $\widehat{\mathbb{C}}$ . Consequently, for any open circular region  $\mathcal{A}$  there is a  $\phi \in \mathcal{M}$  such that  $\phi(\mathcal{H}) = \{\phi(z) : z \in \mathcal{H}\} = \mathcal{A}$ . We henceforth regard circular regions as subsets of  $\widehat{\mathbb{C}}$ . Note that an open circular region  $\mathcal{A}$  is convex if and only if it does not contain  $\infty$ . (The point  $\infty$  is on the boundary of any open half-plane.) In this case, if  $\phi(z) = (az + b)/(cz + d)$  is such that  $\phi(\mathcal{H}) = \mathcal{A}$  then  $cz + d \neq 0$  for all  $z \in \mathcal{H}$ .

Given  $0 \neq f \in \mathbb{C}[x]$  of degree at most  $m$ , consider the polynomial  $\tilde{f}(x) = (cx + d)^m f((ax + b)/(cx + d))$ . If either  $\deg(f) = m$  or  $\mathcal{A}$  is convex, then  $f$  is nonvanishing on  $\mathcal{A}$  if and only if  $\tilde{f}(z)$  is nonvanishing on  $\mathcal{H}$ . Also,

$$\text{Pol}_m \tilde{f}(\mathbf{x}) = \text{Pol}_m f(\phi(x_1), \dots, \phi(x_m)) \cdot \prod_{i=1}^m (cx_i + d).$$

Thus, to prove GWS it suffices to prove the following lemma.

**Lemma 4.2.** *Let  $f \in \mathbb{C}[x]$  be a univariate polynomial of degree at most  $m$ . Then  $\text{Pol}_m f$  is stable if and only if  $f$  is stable.*

Clearly, diagonalization implies that if  $\text{Pol}_m f$  is stable then  $f$  is stable, so only the converse implication needs proof. This is accomplished in the following two easy steps.

**4.2. Partial symmetrization.** The group  $\mathfrak{S}(m)$  of all permutations  $\sigma : [m] \rightarrow [m]$  acts on  $\mathbb{C}[\mathbf{x}]$  by the rule  $\sigma(f)(x_1, \dots, x_m) = f(x_{\sigma(1)}, \dots, x_{\sigma(m)})$ . Notice that

$$\sigma(\mathbf{x}^\alpha) = \prod_{i=1}^m x_{\sigma(i)}^{\alpha(i)} = \prod_{i=1}^m x_i^{\alpha \circ \sigma^{-1}(i)} = \mathbf{x}^{\alpha \circ \sigma^{-1}}.$$

For  $\{i, j\} \subseteq [m]$ , let  $\tau_{ij}$  be the transposition that exchanges  $i$  and  $j$  and fixes all other elements of  $[m]$ .

**Lemma 4.3.** *Let  $0 \leq \lambda \leq 1$  and  $\{i, j\} \subseteq [m]$ , and let  $T_{ij}^{(\lambda)} = (1 - \lambda) + \lambda\tau_{ij}$ . If  $f \in \mathfrak{S}[\mathbf{x}]^{\text{MA}}$  is stable and multiaffine then  $T_{ij}^{(\lambda)}f \in \mathfrak{S}[\mathbf{x}]^{\text{MA}}$  is stable and multiaffine.*

*Proof.* If  $f$  is multiaffine then  $(1 - \lambda)f + \lambda\tau_{ij}(f)$  is also multiaffine. We apply Theorem 3.4 to show that  $T = T_{ij}^{(\lambda)}$  preserves stability of multiaffine polynomials. By permutation we can assume that  $\{i, j\} = \{1, 2\}$ . The algebraic symbol of  $T$  is

$$T((\mathbf{x} + \mathbf{y})^{[m]}) = T((x_1 + y_1)(x_2 + y_2)) \cdot \prod_{i=3}^m (x_i + y_i).$$

Clearly, this is stable if and only if the same is true of  $T((x_1 + y_1)(x_2 + y_2))$ . Exercise 4.4 completes the proof.  $\square$

**Exercise 4.4.** Use the results of Sections 2.4 or 3.1 to show that for  $0 \leq \lambda \leq 1$ , the polynomial

$$x_1x_2 + ((1 - \lambda)x_1 + \lambda x_2)y_2 + (\lambda x_1 + (1 - \lambda)x_2)y_1 + y_1y_2$$

is real stable.

**4.3. Convergence to the polarization.** Let  $0 \neq f(x) \in \mathfrak{S}[x]$  be a univariate stable polynomial of degree at most  $m$ : say  $f(x) = c(x - \xi_1) \cdots (x - \xi_n)$  in which  $c \neq 0$ ,  $n \leq m$ , and  $\xi_i \notin \mathcal{H}$  for all  $i \in [n]$ . Then the polynomial  $F_0 \in \mathbb{C}[\mathbf{x}]$  defined by

$$F_0(x_1, \dots, x_m) = c(x_1 - \xi_1) \cdots (x_n - \xi_n)$$

is multiaffine and stable, and  $F_0(x, \dots, x) = f(x)$ . Let  $\Sigma = (\{i_k, j_k\} : k \in \mathbb{N})$  be a sequence of two-element subsets of  $[m]$ , and for each  $k \in \mathbb{N}$  let  $T_k = T_{i_k j_k}^{(1/2)}$  and define  $F_{k+1} = T_k(F_k)$ . By induction using Lemma 4.3, each  $F_k \in \mathfrak{S}[\mathbf{x}]^{\text{MA}}$  is multiaffine and stable, and  $F_k(x, \dots, x) = f(x)$  for all  $k \in \mathbb{N}$ . We will construct such a sequence  $\Sigma$  for which  $(F_k : k \in \mathbb{N})$  converges to  $\text{Pol}_m f$ .

Let  $P \in \mathbb{C}[\mathbf{x}]^{\text{MA}}$  be multiaffine, say  $P(\mathbf{x}) = \sum_{S \subseteq [m]} c(S)\mathbf{x}^S$ . For  $\{i, j\} \subseteq [m]$  let

$$\omega_{ij}(P) = \sum_{S \subseteq [m]} |c(S) - c(\tau_{ij}(S))|$$

be the  $ij$ -th imbalance of  $P$ , and let  $\|P\| = \sum_{\{i, j\} \subseteq [m]} \omega_{ij}(P)$  be the total imbalance of  $P$ .

**Exercise 4.5.** (a) Let  $(P_k : k \in \mathbb{N})$  be polynomials in  $\mathbb{C}[\mathbf{x}]^{\text{MA}}$  for which there is a  $p \in \mathbb{C}[x]$  such that  $P_k(x, \dots, x) = p(x)$  for all  $k \in \mathbb{N}$ . If  $\|P_k\| \rightarrow 0$  as  $k \rightarrow \infty$ , then  $(P_k : k \in \mathbb{N})$  converges to a limit  $P \in \mathbb{C}[\mathbf{x}]^{\text{MA}}$ , and  $\|P\| = 0$ .

(b) For  $P \in \mathbb{C}[\mathbf{x}]^{\text{MA}}$ ,  $\|P\| = 0$  if and only if  $P$  is invariant under all permutations of  $[m]$ . Thus, in part (a) the limit is  $P = \text{Pol}_m p$ .

**Exercise 4.6.** Let  $P \in \mathbb{C}[\mathbf{x}]^{\text{MA}}$ , let  $\{i, j\} \subseteq [m]$ , and let  $Q = T_{ij}^{(1/2)}P$ .

- (a) Then  $\omega_{ij}(Q) = 0$ .
- (b) If  $h \in [m] \setminus \{i, j\}$  then  $\omega_{hi}(Q) \leq (\omega_{hi}(P) + \omega_{hj}(P))/2$ , and similarly for  $\omega_{hj}(Q)$ .
- (c) If  $\{h, k\} \subseteq [m] \setminus \{i, j\}$  then  $\omega_{hk}(Q) = \omega_{hk}(P)$ .
- (d) Consequently,  $\|Q\| \leq \|P\| - \omega_{ij}(P)$ .

Now we choose the sequence  $\Sigma = (\{i_k, j_k\} : k \in \mathbb{N})$  as follows: for each  $k \in \mathbb{N}$ ,  $\{i_k, j_k\} \subseteq [m]$  is any pair of indices  $\{i, j\}$  for which  $\omega_{ij}(F_k)$  attains its maximum value. Then  $\omega_{i_k j_k}(F_k) \geq \binom{m}{2}^{-1} \|F_k\|$ , so that by Exercise 4.6(d) and induction on  $k \in \mathbb{N}$ ,

$$\|F_{k+1}\| \leq \left(1 - \binom{m}{2}^{-1}\right) \|F_k\| \leq \left(1 - \binom{m}{2}^{-1}\right)^{k+1} \|F_0\|.$$

Thus, by Exercise 4.5,  $F_k$  converges to  $\text{Pol}_m f$ , the  $m$ -th polarization of  $f$ . Finally, since each  $F_k$  is stable (and the limit is a polynomial), Hurwitz's Theorem implies that  $\text{Pol}_m f$  is stable. This completes the proof of Lemma 4.2, and hence of Theorem 4.1.

## 5. POLARIZATION ARGUMENTS AND STABILITY PRESERVERS.

For  $\kappa \in \mathbb{N}^m$  and a set  $\mathcal{S} \subseteq \mathbb{C}[\mathbf{x}]$  of polynomials, let  $\mathcal{S}^{\leq \kappa}$  be the set of all  $f \in \mathcal{S}$  such that  $\deg_i(f) \leq \kappa(i)$  for all  $i \in [m]$ . Let

$$I(\kappa) = \{(i, j) : i \in [m] \text{ and } j \in [\kappa(i)]\}$$

and let  $\mathbf{u} = \{u_{ij} : (i, j) \in I(\kappa)\}$  be indeterminates. For  $f \in \mathbb{C}[\mathbf{x}]^{\leq \kappa}$ , let  $\text{Pol}_{\kappa(i)}^{(i)} f$  denote the  $\kappa(i)$ -th polarization of  $x_i$  in  $f$ : this is the image of  $f$  under the linear transformation  $\text{Pol}_{\kappa(i)}^{(i)}$  defined by  $x_i^j \mapsto \binom{\kappa(i)}{j}^{-1} e_j(u_{i1}, \dots, u_{i\kappa(i)})$  for each  $0 \leq j \leq \kappa(i)$ , and linear extension. Finally, the  $\kappa$ -th polarization of  $f$  is

$$\text{Pol}_{\kappa} f = \text{Pol}_{\kappa(m)}^{(m)} \circ \dots \circ \text{Pol}_{\kappa(1)}^{(1)} f.$$

This defines a linear transformation  $\text{Pol}_{\kappa} : \mathbb{C}[\mathbf{x}]^{\leq \kappa} \rightarrow \mathbb{C}[\mathbf{u}]^{\text{MA}}$ .

### 5.1. The real stability criterion revisited.

**Proposition 5.1.** *Let  $\kappa \in \mathbb{N}^m$  and  $f \in \mathbb{C}[x]^{\leq \kappa}$ . Then  $\text{Pol}_{\kappa} f$  is stable if and only if  $f$  is stable.*

*Proof.* Diagonalization implies that if  $\text{Pol}_{\kappa} f$  is stable then  $f$  is stable, so only the converse implication needs proof. Assume that  $f$  is stable, and let  $z_{ij} \in \mathcal{H}$  for  $(i, j) \in I(\kappa)$ . By induction on  $m$ , repeated application of GWS shows that there are  $\mathbf{z} = (z_1, \dots, z_m) \in \mathcal{H}^m$  such that

$$\text{Pol}_{\kappa} f(z_{ij} : (i, j) \in I(\kappa)) = f(\mathbf{z}).$$

Since  $f$  is stable it follows that  $\text{Pol}_{\kappa} f$  is stable. □

If  $f \in \mathbb{R}[\mathbf{x}]^{\leq \kappa}$  then Theorem 3.1 applies to  $\text{Pol}_{\kappa} f$ . Thus, Proposition 5.1 bootstraps the real stability criterion from multiaffine to arbitrary polynomials. This is a typical application of the GWS Theorem.

### 5.2. Linear transformations preserving stability – polynomial case.

**Theorem 5.2** (Theorem 1.1 of [5]). *Let  $\kappa \in \mathbb{N}^m$ , and let  $T : \mathbb{C}[\mathbf{x}]^{\leq \kappa} \rightarrow \mathbb{C}[\mathbf{x}]$  be a linear transformation. Then  $T$  maps  $\mathfrak{S}[\mathbf{x}]^{\leq \kappa}$  into  $\mathfrak{S}[\mathbf{x}]$  if and only if either*

- (a)  $T(f) = \eta(f) \cdot p$  for some linear functional  $\eta : \mathbb{C}[\mathbf{x}]^{\leq \kappa} \rightarrow \mathbb{C}$  and  $p \in \mathfrak{S}[\mathbf{x}]$ , or
- (b) the polynomial  $T((\mathbf{x} + \mathbf{y})^\kappa)$  is stable in  $\mathfrak{S}[\mathbf{x}, \mathbf{y}]$ .

*Proof.* Let  $\mathbf{u} = \{u_{ij} : (i, j) \in I(\kappa)\}$ , and define a linear transformation  $\tilde{T} : \mathbb{C}[\mathbf{u}]^{\text{MA}} \rightarrow \mathbb{C}[\mathbf{x}]$  as follows. For every  $A \subseteq I(\kappa)$ , define  $\alpha(A) : [m] \rightarrow \mathbb{N}$  by putting  $\alpha(A, i) = |\{j \in [\kappa(i)] : (i, j) \in A\}|$  for each  $i \in [m]$ . Then for each  $A \subseteq I(\kappa)$  define  $\tilde{T}(\mathbf{u}^A) = T(\mathbf{x}^{\alpha(A)})$ , and extend this linearly to all of  $\mathbb{C}[\mathbf{u}]^{\text{MA}}$ . Let  $\Delta : \mathbb{C}[\mathbf{u}]^{\text{MA}} \rightarrow \mathbb{C}[\mathbf{x}]$  be the diagonalization operator defined by  $\Delta(u_{ij}) = x_i$  for all  $(i, j) \in I(\kappa)$ , extended algebraically.

Notice that  $T = \tilde{T} \circ \text{Pol}_\kappa$ , and that  $\tilde{T} = T \circ \Delta$ . By Proposition 5.1 (and Lemma 2.4), it follows that  $T$  preserves stability if and only if  $\tilde{T}$  preserves stability. This is equivalent to one of two cases in Theorem 3.4.

In case (a), if  $\tilde{T} = p \cdot \tilde{\eta}$  for some  $p \in \mathfrak{S}[\mathbf{x}]$  and linear functional  $\tilde{\eta} : \mathbb{C}[\mathbf{y}]^{\text{MA}} \rightarrow \mathbb{C}$  then  $T = p \cdot (\eta \circ \text{Pol}_\kappa)$  is also in case (a). Conversely, if  $T$  is in case (a) then the same is true of  $\tilde{T}$ , by construction.

In case (b), let  $\text{Pol}_\kappa^{(\mathbf{y})} : \mathbb{C}[\mathbf{y}]^{\leq \kappa} \rightarrow \mathbb{C}[\mathbf{v}]^{\text{MA}}$  denote the  $\kappa$ -th polarization of the  $\mathbf{y}$  variables. The symbols of  $T$  and  $\tilde{T}$  are related by

$$\tilde{T}((\mathbf{u} + \mathbf{v})^{I(\kappa)}) = (T \circ \Delta)((\mathbf{u} + \mathbf{v})^{I(\kappa)}) = \text{Pol}_\kappa^{(\mathbf{y})} T((\mathbf{x} + \mathbf{y})^\kappa),$$

and Proposition 5.1 shows that  $T$  is in case (b) if and only if  $\tilde{T}$  is in case (b).  $\square$

### 5.3. Linear transformations preserving stability – transcendental case.

**Exercise 5.3.** Let  $T : \mathbb{C}[\mathbf{x}] \rightarrow \mathbb{C}[\mathbf{x}]$  be a linear transformation.

- (a) Then  $T : \mathfrak{S}[\mathbf{x}] \rightarrow \mathfrak{S}[\mathbf{x}]$  if and only if  $T : \mathfrak{S}[\mathbf{x}]^{\leq \kappa} \rightarrow \mathfrak{S}[\mathbf{x}]$  for all  $\kappa \in \mathbb{N}^m$ .
- (b) Define  $S : \mathbb{C}[\mathbf{x}, \mathbf{y}] \rightarrow \mathbb{C}[\mathbf{x}, \mathbf{y}]$  by  $S(\mathbf{x}^\alpha \mathbf{y}^\beta) = T(\mathbf{x}^\alpha) \mathbf{y}^\beta$  and linear extension. If  $T((\mathbf{x} + \mathbf{u})^\kappa)$  is stable for all  $\kappa \in \mathbb{N}^m$  then  $S((\mathbf{x} + \mathbf{u})^\kappa (\mathbf{y} + \mathbf{v})^\beta)$  is stable for all  $\kappa, \beta \in \mathbb{N}^m$ .

Let  $\overline{\mathfrak{S}[\mathbf{x}]}$  denote the set of all power series in  $\mathbb{C}[[\mathbf{x}]]$  that are obtained as the limit of a sequence of stable polynomials in  $\mathfrak{S}[\mathbf{x}]$  which converges uniformly on compact sets. Theorem 5.4 is an astounding generalization of the Pólya-Schur Theorem. For  $\alpha \in \mathbb{N}^m$ , let  $\alpha! = \prod_{i=1}^m \alpha(i)!$ .

**Theorem 5.4** (Theorem 1.3 of [5]). *Let  $T : \mathbb{C}[\mathbf{x}] \rightarrow \mathbb{C}[\mathbf{x}]$  be a linear transformation. Then  $T$  maps  $\overline{\mathfrak{S}[\mathbf{x}]}$  into  $\overline{\mathfrak{S}[\mathbf{x}]}$  if and only if either*

- (a)  $T(f) = \eta(f) \cdot p$  for some linear functional  $\eta : \mathbb{C}[\mathbf{x}] \rightarrow \mathbb{C}$  and  $p \in \overline{\mathfrak{S}[\mathbf{x}]}$ , or
- (b) the power series

$$T(e^{-\mathbf{x}\mathbf{y}}) = \sum_{\alpha: [m] \rightarrow \mathbb{N}} (-1)^\alpha T(\mathbf{x}^\alpha) \frac{\mathbf{y}^\alpha}{\alpha!}$$

is in  $\overline{\mathfrak{S}[\mathbf{x}, \mathbf{y}]}$

(Theorem 3.5 has a similar extension – see Theorems 1.2 and 1.4 of [5].)

For  $\alpha \leq \beta$  in  $\mathbb{N}^m$ , let  $(\beta)_\alpha = \beta! / (\beta - \alpha)!$ , and for  $\alpha \not\leq \beta$  let  $(\beta)_\alpha = 0$ .

**Theorem 5.5** (Theorem 5.1 of [5]). *Let  $F(\mathbf{x}, \mathbf{y}) = \sum_{\alpha \in \mathbb{N}^m} P_\alpha(\mathbf{x}) \mathbf{y}^\alpha$  be a power series in  $\mathbb{C}[\mathbf{x}][[\mathbf{y}]]$  (so that each  $P_\alpha \in \mathbb{C}[\mathbf{x}]$ ). Then  $F(\mathbf{x}, \mathbf{y})$  is in  $\mathfrak{S}[\mathbf{x}, \mathbf{y}]$  if and only if for all  $\beta \in \mathbb{N}^m$ ,*

$$\sum_{\alpha \leq \beta} (\beta)_\alpha P_\alpha(\mathbf{x}) \mathbf{y}^\alpha$$

*is stable in  $\mathfrak{S}[\mathbf{x}, \mathbf{y}]$ .*

(This implies the analogous result for real stability, since  $\mathfrak{S}_{\mathbb{R}}[\mathbf{x}] = \mathfrak{S}[\mathbf{x}] \cap \mathbb{R}[\mathbf{x}]$ .)

**Exercise 5.6.** Derive Theorem 5.4 from Theorems 5.2 and 5.5. (Hint:  $T((\mathbf{x} + \mathbf{y})^\kappa)$  is stable if and only if  $T((\mathbf{1} - \mathbf{x}\mathbf{y})^\kappa)$  is stable.)

One direction of Theorem 5.5 is relatively straightforward.

**Lemma 5.7** (Lemma 5.2 of [5]). *Fix  $\beta \in \mathbb{N}^m$ . The linear transformation  $T : \mathbf{y}^\alpha \mapsto (\beta)_\alpha \mathbf{y}^\alpha$  on  $\mathbb{C}[\mathbf{y}]$  preserves stability.*

*Proof.* By Theorem 5.2 and Exercise 5.3(a), it suffices to show that for all  $\kappa \in \mathbb{N}^m$ , the polynomial  $T((\mathbf{y} + \mathbf{u})^\kappa)$  is stable. Now

$$T((\mathbf{y} + \mathbf{u})^\kappa) = \prod_{i=1}^m \left[ \sum_{j=0}^{\kappa(i)} j! \binom{\kappa(i)}{j} \binom{\beta(i)}{j} y_i^j u_i^{\kappa(i)-j} \right],$$

so it suffices to show that for all  $k, b \in \mathbb{N}$ , the polynomial  $f(t) = \sum_{j=0}^k j! \binom{k}{j} \binom{b}{j} t^j$  is real stable. Let  $g(t) = (1 + d/dt)^k t^b$ . One can check that  $f(t) = t^b g(1/t)$ . It thus suffices to show that  $1 + d/dt$  preserves stability. For any  $a \in \mathbb{N}$ ,  $(1 + d/dt)(t + u)^a = (t + u + a)(t + u)^{a-1}$  is stable, and so Theorem 5.2 implies the result.  $\square$

Now, let  $F = F(\mathbf{x}, \mathbf{y})$  be as in the statement of Theorem 5.5, and let  $(F_n : n \in \mathbb{N})$  be a sequence of stable polynomials  $F_n(\mathbf{x}, \mathbf{y}) = \sum_{\alpha \in \mathbb{N}^m} P_{n,\alpha}(\mathbf{x}) \mathbf{y}^\alpha$  in  $\mathfrak{S}[\mathbf{x}, \mathbf{y}]$  converging to  $F$  uniformly on compact sets. Fix  $\beta \in \mathbb{N}$  and define a linear transformation  $T : \mathbb{C}[\mathbf{x}, \mathbf{y}] \rightarrow \mathbb{C}[\mathbf{x}, \mathbf{y}]$  by  $T(\mathbf{x}^\gamma \mathbf{y}^\alpha) = (\beta)_\alpha \mathbf{x}^\gamma \mathbf{y}^\alpha$  and linear extension. By Lemma 5.7 and Exercise 5.3,  $T$  preserves stability in  $\mathfrak{S}[\mathbf{x}, \mathbf{y}]$ . Thus,  $(T(F_n) : n \in \mathbb{N})$  is a sequence of stable polynomials converging to  $T(F)$ . Since  $T(F)$  is a polynomial the convergence is uniform on compact sets, and so Hurwitz's Theorem implies that  $T(F)$  is stable.

The converse direction of Theorem 5.5 is considerably more technical, although the idea is simple. With  $F$  as in the theorem, for each  $n \geq 1$  let

$$F_n(\mathbf{x}, \mathbf{y}) = \sum_{\alpha \leq n\mathbf{1}} (n\mathbf{1})_\alpha P_\alpha(\mathbf{x}) \frac{\mathbf{y}^\alpha}{n^\alpha}.$$

The sequence  $(F_n : n \geq 1)$  converges to  $F$ , since for each  $\alpha \in \mathbb{N}^m$ ,  $n^{-\alpha} (n\mathbf{1})_\alpha \rightarrow 1$  as  $n \rightarrow \infty$ . Each  $F_n$  is stable, by hypothesis (and scaling). The hard work is involved with showing that the convergence is uniform on compact sets. To do this, Borcea and Brändén develop a very flexible multivariate generalization of the Szasz Principle [5, Theorem 5.6] – in itself an impressive accomplishment. Unfortunately, we have no space here to develop this result – see Section 5.2 of [5].

6. JOHNSON'S CONJECTURES.

Let  $\mathcal{A} = (A_1, \dots, A_k)$  be a  $k$ -tuple of  $n$ -by- $n$  matrices. Define the *mixed determinant* of  $\mathcal{A}$  to be

$$\text{Det}(\mathcal{A}) = \text{Det}(A_1, \dots, A_k) = \sum_{(S_1, \dots, S_k)} \prod_{i=1}^k \det A_i[S_i],$$

in which the sum is over all ordered sequences of  $k$  pairwise disjoint subsets of  $[n]$  such that  $[n] = S_1 \cup \dots \cup S_k$ , and  $A_i[S_i]$  is the principal submatrix of  $A_i$  supported on rows and columns in  $S_i$ . Let  $A_i(S_i)$  be the complementary principal submatrix supported on rows and columns not in  $S_i$ , and for  $j \in [n]$  let  $A_i(j) = A_i(\{j\})$ .

For example, when  $k = 2$  and  $A_1 = xI$  and  $A_2 = -B$ , this specializes to  $\text{Det}(xI, -B) = \det(xI - B)$ , the characteristic polynomial of  $B$ . In the late 1980s, Johnson made three conjectures about the  $k = 2$  case more generally.

**Johnson's Conjectures.** *Let  $A$  and  $B$  be  $n$ -by- $n$  matrices, with  $A$  positive definite and  $B$  Hermitian.*

- (a) *Then  $\text{Det}(xA, -B)$  has only real roots.*
- (b) *For  $j \in [n]$ , the roots of  $\text{Det}(xA(j), -B(j))$  interlace those of  $\text{Det}(xA, -B)$ .*
- (c) *The inertia of  $\text{Det}(xA, -B)$  is the same as that of  $\det(xI - B)$ .*

In part (c), the *inertia* of a univariate real stable polynomial  $p$  is the triple  $\iota(p) = (\iota_-(p), \iota_0(p), \iota_+(p))$  with entries the number of negative, zero, or positive roots of  $p$ , respectively.

In 2008, Borcea and Brändén [1] proved all three of these statements in much greater generality.

**Theorem 6.1** (Theorem 2.6 of [1]). *Fix integers  $\ell, m, n \geq 1$ . For  $h \in [\ell]$  and  $i \in [m]$  let  $B_h$  and  $A_{hi}$  be  $n$ -by- $n$  matrices, and let*

$$L_h = \sum_{i=1}^m x_i A_{hi} + B_h.$$

- (a) *If all the  $A_{hi}$  are positive semidefinite and all the  $B_h$  are Hermitian, then  $\text{Det}(\mathcal{L}) = \text{Det}(L_1, \dots, L_\ell) \in \mathfrak{S}_{\mathbb{R}}[\mathbf{x}]$  is real stable.*
- (b) *For each  $j \in [n]$ , let  $\mathcal{L}(j) = (L_1(j), \dots, L_\ell(j))$ . With the hypotheses of part (a), the polynomial  $\text{Det}(\mathcal{L}) + y\text{Det}(\mathcal{L}(j)) \in \mathfrak{S}_{\mathbb{R}}[\mathbf{x}, y]$  is real stable.*

*Proof.* Let  $Y = \text{diag}(y_1, \dots, y_n)$  be a diagonal matrix of indeterminates. By Proposition 2.1, for each  $h \in [\ell]$  the polynomial

$$\det(Y + L_h) = \sum_{S \subseteq [n]} \mathbf{y}^S \det L_h(S)$$

is real stable in  $\mathfrak{S}_{\mathbb{R}}[\mathbf{x}, \mathbf{y}]$ . By inversion of all the  $\mathbf{y}$  indeterminates, each

$$\det(I - YL_h) = \sum_{S \subseteq [n]} (-1)^{|S|} \mathbf{y}^S \det L_h[S]$$

is real stable. Since  $\prod_{h=1}^{\ell} \det(I - YL_h)$  is real stable, contraction and specialization imply that

$$\text{Det}(\mathcal{L}) = (-1)^n \frac{\partial^n}{\partial y_1 \cdots \partial y_n} \prod_{h=1}^{\ell} \det(I - YL_h) \Big|_{\mathbf{y}=0}$$

is real stable, proving part (a).

For part (b), let  $V$  be the  $n$ -by- $n$  matrix with all entries zero except for  $V_{jj} = y$ . By part (a),

$$\text{Det}(V, L_1, \dots, L_h) = \text{Det}(\mathcal{L}) + y \text{Det}(\mathcal{L}(j))$$

is real stable.  $\square$

Theorem 6.1 (with Corollary 2.10) clearly settles Conjectures (a) and (b).

*Proof of Conjecture (c).* Let  $A$  and  $B$  be  $n$ -by- $n$  matrices with  $A$  positive definite and  $B$  Hermitian. Let  $(\iota_-, \iota_0, \iota_+)$  be the inertia of  $\det(xI - B)$ . Let  $f(x) = \text{Det}(xA, -B)$ , and let  $(\nu_-, \nu_0, \nu_+)$  be the inertia of  $f$ .

We begin by showing that  $\nu_0 = \iota_0$ . Since  $\iota_0 = \min\{|S| : S \subseteq [n] \text{ and } \det(B(S)) \neq 0\}$ , it follows that  $\nu_0 \geq \iota_0$ . The constant term of  $f(x)$  is  $(-1)^n \det(B)$ , so that if  $\iota_0 = 0$  then  $\nu_0 = 0$ . If  $\iota_0 = k > 0$  then let  $S = \{s_1, \dots, s_k\} \subseteq [n]$  be such that  $\det(B(S)) \neq 0$ . For  $0 \leq i \leq k$  let  $f_i(x) = \text{Det}(A(\{s_1, \dots, s_i\}), -B(\{s_1, \dots, s_i\}))$ , so that  $f_0(x) = f(x)$ . By Theorem 6.1, the roots of  $f_{i-1}$  and of  $f_i$  are interlaced, for each  $i \in [k]$ . Thus,

$$\nu_0 = \iota_0(f_0) \leq \iota_0(f_1) + 1 \leq \iota_0(f_2) + 2 \leq \cdots \leq \iota_0(f_k) + k = k = \iota_0,$$

since  $\iota_0(f_k) = 0$  because  $\det(B(S)) \neq 0$ . Therefore  $\nu_0 = \iota_0$ .

For any positive definite matrix  $A$ ,  $\text{Det}(xA, -B)$  is a polynomial of degree  $n$ . Suppose that  $A$  is such a matrix for which  $\nu_+ \neq \iota_+$ . Consider the matrices  $A_\lambda = (1 - \lambda)I + \lambda A$  for  $\lambda \in [0, 1]$ . Each of these matrices is positive definite. From the paragraph above, each of the polynomials  $g_\lambda(x) = \text{Det}(xA_\lambda, -B)$  has  $\iota_0(g_\lambda) = \iota_0$ . Since  $\iota_+(g_0) = \iota_+ \neq \nu_+ = \iota_+(g_1)$  and the roots of  $g_\lambda$  vary continuously with  $\lambda$ , there is some value  $\mu \in (0, 1)$  for which  $\iota_0(g_\mu) > \iota_0$ . This contradiction shows that  $\nu_+ = \iota_+$ , and hence  $\nu_- = \iota_-$  as well.  $\square$

Borcea and Brändén [1] proceed to derive many inequalities for the principal minors of positive semidefinite matrices, and some for merely Hermitian matrices. These are applications of inequalities valid more generally for real stable polynomials. The simplest of these inequalities are as follows.

For an  $n$ -by- $n$  matrix  $A$ , the  $j$ -th symmetrized Fisher product is

$$\sigma_j(A) = \sum_{S \subseteq [n]: |S|=j} \det(A[S]) \det(A(S)).$$

and the  $j$ -th averaged Fisher product is  $\hat{\sigma}_j(A) = \binom{n}{j}^{-1} \sigma_j(A)$ . Notice that  $\sigma_j(A) = \sigma_{n-j}(A)$  for all  $0 \leq j \leq n$ .

**Corollary 6.2.** *Let  $A$  be an  $n$ -by- $n$  positive semidefinite matrix.*

- (a) *Then  $\hat{\sigma}_j(A)^2 \geq \hat{\sigma}_{j-1}(A) \hat{\sigma}_{j+1}(A)$  for all  $1 \leq j \leq n - 1$ .*
- (b) *Also,  $\hat{\sigma}_0(A) \leq \hat{\sigma}_1(A) \leq \cdots \leq \hat{\sigma}_{\lfloor n/2 \rfloor}$ .*
- (c) *If  $A$  is positive definite and  $\det(A) = d$  then*

$$\frac{\hat{\sigma}_1(A)}{d} \geq \left( \frac{\hat{\sigma}_2(A)}{d} \right)^{1/2} \geq \left( \frac{\hat{\sigma}_3(A)}{d} \right)^{1/3} \geq \cdots \geq \left( \frac{\hat{\sigma}_n(A)}{d} \right)^{1/n} = 1.$$

*Proof.* It suffices to consider positive definite  $A$ . By Theorem 6.1, the polynomial  $\text{Det}(xA, -A) = \sum_{j=0}^n (-1)^j \sigma_j(A) x^j$  has only real roots, and these roots are all positive. Part (a) follows from Newton's Inequalities [12, Theorem 51]. Part (a) and the symmetry  $\sigma_j(A) = \sigma_{n-j}(A)$  for all  $0 \leq j \leq n$  imply part (b). Part (c) follows from Maclaurin's Inequalities [12, Theorem 52].  $\square$

## 7. THE SYMMETRIC EXCLUSION PROCESS.

This section summarizes an application of stable polynomials to probability and statistical mechanics from a 2009 paper of Borcea, Brändén and Liggett [7].

Let  $\Lambda$  be a set of *sites*. A symmetric exclusion process (SEP) is a type of Markov chain with state space a subset of  $\{0, 1\}^\Lambda$ . In a state  $S : \Lambda \rightarrow \{0, 1\}$ , the sites in  $S^{-1}(1)$  are *occupied* and the sites in  $S^{-1}(0)$  are *vacant*. This is meant to model a physical system of particles interacting by means of hard-core exclusions. Such models come in many varieties – to avoid technicalities we discuss only the case of a finite system  $\Lambda$  and continuous time  $t$ . (The results of this section extend to countable  $\Lambda$  under a reasonable finiteness condition on the interaction rates.) Symmetry of the interactions turns out to be crucial, but particle number conservation is unimportant.

Let  $E$  be a set of two-element subsets of  $\Lambda$ . For each  $\{i, j\} \in E$ , let  $\lambda_{ij} > 0$  be a positive real, and let  $\tau_{ij} : \Lambda \rightarrow \Lambda$  be the permutation that exchanges  $i$  and  $j$  and fixes all other sites. Our SEP Markov chain  $\mathcal{M}$  proceeds as follows. Each  $\{i, j\} \in E$  has a Poisson process “clock” of rate  $\lambda_{ij}$ , and these are independent of one another. With probability one, no two clocks ever ring at the same time. When the clock of  $\{i, j\}$  rings, the current state  $S$  is updated to the new state  $S \circ \tau_{ij}$ . In other words, when the  $\{i, j\}$  clock rings, if exactly one of the sites  $\{i, j\}$  is occupied then a particle hops from the occupied to the vacant of these two sites.

Let  $\Lambda = [m]$  and  $\Omega = \{0, 1\}^\Lambda$ , let  $\varphi_0$  be an initial probability distribution on  $\Omega$ , and let  $\varphi_t$  be the distribution of the state of  $\mathcal{M}$ , starting at  $\varphi_0$ , after evolving for time  $t \geq 0$ . We are concerned with properties of the distribution  $\varphi_t$  that hold for all  $t \geq 0$ .

**7.1. Negative correlation and negative association.** Consider a probability distribution  $\varphi$  on  $\Omega$ . An *event*  $\mathcal{E}$  is any subset of  $\Omega$ . The probability of the event  $\mathcal{E}$  is  $\Pr[\mathcal{E}] = \sum_{S \in \mathcal{E}} \varphi(S)$ . An event  $\mathcal{E}$  is *increasing* if whenever  $S \leq S'$  in  $\Omega$  and  $S \in \mathcal{E}$ , then  $S' \in \mathcal{E}$ . For example, if  $K$  is any subset of  $\Lambda$  and  $\mathcal{E}_K$  is the event that all sites in  $K$  are occupied, then  $\mathcal{E}_K$  is an increasing event. Notice that this event has the form  $\mathcal{E}_K = \mathcal{E}' \times \{0, 1\}^{\Lambda \setminus K}$  for some event  $\mathcal{E}' \subseteq \{0, 1\}^K$ . Two events  $\mathcal{E}$  and  $\mathcal{F}$  are *disjointly supported* when one can partition  $\Lambda = A \cup B$  with  $A \cap B = \emptyset$  and  $\mathcal{E} = \mathcal{E}' \times \{0, 1\}^B$  and  $\mathcal{F} = \{0, 1\}^A \times \mathcal{F}'$  for some events  $\mathcal{E}' \subseteq \{0, 1\}^A$  and  $\mathcal{F}' \subseteq \{0, 1\}^B$ .

A probability distribution on  $\Omega$  is *negatively associated* (NA) when  $\Pr[\mathcal{E} \cap \mathcal{F}] \leq \Pr[\mathcal{E}] \cdot \Pr[\mathcal{F}]$  for any two increasing events that are disjointly supported. It is *negatively correlated* (NC) when  $\Pr[\mathcal{E}_{\{i,j\}}] \leq \Pr[\mathcal{E}_{\{i\}}] \cdot \Pr[\mathcal{E}_{\{j\}}]$  for any two distinct sites  $\{i, j\} \subseteq \Lambda$ . Clearly NA implies NC.

It is useful to find conditions under which NC implies NA, since NC is so much easier to check. The following originates with Feder and Mihail, but many others have contributed their insights – see Section 4.2 of [7]. The *partition function* of



any  $\varphi : \Omega \rightarrow \mathbb{R}$  is the real multiaffine polynomial

$$Z(\varphi) = Z(\varphi; \mathbf{x}) = \sum_{S \in \Omega} \varphi(S) \mathbf{x}^S$$

in  $\mathbb{R}[\mathbf{x}]^{\text{MA}}$ . If  $\varphi$  is nonzero and nonnegative, then for any  $\mathbf{a} \in \mathbb{R}^\Lambda$  with  $\mathbf{a} > \mathbf{0}$ , this defines a probability distribution  $\varphi^{\mathbf{a}} : \Omega \rightarrow [0, 1]$  by setting  $\varphi^{\mathbf{a}}(S) = \varphi(S) \mathbf{a}^S / Z(\varphi; \mathbf{a})$  for all  $S \in \Omega$ .

**Feder-Mihail Theorem** (Theorem 4.8 of [7]). *Let  $\mathcal{S}$  be a class of nonzero non-negative functions satisfying the following conditions.*

- (i) *Each  $\varphi \in \mathcal{S}$  has domain  $\{0, 1\}^\Lambda$  for some finite set  $\Lambda = \Lambda(\varphi)$ .*
  - (ii) *For each  $\varphi \in \mathcal{S}$ ,  $Z(\varphi)$  is a homogeneous polynomial.*
  - (iii) *For each  $\varphi \in \mathcal{S}$  and  $i \in \Lambda(\varphi)$ ,  $Z(\varphi)|_{x_i=0}$  and  $\partial_i Z(\varphi)$  are partition functions of members of  $\mathcal{S}$ .*
  - (iv) *For each  $\varphi \in \mathcal{S}$  and  $\mathbf{a} \in \mathbb{R}^{\Lambda(\varphi)}$  with  $\mathbf{a} > \mathbf{0}$ ,  $\varphi^{\mathbf{a}}$  is NC.*
- Then for every  $\varphi \in \mathcal{S}$  and  $\mathbf{a} \in \mathbb{R}^{\Lambda(\varphi)}$  with  $\mathbf{a} > \mathbf{0}$ ,  $\varphi^{\mathbf{a}}$  is NA.*

**7.2. A conjecture of Liggett and Pemantle.** In the early 2000s, Liggett and Pemantle arrived independently at the following conjecture, now a theorem.

**Theorem 7.1** (Theorem 5.2 of [7]). *If the initial distribution  $\varphi_0$  of a SEP is deterministic (i.e. concentrated on a single state) then  $\varphi_t$  is NA for all  $t \geq 0$ .*

*Proof.* This amounts to finding a class  $\mathfrak{S}$  of probability distributions such that:

- (1) deterministic distributions are in  $\mathfrak{S}$ ,
- (2) being in  $\mathfrak{S}$  implies NA, and
- (3) time evolution of the SEP preserves membership in  $\mathfrak{S}$ .

Borcea, Brändén, and Liggett [7] identified such a class:  $\varphi$  is in  $\mathfrak{S}$  if and only if the partition function  $Z(\varphi)$  is homogeneous, multiaffine, and real stable. (Notice that if  $\varphi$  is in  $\mathfrak{S}$  then  $\varphi^{\mathbf{a}}$  is in  $\mathfrak{S}$  for all  $\mathbf{a} \in \mathbb{R}^\Lambda$  with  $\mathbf{a} > \mathbf{0}$ , by scaling.) We proceed to check the three claims above.

Claim (1) is trivial, since if  $\varphi(S) = 1$  then  $Z(\varphi) = \mathbf{x}^S$ , which is clearly homogeneous, multiaffine, and real stable.

To check claim (2) we verify the hypotheses of the Feder-Mihail Theorem. Hypotheses (i) and (ii) hold since  $Z(\varphi)$  is multiaffine and homogeneous. By specialization and contraction, (iii) holds. To check (iv), let  $\mathbf{a} \in \mathbb{R}^\Lambda$  with  $\mathbf{a} > \mathbf{0}$ , let  $\{i, j\} \subseteq \Lambda$ , and consider the probability distribution  $\varphi^{\mathbf{a}}$  on  $\Omega$ . The occupation probability for site  $i$  is

$$\Pr[\mathcal{E}_{\{i\}}] = \sum_{S \in \{0,1\}^\Lambda: S(i)=1} \frac{\varphi(S) \mathbf{a}^S}{Z(\varphi; \mathbf{a})} = a_i \frac{\partial_i Z(\varphi; \mathbf{a})}{Z(\varphi; \mathbf{a})},$$

and similarly for  $\Pr[\mathcal{E}_{\{j\}}]$ . Likewise,  $\Pr[\mathcal{E}_{\{i,j\}}] = a_i a_j Z(\varphi; \mathbf{a})^{-1} \cdot \partial_i \partial_j Z(\varphi; \mathbf{a})$ . Now

$$\Pr[\mathcal{E}_{\{i,j\}}] - \Pr[\mathcal{E}_{\{i\}}] \cdot \Pr[\mathcal{E}_{\{j\}}] = -\frac{a_i a_j}{Z(\varphi; \mathbf{a})^2} \cdot \Delta_{ij} Z(\varphi; \mathbf{a}) \leq 0,$$

by Theorem 3.1. Thus  $\varphi^{\mathbf{a}}$  is NC. By the Feder-Mihail Theorem, every  $\varphi$  in  $\mathfrak{S}$  is NA.

To check claim (3) we need some of the theory of continuous time Markov chains. The time evolution of a Markov chain  $\mathcal{M}$  with finite state space  $\Omega$  is governed by a one-parameter semigroup  $T(t)$  of transformations of  $\mathbb{R}^\Omega$ . For a function  $F \in \mathbb{R}^\Omega$

and time  $t \geq 0$  and state  $S \in \Omega$ ,  $(T(t)F)(S)$  is the expected value of  $F$  at time  $t$ , given that the initial distribution of  $\mathcal{M}$  is concentrated at  $S$  with probability one at time 0. In particular,  $\varphi_t = T(t)\varphi_0$  for all  $t \geq 0$ , and all initial distributions  $\varphi_0$ . In the case of the SEP we are considering, the infinitesimal generator  $\mathcal{L}$  of the semigroup  $T(t)$  is given by

$$\mathcal{L} = \sum_{\{i,j\} \in E} \lambda_{ij} (\tau_{ij} - 1).$$

For each  $\{i, j\} \in E$ , this replaces each  $S \in \Omega$  by  $S \circ \tau_{ij}$  at the rate  $\lambda_{ij}$ .

In preparation for Section 7.3, it is useful to regard  $\mathcal{L}$  as an element of the real semigroup algebra  $\mathfrak{A} = \mathbb{R}[\mathfrak{E}]$  of the semigroup  $\mathfrak{E}$  of all endofunctions  $f : \Omega \rightarrow \Omega$  (with the operation of functional composition). The left action of  $\mathfrak{E}$  on  $\Omega$  is extended to a left action of  $\mathfrak{A}$  on  $\mathbb{C}[\mathbf{x}]$  as usual: for  $f \in \mathfrak{E}$  and  $S \in \Omega$ ,  $f(\mathbf{x}^S) = \mathbf{x}^{f(S)}$ , extended bilinearly to all of  $\mathfrak{A}$  and  $\mathbb{C}[\mathbf{x}]$ . A permutation  $\sigma \in \mathfrak{S}(\Lambda)$  is identified with the endofunction  $f_\sigma : S \mapsto S \circ \sigma^{-1}$ , so this action of  $\mathfrak{A}$  agrees with the action of  $\mathfrak{S}(m)$  in Section 4.2. A left action of  $\mathfrak{A}$  on  $\mathbb{R}^\Omega$  is defined by  $Z(f(F)) = f(Z(F))$  for all  $f \in \mathfrak{E}$  and  $F \in \mathbb{R}^\Omega$ , and linear extension. More explicitly, for  $f \in \mathfrak{E}$ ,  $F \in \mathbb{R}^\Omega$ , and  $S \in \Omega$ ,

$$(f(F))(S) = F(f^{-1}(S)) = \sum \{F(S') : S' \in \Omega \text{ and } f(S') = S\}.$$

Consider an element of  $\mathfrak{A}$  of the form  $\mathcal{L} = \sum_{i=1}^N \lambda_i (f_i - 1)$  with all  $\lambda_i > 0$ . Let  $\lambda_i \leq L$  for all  $i \in [N]$ , and let  $K = \sum_{i=1}^N \lambda_i$ . The power series

$$\exp(t\mathcal{L}) = e^{-Kt} \sum_{n=0}^{\infty} \frac{t^n}{n!} \left[ \sum_{i=1}^N \lambda_i f_i \right]^n = \sum_{f \in \mathfrak{E}} P_f(t) \cdot f$$

in  $\mathfrak{A}[[t]]$  is such that for each  $f \in \mathfrak{E}$ ,  $P_f(t) \in \mathbb{R}[[t]]$  is dominated coefficientwise by  $\exp((LN - K)t)$ . Thus  $\exp(t\mathcal{L}) \in \mathfrak{A}[[t]]$  converges for all  $t \geq 0$ . The semigroup of transformations generated by  $\mathcal{L}$  is  $\exp(t\mathcal{L})$ .

To check claim (3) we will show that the semigroup  $T(t)$  of the SEP preserves stability for all  $t \geq 0$ : that is, if  $Z(\varphi_0)$  is stable then  $Z(\varphi_t) = T(t)Z(\varphi_0)$  is stable for all  $t \geq 0$ . This reduces to the case of a single pair  $\{i, j\} \in E$ , as follows. If  $\mathcal{M}_1$  and  $\mathcal{M}_2$  are Markov chains on the same finite state space, with semigroups  $T_1(t)$  and  $T_2(t)$  generated by  $\mathcal{L}_1$  and  $\mathcal{L}_2$ , then the semigroup generated by  $\mathcal{L}_1 + \mathcal{L}_2$  is

$$T(t) = \lim_{n \rightarrow \infty} [T_1(t/n)T_2(t/n)]^n,$$

by the Trotter product formula. By Hurwitz's Theorem, It follows that if  $T_i(t)$  preserves stability for all  $t \geq 0$  and  $i \in \{1, 2\}$ , then  $T(t)$  preserves stability for all  $t \geq 0$ . By repeated application of this argument, in order to show that the SEP semigroup  $T(t) = \exp(t\mathcal{L})$  preserves stability for all  $t \geq 0$  it is enough to show that for each  $\{i, j\} \in E$ ,  $T_{ij}(t) = \exp(t\lambda_{ij}(\tau_{ij} - 1))$  preserves stability for all  $t \geq 0$ . Now, since  $\tau_{ij}^2 = 1$ ,

$$T_{ij}(t) = \left( \frac{1 + e^{-2\lambda_{ij}t}}{2} \right) + \left( \frac{1 - e^{-2\lambda_{ij}t}}{2} \right) \cdot \tau_{ij}.$$

By Lemma 4.3, this preserves stability for all  $t \geq 0$ . This proves Theorem 7.1.  $\square$

**7.3. Further observations.** In verifying the hypotheses of the Feder-Mihail Theorem we used the fact that if  $f \in \mathfrak{S}_{\mathbb{R}}[\mathbf{x}]^{\text{MA}}$  is multiaffine and real stable, then  $\Delta_{ij}f(\mathbf{a}) \geq 0$  for all  $\{i, j\} \subseteq E$  and  $\mathbf{a} \in \mathbb{R}^m$ , by Theorem 3.1. In fact, we only needed the weaker hypothesis that  $\Delta_{ij}f(\mathbf{a}) \geq 0$  for all  $\{i, j\} \subseteq E$  and  $\mathbf{a} \in \mathbb{R}^m$  with  $\mathbf{a} > \mathbf{0}$ . A multiaffine real polynomial satisfying this weaker condition is a *Rayleigh* polynomial. (This terminology is by analogy with the Rayleigh monotonicity property of electrical networks – see Definition 2.5 of [7] and the references cited there. Multiaffine real stable polynomials are also called *strongly Rayleigh*.) The class of probability distributions  $\varphi$  such that  $Z(\varphi)$  is homogeneous, multiaffine, and Rayleigh meets all the conditions of the Feder-Mihail Theorem. It follows that all such distributions are NA.

Claim (2) above can be generalized in another way – the hypothesis of homogeneity can be removed, as follows. Let  $\mathbf{y} = (y_1, \dots, y_m)$  and let  $e_j(\mathbf{y})$  be the  $j$ -th elementary symmetric function of the  $\mathbf{y}$ . Given a multiaffine polynomial  $f = \sum_{S \subseteq [m]} c(S)\mathbf{x}^S$ , the *symmetric homogenization* of  $f$  is the polynomial  $f_{\text{sh}}(\mathbf{x}, \mathbf{y}) \in \mathbb{C}[\mathbf{x}, \mathbf{y}]^{\text{MA}}$  defined by

$$f_{\text{sh}}(\mathbf{x}, \mathbf{y}) = \sum_{S \subseteq [m]} c(S)\mathbf{x}^S \binom{m}{|S|}^{-1} e_{m-|S|}(\mathbf{y}).$$

Note that  $f_{\text{sh}}$  is homogeneous of degree  $m$ , and  $f_{\text{sh}}(\mathbf{x}, \mathbf{1}) = f(\mathbf{x})$ .

**Proposition 7.2** (Theorem 4.2 of [7]). *If  $f \in \mathfrak{S}_{\mathbb{R}}[\mathbf{x}]^{\text{MA}}$  is multiaffine and real stable then  $f_{\text{sh}} \in \mathfrak{S}_{\mathbb{R}}[\mathbf{x}, \mathbf{y}]^{\text{MA}}$  is homogeneous, multiaffine and real stable.*

(We omit the proof.)

**Corollary 7.3** (Theorem 4.9 of [7]). *Let  $\varphi : \Omega \rightarrow [0, \infty)$  be such that  $Z(\varphi)$  is nonzero, multiaffine, and real stable. Then for all  $\mathbf{a} \in \mathbb{R}^m$  with  $\mathbf{a} > \mathbf{0}$ ,  $\varphi^{\mathbf{a}}$  is NA.*

*Proof.* By Proposition 7.2,  $Z_{\text{sh}}(\varphi; \mathbf{x}, \mathbf{y})$  is nonzero, homogeneous, multiaffine, and real stable. This is the partition function for  $\psi : \{0, 1\}^{[2m]} \rightarrow [0, \infty)$  given by  $\psi(S) = \binom{m}{|S \cap [m]|}^{-1} \varphi(S \cap [m])$ . By claim (2) above,  $\psi^{\mathbf{a}}$  is NA for all  $\mathbf{a} \in \mathbb{R}^{2m}$  with  $\mathbf{a} > \mathbf{0}$ . By considering those  $\mathbf{a} \in \mathbb{R}^{2m}$  for which  $a_i = 1$  for all  $m+1 \leq i \leq 2m$ , it follows that  $\varphi^{\mathbf{a}}$  is NA for all  $\mathbf{a} \in \mathbb{R}^m$  with  $\mathbf{a} > \mathbf{0}$ .  $\square$

**Corollary 7.4** (Theorem 5.2 of [7]). *If the initial distribution  $\varphi_0$  of a SEP is such that  $Z(\varphi)$  is stable (but not necessarily homogeneous), then  $Z(\varphi_t)$  is stable, and hence  $\varphi_t$  is NA, for all  $t \geq 0$ .*

It is natural to try extending these results to asymmetric exclusion processes. For  $(i, j) \in \Lambda^2$  define  $\mathbf{t}_{ij} \in \mathfrak{E}$  by  $\mathbf{t}_{ij}(S) = S \circ \tau_{ij}$  if  $S(i) = 1$  and  $S(j) = 0$ , and  $\mathbf{t}_{ij}(S) = S$  otherwise, for all  $S \in \Omega$ . That is,  $\mathbf{t}_{ij}$  makes a particle hop from site  $i$  to site  $j$ , if possible. Let  $E$  be a set of ordered pairs in  $\Lambda^2$ , and for  $(i, j) \in E$  let  $\lambda_{ij} > 0$ . An asymmetric exclusion process is a Markov chain on  $\Omega$  with semigroup  $T(t) = \exp(t\mathcal{L})$  generated by something of the form

$$\mathcal{L} = \sum_{(i,j) \in E} \lambda_{ij}(\mathbf{t}_{ij} - 1).$$

By the argument for claim (3) above, in order to show that  $T(t)$  preserves stability for all  $t \geq 0$ , it suffices to do so for the two-site semigroup  $T_{\{1,2\}}(t) =$

$\exp(t\mathcal{L}_{\{1,2\}})$  generated by

$$\mathcal{L}_{\{1,2\}} = \lambda_{12}(\mathbf{t}_{12} - 1) + \lambda_{21}(\mathbf{t}_{21} - 1).$$

**Exercise 7.5** (Strengthening Remark 5.3 of [7]). With the notation above, let  $\lambda = \lambda_{12} + \lambda_{21}$ ,  $\beta_{12} = \lambda_{12}/\lambda$ , and  $\beta_{21} = \lambda_{21}/\lambda$ .

(a) In  $\mathfrak{A}$ ,  $\mathbf{t}_{12} + \mathbf{t}_{21} = 1 + \tau_{12}$ .

(b) If  $\omega$  is any word in  $\{\mathbf{t}_{12}, \mathbf{t}_{21}\}^n$ , then  $\mathbf{t}_{12}\omega = \mathbf{t}_{12}$  and  $\mathbf{t}_{21}\omega = \mathbf{t}_{21}$ .

(c) The semigroup generated by  $\mathcal{L}_{\{1,2\}}$  is

$$T_{\{1,2\}}(t) = e^{-\lambda t} + (1 - e^{-\lambda t})(\beta_{12}\mathbf{t}_{12} + \beta_{21}\mathbf{t}_{21})$$

(d) The semigroup  $T_{\{1,2\}}(t)$  preserves stability for all  $t \geq 0$  if and only if  $\beta_{12} = \beta_{21} = 1/2$ , in which case it reduces to the SEP (of rate  $\lambda/2$ ).

Thus, even the slightest asymmetry ruins preservation of stability by the SEP!

Finally, we consider a SEP in which particle number is not conserved. For  $i \in \Lambda$  define  $\mathbf{a}_i, \mathbf{a}_i^* \in \mathfrak{E}$  as follows: for  $S \in \Omega$  and  $j \in \Lambda$ , let  $(\mathbf{a}_i(S))(j) = (\mathbf{a}_i^*(S))(j) = S(j)$  if  $j \neq i$ , and  $(\mathbf{a}_i(S))(i) = 0$  and  $(\mathbf{a}_i^*(S))(i) = 1$ . That is,  $\mathbf{a}_i$  annihilates a particle at site  $i$ , and  $\mathbf{a}_i^*$  creates a particle at site  $i$ , if possible.

A SEP with particle creation and annihilation is a Markov chain on  $\Omega$  with semigroup  $T(t) = \exp(t\mathcal{L})$  generated by something of the form

$$\mathcal{L} = \sum_{\{i,j\} \in E} \lambda_{ij}(\tau_{ij} - 1) + \sum_{i \in \Lambda} [\theta_i(\mathbf{a}_i - 1) + \theta_i^*(\mathbf{a}_i^* - 1)],$$

in which the first sum is the generator of the SEP in Theorem 7.1 and  $\theta_i, \theta_i^* \geq 0$  for each  $i \in \Lambda$ .

By the argument for claim (3) above, to show that this  $T(t)$  preserves stability for all  $t \geq 0$ , it suffices to do so for the one-site semigroups generated by  $\mathcal{L}_1 = \theta(\mathbf{a}_1 - 1)$  and  $\mathcal{L}_1^* = \theta(\mathbf{a}_1^* - 1)$ , respectively.

**Exercise 7.6.** The semigroups generated by  $\mathcal{L}_1$  and  $\mathcal{L}_1^*$  are

$$T_1(t) = e^{-\theta t} + (1 - e^{-\theta t})\mathbf{a}_1 \quad \text{and} \quad T_1^*(t) = e^{-\theta t} + (1 - e^{-\theta t})\mathbf{a}_1^*,$$

respectively. Both  $T_1(t)$  and  $T_1^*(t)$  preserve stability.

**Corollary 7.7.** *If the initial distribution  $\varphi_0$  of a SEP with particle creation and annihilation is such that  $Z(\varphi)$  is stable, then  $Z(\varphi_t)$  is stable, and hence  $\varphi_t$  is NA, for all  $t \geq 0$ .*

## 8. INEQUALITIES FOR MIXED DISCRIMINANTS.

This section summarizes a powerful application of stable polynomials from a 2008 paper of Gurvits [11].

We will use without mention the facts that  $\log$  and  $\exp$  are strictly increasing functions on  $(0, \infty)$ . A function  $\rho : I \rightarrow \mathbb{R}$  defined on an interval  $I \subseteq \mathbb{R}$  is *convex* provided that for all  $a_1, a_2 \in I$ ,  $\rho((a_1 + a_2)/2) \leq (\rho(a_1) + \rho(a_2))/2$ . It is *strictly convex* if it is convex and equality holds here only when  $a_1 = a_2$ . A function  $\rho : I \rightarrow \mathbb{R}$  is *(strictly) concave* if  $-\rho$  is (strictly) convex. For example, for positive reals  $a_1, a_2 > 0$  one has  $(\sqrt{a_1} - \sqrt{a_2})^2 \geq 0$ , with equality only if  $a_1 = a_2$ . It follows that  $\log((a_1 + a_2)/2) \geq (\log(a_1) + \log(a_2))/2$ , with equality only if  $a_1 = a_2$ . That is,  $\log$  is strictly concave.

**Jensen's Inequality** (Theorem 90 of [12]). *Let  $\rho : I \rightarrow \mathbb{R}$  be defined on an interval  $I \subseteq \mathbb{R}$ , let  $a_i \in I$  for  $i \in [n]$ , and let  $b_i > 0$  for  $i \in [n]$  be such that  $\sum_{i=1}^n b_i = 1$ . If  $\rho$  is convex then*

$$\rho\left(\sum_{i=1}^n b_i a_i\right) \leq \sum_{i=1}^n b_i \rho(a_i).$$

*If  $\rho$  is strictly convex and equality holds, then  $a_1 = a_2 = \cdots = a_n$ .*

For integer  $d \geq 1$ , let  $G(d) = (1 - 1/d)^{d-1}$ , and let  $G(0) = 1$ . Note that  $G(1) = 0^0 = 1$ , and that  $G(d)$  is a strictly decreasing function for  $d \geq 1$ . For homogeneous  $f \in \mathbb{R}[x]$  with nonnegative coefficients, define the *capacity* of  $f$  to be

$$\text{cap}(f) = \inf_{\mathbf{c} > \mathbf{0}} \frac{f(\mathbf{c})}{c_1 \cdots c_m},$$

with the infimum over the set of all  $\mathbf{c} \in \mathbb{R}^m$  with  $c_i > 0$  for all  $i \in [m]$ .

**Lemma 8.1** (Lemma 3.2 of [11]). *Let  $f = \sum_{i=0}^d b_i x^i \in \mathbb{R}[x]$  be a nonzero univariate polynomial of degree  $d$  with nonnegative coefficients. If  $f$  is real stable then  $b_1 = f'(0) \geq G(d)\text{cap}(f)$ , and if  $\text{cap}(f) > 0$  then equality holds if and only if  $d \leq 1$  or  $f(x) = b_d(x + \xi)^d$  for some  $\xi > 0$ .*

*Proof.* If  $\text{cap}(f) = 0$  then there is nothing to prove, so assume that  $\text{cap}(f) > 0$ . If  $d = 0$  then  $f'(0) = b_1 = 0 = G(0)\text{cap}(f)$ , and if  $d = 1$  then  $f'(0) = b_1 = G(1)\text{cap}(f)$ , so assume that  $d \geq 2$ . If  $f(0) = 0$  then  $f'(0) = \lim_{c \rightarrow 0} f(c)/c \geq \text{cap}(f) > G(d)\text{cap}(f)$ . Thus, assume that  $d \geq 2$  and  $f(0) = b_0 > 0$ . We may rescale the polynomial so that  $b_0 = 1$ . Now there are  $a_i > 0$  for  $i \in [d]$  such that

$$f(x) = \prod_{i=1}^d (1 + a_i x),$$

and  $b_1 = a_1 + \cdots + a_d$ . For any  $c > 0$  we have

$$\frac{\log(\text{cap}(f)c)}{d} \leq \frac{\log(f(c))}{d} = \frac{1}{d} \sum_{i=1}^d \log(1 + a_i c) \leq \log\left(1 + \frac{b_1 c}{d}\right),$$

by Jensen's Inequality. It follows that  $\text{cap}(f)c \leq (1 + b_1 c/d)^d$  for all  $c > 0$ . Let  $g(x) = (1 + b_1 x/d)^d$ . Elementary calculus shows that

$$\text{cap}(g) = \inf_{c > 0} \frac{g(c)}{c} = \frac{g(c_*)}{c_*} = \frac{b_1}{G(d)}, \quad \text{in which } c_* = \frac{d}{b_1(d-1)}.$$

Since  $\text{cap}(f) \leq \text{cap}(g)$ , this yields the stated inequality. If equality holds, then equality holds in the application of Jensen's Inequality, and so  $f$  has the stated form.  $\square$

**Lemma 8.2** (Theorem 4.10 of [11]). *Let  $f \in \mathfrak{S}_{\mathbb{R}}[x_1, \dots, x_m]$  be real stable, with nonnegative coefficients, and homogeneous of degree  $m$ . Let  $g = \partial_m f|_{x_m=0}$ . Then*

$$\text{cap}(g) \geq G(\deg_m(d))\text{cap}(f).$$

*Proof.* We may assume that  $d = \deg_m(f) \geq 1$ . Let  $c_i > 0$  for  $i \in [m-1]$ , and let  $p_{\mathbf{c}}(x) = f(c_1, \dots, c_{m-1}, x)$ . Since  $f$  has nonnegative coefficients,  $p_{\mathbf{c}} \not\equiv 0$ . As in the proof of Lemma 2.4(f),  $p_{\mathbf{c}}$  has degree  $d$ . By specialization,  $p_{\mathbf{c}}$  is real stable. Lemma 8.1 implies that

$$g(\mathbf{c}) = p'_{\mathbf{c}}(0) \geq G(d)\text{cap}(p_{\mathbf{c}}) \geq G(d)\text{cap}(f)$$

for all  $\mathbf{c} \in \mathbb{R}^{m-1}$  with  $\mathbf{c} > \mathbf{0}$ . If  $m = 1$  then  $g = \text{cap}(g)$  is a constant. If  $m \geq 2$  then for any such  $\mathbf{c}$  let  $b = (c_1 \cdots c_{m-1})^{-1/(m-1)}$ . Since  $g$  is homogeneous of degree  $m - 1$ ,

$$\frac{g(\mathbf{c})}{c_1 \cdots c_{m-1}} = g(bc_1, \dots, bc_{m-1}) \geq G(d)\text{cap}(f).$$

It follows that  $\text{cap}(g) \geq G(d)\text{cap}(f)$ .  $\square$

**Theorem 8.3** (Theorem 2.4 of [11]). *Let  $f \in \mathfrak{S}_{\mathbb{R}}[x_1, \dots, x_m]$  be real stable, with nonnegative coefficients, and homogeneous of degree  $m$ . Let  $\deg_i(f) = d_i$  and  $e_i = \min\{i, d_i\}$  for each  $i \in [m]$ . Then*

$$\partial^{\mathbf{1}} f(\mathbf{0}) \geq \text{cap}(f) \prod_{i=2}^m G(e_i).$$

*Proof.* Let  $g_m = f$  and let  $g_{i-1} = \partial_i g_i|_{x_i=0}$  for all  $i \in [m]$ . By contraction and specialization,  $g_i$  is real stable for each  $i \in [m]$ . Notice that  $g_0 = \partial^{\mathbf{1}} f(\mathbf{0}) = \text{cap}(g_0)$ . By Lemma 8.2,  $\text{cap}(g_{i-1}) \geq \text{cap}(g_i) \cdot G(\deg_i g_i)$  for each  $i \in [m]$ . But  $\deg_i g_i \leq \deg_i f = d_i$ , and  $\deg_i g_i$  is at most the total degree of  $g_i$ , which is  $i$ . Hence  $\deg_i g_i \leq e_i$ , and thus  $G(\deg_i g_i) \geq G(e_i)$ . Thus  $\text{cap}(g_{i-1}) \geq \text{cap}(g_i) \cdot G(e_i)$  for each  $i \in [m]$ . Combining these inequalities (and  $G(e_1) = 1$ ) gives the result.  $\square$

With the notation of Theorem 8.3, since  $e_i \leq i$  for all  $i \in [m]$  and  $G(d)$  is a decreasing function of  $d$ , one has the inequality

$$\prod_{i=2}^m G(e_i) \geq \prod_{i=2}^m G(i) = \prod_{i=2}^m \left( \frac{i-1}{i} \right)^{i-1} = \frac{m!}{m^m}.$$

Thus, the following corollary is immediate.

**Corollary 8.4.** *Let  $f \in \mathfrak{S}_{\mathbb{R}}[x_1, \dots, x_m]$  be real stable, with nonnegative coefficients, and homogeneous of degree  $m$ . Then*

$$\partial^{\mathbf{1}} f(\mathbf{0}) \geq \frac{m!}{m^m} \cdot \text{cap}(f).$$

**Theorem 8.5** (Theorem 5.7 of [11]). *Let  $f \in \mathfrak{S}_{\mathbb{R}}[x_1, \dots, x_m]$  be real stable, with nonnegative coefficients, and homogeneous of degree  $m$ . Equality holds in the bound of Corollary 8.4 if and only if there are nonnegative reals  $a_i \geq 0$  for  $i \in [m]$  such that*

$$f(\mathbf{x}) = (a_1 x_1 + \cdots + a_m x_m)^m.$$

(We omit the proof.)

**Lemma 8.6** (Fact 2.2 of [11]). *Let  $f \in \mathbb{R}[x_1, \dots, x_m]$  be homogeneous of degree  $m$ , with nonnegative coefficients. Assume that  $\partial_i f(\mathbf{1}) = 1$  for all  $i \in [m]$ . Then  $\text{cap}(f) = 1$ .*

*Proof.* Let  $f = \sum_{\alpha} b(\alpha) \mathbf{x}^{\alpha}$ , so that if  $b(\alpha) \neq 0$  then  $|\alpha| = \sum_{i=1}^m \alpha(i) = m$ . By hypothesis, for all  $i \in [m]$ ,  $\sum_{\alpha} b(\alpha) \alpha(i) = 1$ . Averaging these over all  $i \in [m]$  yields  $f(\mathbf{1}) = \sum_{\alpha} b(\alpha) = 1$ , so that  $\text{cap}(f) \leq 1$ . Conversely, let  $\mathbf{c} \in \mathbb{R}^m$  with  $\mathbf{c} > \mathbf{0}$ .

Jensen's Inequality implies that

$$\begin{aligned} \log(f(\mathbf{c})) &= \log\left(\sum_{\alpha} b(\alpha)\mathbf{c}^{\alpha}\right) \\ &\geq \sum_{\alpha} b(\alpha)\log(\mathbf{c}^{\alpha}) = \sum_{i=1}^m \log(c_i) \sum_{\alpha} b(\alpha)\alpha(i) = \log(c_1 \cdots c_m). \end{aligned}$$

It follows that  $\text{cap}(f) \geq 1$ .  $\square$

**Example 8.7** (van der Waerden Conjecture). An  $m$ -by- $m$  matrix  $A = (a_{ij})$  is *doubly stochastic* if all entries are nonnegative reals and every row and column sums to one. In 1926, van der Waerden conjectured that if  $A$  is an  $m$ -by- $m$  doubly stochastic matrix then  $\text{per}(A) \geq m!/m^m$ , with equality if and only if  $A = (1/m)J$ , the  $m$ -by- $m$  matrix in which every entry is  $1/m$ . In 1981 this lower bound was proved by Falikman, and the characterization of equality was proved by Egorychev. These results follow immediately from Corollary 8.4 and Theorem 8.5, as follows. It suffices to prove the result for an  $m$ -by- $m$  doubly stochastic matrix  $A = (a_{ij})$  with no zero entries, by a routine limit argument. The polynomial

$$f_A(\mathbf{x}) = \prod_{j=1}^m (a_{1j}x_1 + \cdots + a_{mj}x_m)$$

is clearly homogeneous and real stable, with nonnegative coefficients and of degree  $m$ , and such that  $\deg_i(f_A) = m$  for all  $i \in [m]$ . Since  $A$  is doubly stochastic, Lemma 8.6 implies that  $\text{cap}(f_A) = 1$ . Since

$$\text{per}(A) = \partial^{\mathbf{1}} f_A(\mathbf{0}),$$

Corollary 8.4 and Theorem 8.5 imply the results of Falikman and Egorychev, respectively. Gurvits [11] also uses a similar argument to prove a refinement of the van der Waerden conjecture due to Schrijver and Valiant – see also [13].

Given  $n$ -by- $n$  matrices  $A_1, \dots, A_m$ , the *mixed discriminant* of  $\mathbf{A} = (A_1, \dots, A_m)$  is

$$\text{Disc}(\mathbf{A}) = \partial^{\mathbf{1}} \det(x_1 A_1 + \cdots + x_m A_m) \Big|_{\mathbf{x}=\mathbf{0}}.$$

This generalizes the permanent of an  $m$ -by- $m$  matrix  $B = (b_{ij})$  by considering the collection of matrices  $\mathbf{A}(B) = (A_1, \dots, A_m)$  defined by  $A_h = \text{diag}(a_{h1}, \dots, a_{hm})$  for each  $h \in [m]$ . In this case one sees that

$$\det(x_1 A_1 + \cdots + x_m A_m) = f_B(\mathbf{x})$$

with the notation of Example 8.7, and it follows that  $\text{Disc}(\mathbf{A}(B)) = \text{per}(B)$ .

**Example 8.8** (Bapat's Conjecture). Generalizing the van der Waerden conjecture, in 1989 Bapat considered the set  $\Omega(m)$  of  $m$ -tuples of  $m$ -by- $m$  matrices  $\mathbf{A} = (A_1, \dots, A_m)$  such that each  $A_i$  is positive semidefinite with trace  $\text{tr}(A_i) = 1$ , and  $\sum_{i=1}^m A_i = I$ . For any doubly stochastic matrix  $B$ ,  $\mathbf{A}(B)$  is in this set. The natural conjecture is that for all  $\mathbf{A} \in \Omega(m)$ ,  $\text{Disc}(\mathbf{A}) \geq m!/m^m$ , and equality is attained if and only if  $\mathbf{A} = \mathbf{A}((1/m)J)$ . This was proved by Gurvits in 2006 – again, it follows directly from Corollary 8.4 and Theorem 8.5. It suffices to prove the result for  $\mathbf{A} \in \Omega(m)$  such that each  $A_i$  is positive definite, by a routine limit argument. By Proposition 2.1, for  $\mathbf{A} \in \Omega(m)$ , the polynomial

$$f_{\mathbf{A}}(\mathbf{x}) = \det(x_1 A_1 + \cdots + x_m A_m)$$

is real stable. Since each  $A_i$  is positive definite, all coefficients of  $f_{\mathbf{A}}$  are nonnegative,  $f_{\mathbf{A}}$  is homogeneous of degree  $m$ , and  $\deg_i(f_{\mathbf{A}}) = m$  for all  $i \in [m]$ . Since  $\mathbf{A} \in \Omega(m)$ , Lemma 8.6 implies that  $\text{cap}(f_{\mathbf{A}}) = 1$ . Thus,  $f_{\mathbf{A}}$  satisfies the hypothesis of Theorems 8.3 and 8.5, and since  $\text{Disc}(\mathbf{A}) = \partial^1 f_{\mathbf{A}}(\mathbf{0})$ , the result follows.

## 9. FURTHER DIRECTIONS.

**9.1. Other circular regions.** Let  $\Omega \subseteq \mathbb{C}^m$ . A polynomial  $f \in \mathbb{C}[\mathbf{x}]$  is  $\Omega$ -stable if either  $f \equiv 0$  identically, or  $f(\mathbf{z}) \neq 0$  for all  $\mathbf{z} \in \Omega$ . At this level of generality little can be said. If  $\Omega = \mathcal{A}_1 \times \cdots \times \mathcal{A}_m$  is a product of open circular regions then there are Möbius transformations  $z \mapsto \phi_i(z) = (a_i z + b_i)/(c_i z + d_i)$  such that  $\phi_i(\mathcal{H}) = \mathcal{A}_i$  for all  $i \in [m]$ . The argument in Section 4.1 shows that  $f \in \mathbb{C}[\mathbf{x}]$  is  $\Omega$ -stable if and only if

$$\tilde{f} = (\mathbf{c}\mathbf{z} + \mathbf{d})^{\deg f} \cdot f(\phi_1(z_1), \dots, \phi_m(z_m))$$

is stable. In this way results about stable polynomials can be translated into results about  $\Omega$ -stable polynomials for any  $\Omega$  that is a product of open circular regions.

Theorem 6.3 of [5] is the  $\Omega$ -stability analogue of Theorem 5.2. We mention only two consequences of this. Let  $\mathcal{D} = \{z \in \mathbb{C} : |z| < 1\}$  be the open unit disc, and for  $\theta \in \mathbb{R}$  let  $\mathcal{H}_\theta = \{e^{-i\theta} z : z \in \mathcal{H}\}$ . Thus  $\mathcal{H}_0 = \mathcal{H}$ , and  $\mathcal{H}_{\pi/2}$  is the open right half-plane. A  $\mathcal{D}^m$ -stable polynomial is called *Schur stable*, and a  $\mathcal{H}_{\pi/2}^m$ -stable polynomial is called *Hurwitz stable*.

**Proposition 9.1** (Remark 6.1 of [5]). *Fix  $\kappa \in \mathbb{N}^m$ , and let  $T : \mathbb{C}[\mathbf{x}]^{\leq \kappa} \rightarrow \mathbb{C}[\mathbf{x}]$  be a linear transformation. The following are equivalent:*

- (a)  *$T$  preserves Schur stability.*
- (b)  *$T((\mathbf{1} + \mathbf{xy})^\kappa)$  is Schur stable in  $\mathbb{C}[\mathbf{x}, \mathbf{y}]$ .*

**Proposition 9.2** (Remark 6.1 of [5]). *Fix  $\kappa \in \mathbb{N}^m$ , and let  $T : \mathbb{C}[\mathbf{x}]^{\leq \kappa} \rightarrow \mathbb{C}[\mathbf{x}]$  be a linear transformation. The following are equivalent:*

- (a)  *$T$  preserves Hurwitz stability.*
- (b)  *$T((\mathbf{1} + \mathbf{xy})^\kappa)$  is Hurwitz stable in  $\mathbb{C}[\mathbf{x}, \mathbf{y}]$ .*

**9.2. Applications of Theorem 5.4.** It is natural to consider a multivariate analogue of the multiplier sequences studied by Pólya and Schur. Let  $\lambda : \mathbb{N}^m \rightarrow \mathbb{R}$ , and define a linear transformation  $T_\lambda : \mathbb{C}[\mathbf{x}] \rightarrow \mathbb{C}[\mathbf{x}]$  by  $T_\lambda(\mathbf{x}^\alpha) = \lambda(\alpha)\mathbf{x}^\alpha$  for all  $\alpha \in \mathbb{N}^m$ , and linear extension. For which  $\lambda$  does  $T_\lambda$  preserve real stability? The answer: just the ones you get from the Pólya-Schur Theorem, and no more.

**Theorem 9.3** (Theorem 1.8 of [4]). *Let  $\lambda : \mathbb{N}^m \rightarrow \mathbb{R}$ . Then  $T_\lambda$  preserves real stability if and only if there are univariate multiplier sequences  $\lambda_i : \mathbb{N} \rightarrow \mathbb{R}$  for  $i \in [m]$  and  $\epsilon \in \{-1, +1\}$  such that*

$$\lambda(\alpha) = \lambda_1(\alpha(1)) \cdots \lambda_m(\alpha(m))$$

*for all  $\alpha \in \mathbb{N}^m$ , and either  $\epsilon^{|\alpha|}\lambda(\alpha) \geq 0$  for all  $\alpha \in \mathbb{N}^m$ , or  $\epsilon^{|\alpha|}\lambda(\alpha) \leq 0$  for all  $\alpha \in \mathbb{N}^m$ .*

Theorem 5.4 (and similarly Propositions 9.1 and 9.2) can be used to derive a wide variety of results of the form: such-and-such an operation preserves stability (or Schur or Hurwitz stability). Here is a short account of Hinkkanen's proof of the Lee-Yang Circle Theorem, taken from Section 8 of [6].



For  $f, g \in \mathbb{C}[\mathbf{x}]^{\text{MA}}$ , say  $f = \sum_{S \subseteq [m]} a(S) \mathbf{x}^S$  and  $g = \sum_{S \subseteq [m]} b(S) \mathbf{x}^S$ , let

$$f \bullet g = \sum_{S \subseteq [m]} a(S) b(S) \mathbf{x}^S$$

be the *Schur-Hadamard product* of  $f$  and  $g$ .

**Theorem 9.4** (Hinkkanen, Theorem 8.5 of [6]). *If  $f, g \in \mathbb{C}[\mathbf{x}]^{\text{MA}}$  are Schur stable then  $f \bullet g$  is Schur stable.*

*Proof.* Let  $T_g : \mathbb{C}[\mathbf{x}]^{\text{MA}} \rightarrow \mathbb{C}[\mathbf{x}]^{\text{MA}}$  be defined by  $f \mapsto f \bullet g$ . By Proposition 9.1, to show that  $T_g$  preserves Schur stability it suffices to show that  $T_g((\mathbf{1} + \mathbf{xy})^{[m]})$  is Schur stable. Clearly  $T_g((\mathbf{1} + \mathbf{xy})^{[m]}) = g(x_1 y_1, \dots, x_m y_m)$  is Schur stable since  $g(\mathbf{x})$  is. Hence  $T_g$  preserves Schur stability, and so  $f \bullet g$  is Schur stable.  $\square$

**Theorem 9.5** (Lee-Yang Circle Theorem, Theorem 8.4 of [6]). *Let  $A = (a_{ij})$  be a Hermitian  $m$ -by- $m$  matrix with  $|a_{ij}| \leq 1$  for all  $i, j \in [m]$ . Then the polynomial*

$$f(\mathbf{x}) = \sum_{S \subseteq [m]} \mathbf{x}^S \prod_{i \in S} \prod_{j \notin S} a_{ij}$$

*is Schur stable. The diagonalization  $g(x) = f(x, \dots, x)$  is such that  $x^m g(1/x) = g(x)$ , and it follows that all roots of  $g(x)$  are on the unit circle.*

*Proof.* For  $i < j$  in  $[m]$  let

$$f_{ij} = (1 + a_{ij} x_i + \overline{a_{ij}} x_j + x_i x_j) \prod_{h \in [m] \setminus \{i, j\}} (1 + x_h).$$

One can check that each  $f_{ij}$  is Schur stable. The polynomial  $f(\mathbf{x})$  is the Schur-Hadamard product of all the  $f_{ij}$  for  $\{i, j\} \subseteq [m]$ . By Theorem 9.4,  $f(\mathbf{x})$  is Schur stable.  $\square$

Section 8 of [6] contains many many more results of this nature.

**9.3. A converse to the Grace-Walsh-Szegő Theorem.** The argument of Sections 4.2 and 4.3 can be used to prove the following.

**Exercise 9.6.** If  $f \in \mathfrak{S}[\mathbf{x}]^{\text{MA}}$  is multiaffine and stable then

$$T_{\mathfrak{S}(m)}(f) = \frac{1}{m!} \sum_{\sigma \in \mathfrak{S}(m)} \sigma(f)$$

is multiaffine and stable.

This is in fact equivalent to the GWS Theorem, since for all  $f \in \mathbb{C}[\mathbf{x}]^{\text{MA}}$ ,  $T_{\mathfrak{S}(m)} f(\mathbf{x}) = \text{Pol}_m f(x, \dots, x)$ . For which transitive permutation groups  $G \leq \mathfrak{S}(m)$  does the linear transformation  $T_G = |G|^{-1} \sum_{\sigma \in G} \sigma$  preserve stability? The answer: not many, and they give nothing new.

**Theorem 9.7** (Theorem 6 of [9]). *Let  $G \leq \mathfrak{S}(m)$  be a transitive permutation group such that  $T_G$  preserves stability. Then  $T_G = T_{\mathfrak{S}(m)}$ .*

**9.4. Phase and support theorems.** A polynomial  $f \in \mathbb{C}[\mathbf{x}]$  has *definite parity* if every monomial  $\mathbf{x}^\alpha$  occurring in  $f$  has total degree of the same parity: all are even, or all are odd.

**Theorem 9.8** (Theorem 6.2 of [10]). *Let  $f \in \mathbb{C}[\mathbf{x}]$  be Hurwitz stable and with definite parity. Then there is a phase  $0 \leq \theta < 2\pi$  such that  $e^{-i\theta} f(\mathbf{x})$  has only real nonnegative coefficients.*

The *support* of  $f = \sum_{\alpha} c(\alpha) \mathbf{x}^\alpha$  is  $\text{supp}(f) = \{\alpha \in \mathbb{N}^m : c(\alpha) \neq 0\}$ . Let  $\delta_i$  denote the unit vector with a one in the  $i$ -th coordinate, and for  $\alpha \in \mathbb{Z}^n$  let  $|\alpha| = \sum_{i=1}^m |\alpha(i)|$ . A *jump system* is a subset  $\mathcal{J} \subseteq \mathbb{Z}^m$  satisfying the following *two-step axiom*:

**(J)** If  $\alpha, \beta \in \mathcal{J}$  and  $i \in [m]$  and  $\epsilon \in \{-1, +1\}$  are such that  $\alpha' = \alpha + \epsilon \delta_i$  satisfies  $|\alpha' - \beta| < |\alpha - \beta|$ , then either  $\alpha' \in \mathcal{J}$  or there exists  $j \in [m]$  and  $\varepsilon \in \{-1, +1\}$  such that  $\alpha'' = \alpha' + \varepsilon \delta_j \in \mathcal{J}$  and  $|\alpha'' - \beta| < |\alpha' - \beta|$ .

Jump systems generalize some more familiar combinatorial objects. A jump system contained in  $\{0, 1\}^m$  is a *delta-matroid*. A delta-matroid  $\mathcal{J}$  for which  $|\alpha|$  is constant for all  $\alpha \in \mathcal{J}$  is the set of bases of a *matroid*. For bases of matroids, the two-step axiom (J) reduces to the basis exchange axiom familiar from linear algebra: if  $A, B \in \mathcal{J}$  and  $a \in A \setminus B$ , then there exists  $b \in B \setminus A$  such that  $(A \setminus \{a\}) \cup \{b\}$  is in  $\mathcal{J}$ .

**Theorem 9.9** (Theorem 3.2 of [8]). *If  $f \in \mathfrak{S}[\mathbf{x}]$  is stable then the support  $\text{supp}(f)$  is a jump system.*

Recall from Section 7 that for multiaffine polynomials with nonnegative coefficients, real stability implies the Rayleigh property. A set system  $\mathcal{J}$  is *convex* when  $A, B \in \mathcal{J}$  and  $A \subseteq B$  imply that  $C \in \mathcal{J}$  for all  $A \subseteq C \subseteq B$ .

**Theorem 9.10** (Section 4 of [15]). *Let  $f = \sum_{S \subseteq [m]} c(S) \mathbf{x}^S$  be multiaffine with real nonnegative coefficients, and assume that  $f$  is Rayleigh.*

- (a) *The support  $\text{supp}(f)$  is a convex delta-matroid.*
- (b) *The coefficients are log-submodular: for all  $A, B \subseteq [m]$ ,*

$$c(A \cap B)c(A \cup B) \leq c(A)c(B).$$

#### REFERENCES

1. J. Borcea and P. Brändén, *Applications of stable polynomials to mixed determinants: Johnson's conjectures, unimodality, and symmetrized Fischer products*, Duke Math. J. **143** (2008), 205–223.
2. J. Borcea and P. Brändén, *Lee–Yang problems and the geometry of multivariate polynomials*, Lett. Math. Phys. **86** (2008), 53–61.
3. J. Borcea and P. Brändén, *Pólya–Schur master theorems for circular domains and their boundaries*, Ann. of Math. **170** (2009), 465–492.
4. J. Borcea and P. Brändén, *Multivariate Pólya–Schur classification problems in the Weyl algebra*, to appear in Proc. London Math. Soc.
5. J. Borcea and P. Brändén, *The Lee–Yang and Pólya–Schur programs I: linear operators preserving stability*, Invent. Math. **177** (2009), 541–569.
6. J. Borcea and P. Brändén, *The Lee–Yang and Pólya–Schur programs II: theory of stable polynomials and applications* Comm. Pure Appl. Math. **62** (2009), 1595–1631.
7. J. Borcea, P. Brändén, and T.M. Liggett, *Negative dependence and the geometry of polynomials*, J. Amer. Math. Soc. **22** (2009), 521–567.
8. P. Brändén, *Polynomials with the half-plane property and matroid theory*, Adv. Math. **216** (2007), 302–320.

9. P. Brändén and D.G. Wagner, *A converse to the Grace–Walsh–Szegő theorem*, Math. Proc. Camb. Phil. Soc. **147** (2009), 447–453.
10. Y.-B. Choe, J.G. Oxley, A.D. Sokal, and D.G. Wagner, *Homogeneous polynomials with the half-plane property*, Adv. in Appl. Math. **32** (2004), 88–187.
11. L. Gurvits, *Van der Waerden/Schrijver–Valiant like conjectures and stable (aka hyperbolic) homogeneous polynomials: one theorem for all. With a corrigendum*, Electron. J. Combin. **15** (2008), R66 (26 pp).
12. G. Hardy, J.E. Littlewood, and G. Pólya, “Inequalities (Second Edition),” Cambridge U.P., Cambridge UK, 1952.
13. M. Laurent and A. Schrijver, *On Leonid Gurvits’ proof for permanents*, <http://homepages.cwi.nl/~lex/files/perma5.pdf>
14. Q.I. Rahman and G. Schmeisser, “Analytic Theory of Polynomials,” London Math. Soc. Monographs (N.S.) **26**, Oxford U.P., New York NY, 2002.
15. D.G. Wagner, *Negatively correlated random variables and Mason’s conjecture for independent sets in matroids*, Ann. of Combin. **12** (2008), 211–239.
16. D.G. Wagner and Y. Wei, *A criterion for the half-plane property*, Discrete Math. **309** (2009), 1385–1390.

DEPARTMENT OF COMBINATORICS AND OPTIMIZATION, UNIVERSITY OF WATERLOO, WATERLOO, ONTARIO, CANADA N2L 3G1

*E-mail address:* `dgwagner@math.uwaterloo.ca`

# THE CONFORMAL GEOMETRY OF BILLIARDS

LAURA DE MARCO

## 1. INTRODUCTION

In this note, we examine the dynamics of billiards on polygonal tables. This is intended to be neither new research nor a survey, but rather a snapshot of recent work in one corner of the billiard-dynamics arena. We will concentrate on billiard tables where all interior angles are rational multiples of  $\pi$ . This class of billiard tables is closely related to the study of *translation surfaces*, Riemann surfaces  $X$  equipped with a holomorphic 1-form  $\omega$ , thus endowing  $X$  with a flat Euclidean metric structure away from finitely many cone-type singularities. Many recent results about billiard tables of this type come from general statements about moduli spaces of translation surfaces. The theme of this note is the search for *dynamically optimal* billiard tables: tables on which any billiard trajectory (which avoids the corners) is either periodic or it covers the table uniformly. Figure 1.1 shows an example of a dynamically optimal table. Careful definitions and examples are given in later sections; the following is an overview of the presentation.

A polygonal billiard table  $T$ , with all angles equal to rational multiples of  $\pi$ , gives rise to

- a translation surface  $(X_T, \omega_T)$  with genus  $g(X_T) \geq 1$ , via a process called *unfolding*; and
- a discrete subgroup  $\Gamma_T \subset \mathrm{SL}_2\mathbb{R}$ , the stabilizer of  $(X_T, \omega_T)$  under a *stretching* operation.

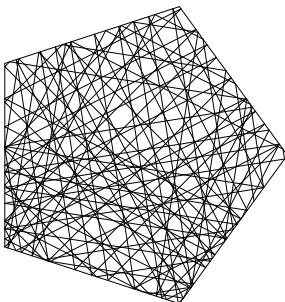


FIGURE 1.1. A billiard trajectory on the regular pentagon.

The genus  $g(X_T)$  is easily computable from the table  $T$ ; the formula is stated below in equation (3.1). Every table with genus  $g(X_T) = 1$  has optimal dynamics as we discuss in §4.2. In genus 2, by studying the orbits of an  $\mathrm{SL}_2\mathbb{R}$ -action on the moduli space of translation surfaces, McMullen showed:

**Theorem 1.1.** [Mc2] *For  $g(X_T) = 2$ , a table  $T$  has optimal dynamics if and only if  $\Gamma_T$  is a lattice in  $\mathrm{SL}_2\mathbb{R}$ .*

Billiard tables with  $\Gamma_T$  a lattice are in fact quite rare. McMullen established a complete list of tables  $T$  with optimal dynamics in genus 2 [Mc3]; see Theorem 4.3.

It was first observed by Veech that for every genus  $g(X_T)$ , if  $\Gamma_T$  is a lattice in  $\mathrm{SL}_2\mathbb{R}$ , then the table  $T$  has optimal dynamics [Ve]. In fact, his statement was much stronger: the geodesic flow on a translation surface  $(X, \omega)$  (which is not necessarily the translation surface for a billiard table) is dynamically optimal whenever  $\Gamma$  is a lattice, where  $\Gamma = \Gamma(X, \omega)$  is the so-called *Veech group* of the surface. Theorem 1.1 is itself a consequence of a more general statement about the geodesic flow on translation surfaces of genus 2; see §4. It is reasonable to guess that the equivalence of Theorem 1.1 holds for translation surfaces in every genus. However, Smillie and Weiss have shown recently:

**Theorem 1.2.** [SW] *There exist translation surfaces  $(X, \omega)$  which have optimal flow dynamics but for which the Veech group  $\Gamma$  is not a lattice.*

Theorem 1.2 leaves open the existence of billiard tables with optimal dynamics but non-lattice Veech group; billiard surfaces  $(X_T, \omega_T)$  and their  $\mathrm{SL}_2\mathbb{R}$ -orbits form only a small (measure 0) subset of the moduli space of translation surfaces in any genus  $> 1$ .

The Smillie-Weiss examples rely on a covering construction of Hubert and Schmidt [HS]: there exist surfaces  $(X, \omega)$  with lattice Veech group and holomorphic branched coverings of finite degree

$$f : Y \rightarrow X$$

so that the Veech group of the translation surface  $(Y, f^*\omega)$  is not a lattice nor even finitely generated. In certain cases, the dynamical properties of the geodesic flow on the surface  $(X, \omega)$  are preserved when passing to the branched cover.

The Hubert-Schmidt construction has led to further collections of interesting examples. We conclude this note with a discussion of the following recent result of Cheung, Hubert, and Masur:

**Theorem 1.3.** [CHM] *The billiard dynamics on the isosceles triangle with angles  $(2\pi/5, 3\pi/10, 3\pi/10)$  satisfy a topological dichotomy but are non-optimal: for each direction, either all billiard trajectories are closed or all are dense, but there exist directions in which billiard trajectories are dense but not uniformly distributed.*

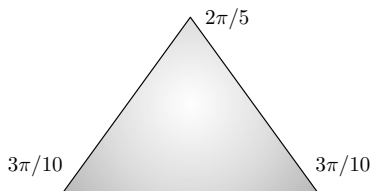


FIGURE 1.2. A triangular table with non-optimal dynamics but satisfying the topological dichotomy.

In standard dynamical language, the billiard flow in each direction is either completely periodic or minimal, while certain minimal directions are not uniquely ergodic. The triangle of Theorem 1.3 is the only known billiard table with this type of dynamics.

**Acknowledgements.** I am greatly indebted to the experts in translation surfaces for helping me prepare this note. Special thanks go to Yitwah Cheung, Howard Masur, and Curt McMullen for the conversations and numerous email exchanges about their work. Also, Curt McMullen generated the images of Figures 1.1 and 2.2. I would like to thank Jayadev Athreya, Matthew Bainbridge, Alex Eskin, and John Smillie for helpful and lengthy discussions about other recent results in this area. This note does not do justice to their beautiful mathematics. I thank the AMS for giving me this opportunity to learn about billiards; my research is also supported by the National Science Foundation and the Sloan Foundation.

## 2. BILLIARD TABLES

For this article, a *billiard table* means a polygon in  $\mathbb{R}^2$  with all angles a rational multiple of  $\pi$ . See Figure 2.1. A *billiard trajectory* in direction  $\theta$  is a straight-line path which begins at some point in the interior of the table, at angle  $\theta$  as measured from the positive real axis, and bounces off the edges with angle of reflection equal to the angle of incidence. If a billiard trajectory hits a vertex of the polygon, it stops. As the angles are rational multiples of  $\pi$ , the billiard path will travel again in direction  $\theta$  after finitely many reflections off the sides of the table.

**2.1. The square table.** The simplest example is the square table of side length 1, with sides parallel to the coordinate axes. In this table, it is easy to see that any billiard trajectory of angle  $\theta = p\pi/q$  for integers  $p$  and  $q$  will either hit a vertex or eventually return to its original configuration (position and angle). If the trajectory encounters a vertex, then it must encounter a vertex also in backward time (traveling in the opposite direction). On the other hand, for all other angles and any initial point, if the trajectory encounters a vertex then it will never hit a vertex in backwards time; all infinite trajectories bounce around the table spending equal time in parts with equal area.

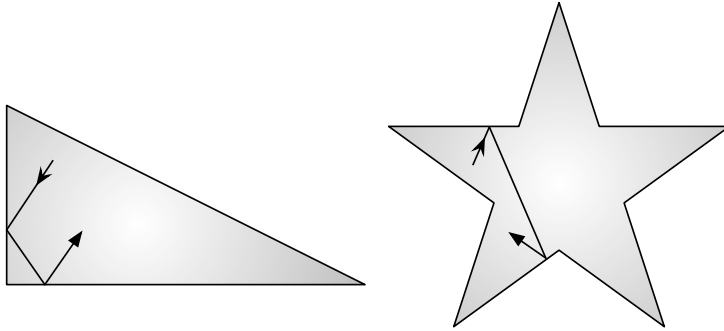


FIGURE 2.1. Polygonal billiard tables and trajectories.

**2.2. Optimal billiard dynamics.** We say a billiard table has *optimal dynamics* if for each direction  $\theta$ , one of the following holds:

- (1) every trajectory is either periodic or encounters a vertex in both forward and backward time; or
- (2) every trajectory is infinite in either forward or backward time (or both), and every infinite trajectory is uniformly distributed.

We must take care in our meaning of uniform distribution. Because the table has rational angles, every trajectory points in only finitely many different directions under reflections off the sides of the tables. For each direction  $\theta$ , we may take multiple copies of the polygonal table, one for each direction arising by reflection of trajectories in direction  $\theta$ . We say a trajectory is uniformly distributed if it equidistributes with respect to Lebesgue measure on this union of tables. In other words, for any trajectory of infinite length in direction  $\theta$ , let  $\gamma(t)$ ,  $t \geq 0$ , be a parametrization of this trajectory with unit speed (so with each reflection in a side,  $\gamma(t)$  jumps to another copy of the table). For each time  $s > 0$ , we can define probability measure on the union of tables by

$$\frac{1}{s} \gamma_* m_s$$

where  $m_s$  is arc-length measure on the interval  $[0, s]$  in  $\mathbb{R}$ . Uniform distribution means that this family of measures converges weakly as  $s \rightarrow \infty$  to normalized area measure on the finite union of polygonal tables. The property of optimal dynamics is also called *Veech dichotomy* in the literature.

**2.3. Examples and non-examples.** As explained in §2.1, the unit square table has optimal dynamics. In fact, any polygon which is tiled by a square (so that all vertices of the polygon coincide with vertices of the square tiles) will also have optimal dynamics [GJ]. A distinctly different class of examples was studied by Veech, who showed:

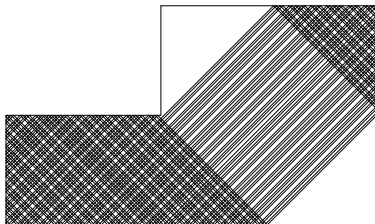


FIGURE 2.2. A billiard trajectory on an  $L$ -shaped table, neither closed nor dense.

**Theorem 2.1.** [Ve] *For every  $n \geq 3$ , the regular  $n$ -gon is a dynamically optimal billiard table.*

On the other hand, it is easy to construct tables with non-optimal dynamics. For example, begin with a square table and attach a rectangle to one side with side lengths  $a$  and  $b$  where  $a/b \notin \mathbb{Q}$ . See Figure 2.2. Taking the direction  $\theta = \pi/4$ , we see that the trajectories in direction  $\theta$  which enter the smaller rectangle are neither closed nor dense.

**2.4. Topologically optimal tables.** There is a notion called topological dichotomy for billiard tables which is weaker than optimal dynamics. A billiard table satisfies the topological dichotomy if for each direction  $\theta$ ,

- (1) every trajectory is either periodic or encounters a vertex in both forward and backward time; or
- (2) every trajectory is infinite in either forward or backward time (or both), and every infinite trajectory is dense.

As with uniform distribution, we require that the trajectory be dense on the finite union of tables corresponding to different directions under reflection.

It is a non-trivial task to find billiard table examples which have dense but non-uniformly distributed trajectories. The following examples were studied by Masur and Smillie, following a construction of Veech; see [MT], [MS]. Consider a rectangular table with barrier as in Figure 2.3: begin with a rectangular table of side lengths 1 and 2 and build a perpendicular wall in the middle of the long side of length  $\ell < 1$ . When  $\ell$  is rational, the table is dynamically optimal. When  $\ell$  is Diophantine (so it is not too closely approximated by rationals), the table is neither dynamically optimal nor topologically optimal, but it has billiard trajectories which are dense and non-uniformly distributed. In this case, the set of directions  $\theta \in [0, 2\pi)$  with dense but non-uniformly distributed trajectories is as large as possible, having Hausdorff dimension  $1/2$  [Ch].



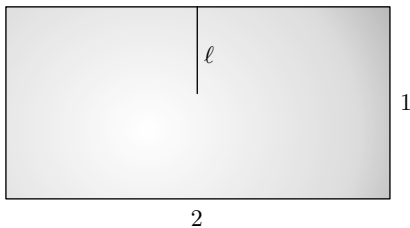


FIGURE 2.3. A rectangular billiard table with barrier of length  $\ell$ .

### 3. THE TRANSLATION SURFACE OF A BILLIARD TABLE

In this section we describe the process of *unfolding*, passing from a polygonal billiard table to a Riemann surface equipped with a holomorphic 1-form. In this way, a billiard trajectory which bounces off the walls of the table *unfolds* into a straight line on the surface.

**3.1. Unfolding a billiard table.** Fix a polygon  $T$  in  $\mathbb{R}^2$  and assume that all of its angles are rational multiples of  $\pi$ . Let  $G \subset O_2(\mathbb{R})$  be the group generated by reflections in the sides of  $T$ . Because of the rational angles, the group  $G$  is finite; let  $N = |G|$ . If the interior angles of  $T$  are expressed as  $m_i\pi/n_i$ , where the integers  $m_i$  and  $n_i$  have no common factors, then the number  $N$  is equal to twice the least common multiple of the  $n_i$ .

Take  $N$  copies of  $T$ , one for each reflected image  $gT$  with  $g \in G$ . Glue edges of distinct copies according to the reflection rules: if  $h \in G$  is represented by reflection across an edge  $e$  of  $gT$ , then  $e$  is glued to its image in  $hgT$ . The genus of the resulting surface  $X_T$  is given by the formula

$$(3.1) \quad g(X_T) = 1 + \frac{N}{4} \left( k - 2 - \sum_{i=1}^k \frac{1}{n_i} \right)$$

where  $k$  is the number of vertices of  $T$ ; see [MT].

It is easy to see that the unit square unfolds into a torus, as does an equilateral triangle. For the regular pentagon depicted in Figure 1.1, the reflection group has 10 elements, and the table unfolds into a surface of genus 6. On the other hand, the  $(2\pi/5, 3\pi/10, 3\pi/10)$  triangle shown in Figure 1.2 tiles the regular pentagon, but it unfolds into a surface of genus 4.

The Euclidean coordinates on the polygon  $T$  induce a flat conformal structure on the resulting surface, together with a finite collection of cone-point singularities (where the total angle at a point exceeds  $2\pi$ ). This structure can be recorded by the holomorphic 1-form  $dz$  on  $T$ , glued up to define a 1-form  $\omega_T$  on the unfolded surface

$X_T$ . The cone points are simply the zeroes of  $\omega_T$ . The pair  $(X_T, \omega_T)$  defines the *translation surface* associated to the table  $T$ .

In fact, every compact Riemann surface  $X$  equipped with a holomorphic 1-form can be obtained by gluing polygons in this way, though the polygons will not generally be reflections of a single polygonal shape  $T$ ; see the discussion in [Ma2].

**3.2. Geodesic flow on a translation surface.** The notions of trajectories and optimal dynamics can be defined on general translation surfaces  $(X, \omega)$ . Indeed, the 1-form  $\omega$  gives a natural way to choose local Euclidean coordinates on  $X$  away from the zeroes of  $\omega$ . Namely, for any point  $z_0 \in X$  with  $\omega_{z_0} \neq 0$ , we can integrate  $\omega$  to define a coordinate chart near  $z_0$  by

$$\varphi(z) = \int_{z_0}^z \omega$$

which is locally invertible and locally independent of path. In fact, the transition functions for these coordinate charts are given by translations, which explains the term “translation” surface. The Euclidean charts induce a flat metric on the surface, away from the zeroes of  $\omega$ , and the geodesics in this metric are simply the straight lines in these coordinates. The charts glue up at the zeroes of  $\omega$  to form the cone-like singularities.

When a surface comes from unfolding a billiard table, the straight lines are precisely the unfolded billiard trajectories. Thus, we can discuss the geodesics on a general translation surface to make conclusions about billiard trajectories. We say a translation surface  $(X, \omega)$  has *optimal flow dynamics* if its geodesics in each direction satisfy the dichotomy of §2.2. We say the surface satisfies the *topological dichotomy* if its geodesics in each direction satisfy the dichotomy of §2.4.

**3.3. Most geodesics are uniformly distributed.** The idea of unfolding seems to have first appeared in [KZ], where the authors studied topological transitivity of the geodesic flow on the associated translation surface, concluding that most directions on a table are minimal (all orbits are dense). General results about differentials on closed surfaces imply that the periodic directions (where all geodesics are either closed or travel between zeroes of  $\omega$ ) are dense in the circle; see [MT]. On the other hand, almost every direction gives rise to uniformly distributed geodesics [KMS]. Back on the table  $T$ , we get uniform distribution of the billiard trajectories in those directions.

## 4. STRETCHING THE BILLIARD TABLES

In this section, we describe the stretching deformation of billiard tables and translation surfaces, and we define the Veech group associated to a table. We conclude the section with McMullen’s classification of dynamically optimal billiard tables which unfold into surfaces of genus 2.

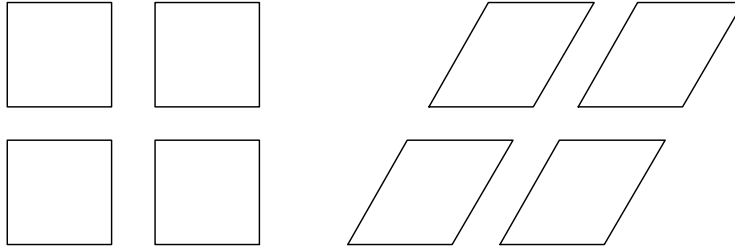


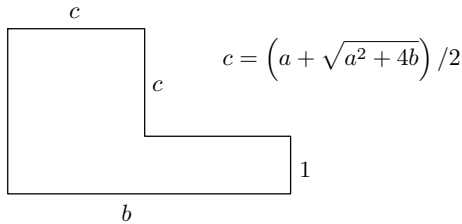
FIGURE 4.1. The square table with its reflected copies, and their images under shearing.

**4.1. An action of  $\mathrm{SL}_2\mathbb{R}$ .** The group  $\mathrm{SL}_2\mathbb{R}$  consists of  $2 \times 2$  matrices with real entries and determinant 1. These linear transformations act on the set of polygons in the plane: for a polygon  $P \subset \mathbb{R}^2$ , the matrix  $A \in \mathrm{SL}_2\mathbb{R}$  sends  $P$  to the new polygon  $A(P)$ . This action induces a deformation of billiard tables and associated translation surfaces: let a matrix  $A$  act on each of the reflected copies of  $T$ , and glue the stretched polygons according to the same rules. The result is a family of translation surfaces  $(X_T^A, \omega_T^A)$  parametrized by  $A \in \mathrm{SL}_2\mathbb{R}$ . See Figure 4.1.

A stretched surface  $(X_T^A, \omega_T^A)$  is conformally isomorphic to  $(X_T, \omega_T)$  if there exists a biholomorphic map  $c : X_T^A \rightarrow X_T$  which pulls  $\omega_T$  back to  $\omega_T^A$ . This can be seen in terms of the polygons: the surfaces are isomorphic if the  $N$  polygons making up  $A \cdot T$  can be cut into smaller polygons and reglued (without violating the reflection rules) to obtain  $(X, \omega)$ . If so, we say that  $A$  lies in the *Veech group*  $\Gamma_T$ . As an example, the Veech group of the square table is the lattice  $\mathrm{SL}_2\mathbb{Z}$  consisting of all  $2 \times 2$  matrices with integer entries and determinant 1.

This stretching action on billiard tables extends to an action of  $\mathrm{SL}_2\mathbb{R}$  on all of  $\Omega^1\mathcal{M}_g$ , the moduli space of translation surfaces, and the *Veech group*  $\Gamma(X, \omega)$  is the stabilizer of  $(X, \omega)$ . The action is again defined by the linear stretching of polygons: as mentioned before, any translation surface can be represented by a finite collection of polygons in the plane, with parallel sides glued by a translation, equipped with the 1-form  $dz$ . The Veech group  $\Gamma(X, \omega)$  is easily seen to be a discrete subgroup of  $\mathrm{SL}_2\mathbb{R}$ , but it is a lattice only in special cases; see [Ve], [MT]. The parametrization of an orbit  $\mathrm{SL}_2\mathbb{R} \rightarrow \Omega^1\mathcal{M}_g$  descends to a map  $\mathbb{H} \rightarrow \mathcal{M}_g$  which is a local isometry with respect to the Poincaré metric on the upper half-plane  $\mathbb{H}$  and the Teichmüller metric on  $\mathcal{M}_g$ ; see e.g. [KMS], [Ve], [Mc1]. Translation surfaces with lattice Veech group then correspond to the so-called *Teichmüller curves*, isometrically embedded algebraic curves in  $\mathcal{M}_g$ .

**4.2. Genus 1.** There is a unique holomorphic 1-form (up to scaling) on a torus, coming from the form  $dz$  on  $\mathbb{C}$ , when representing the torus as the quotient of  $\mathbb{C}$  by a lattice. The translation structure from  $dz$  is just the usual flat metric from the plane with no singularities. It is well-known that geodesics on a flat torus satisfy

FIGURE 4.2. The  $L$ -shaped table  $L(a, b)$ .

the optimal dichotomy described in §2.2. Thus any billiard table which unfolds into a genus 1 translation surface must have optimal dynamics. In fact, using the genus formula (3.1), we see that there are only 4 such tables: the three triangles with angles  $(\pi/3, \pi/3, \pi/3)$ ,  $(\pi/2, \pi/4, \pi/4)$ ,  $(\pi/2, \pi/3, \pi/6)$ , and the square.

**4.3. Genus 2.** The story in genus 2 is significantly more complicated, and a complete discussion of billiard tables has involved a sophisticated understanding of the  $\mathrm{SL}_2\mathbb{R}$  action on  $\Omega^1\mathcal{M}_2$ . McMullen found that translation surfaces in genus 2 for which  $\Gamma(X, \omega)$  is *not* a lattice have geodesics as in Figure 2.2, showing:

**Theorem 4.1.** [Mc2] *If  $(X, \omega)$  is a translation surface of genus 2 and  $\Gamma(X, \omega)$  is not a lattice, then there exists a geodesic which is neither dense nor closed.*

For a different proof when  $\omega$  has a double zero, see [Ca]. Consequently, we have:

**Corollary 4.2.** *If  $X$  has genus 2, then the following are equivalent:*

- (1) *the translation surface  $(X, \omega)$  is dynamically optimal;*
- (2) *the translation surface  $(X, \omega)$  satisfies the topological dichotomy; and*
- (3) *the Veech group  $\Gamma(X, \omega)$  is a lattice in  $\mathrm{SL}_2\mathbb{R}$ .*

McMullen went on to describe all Teichmüller curves and all dynamically optimal billiard tables in genus 2. To clarify the following statement, we need a few definitions. For any pair of integers  $a$  and  $b$  with  $b > 0$ , the billiard table  $L(a, b)$  is shown in Figure 4.2. Two tables are equivalent if their unfolded surfaces lie in the same  $\mathrm{SL}_2\mathbb{R}$  orbit.

**Theorem 4.3.** [Mc3] *Let  $T$  be a table which unfolds into a surface  $(X_T, \omega_T)$  of genus 2. Then  $T$  is dynamically optimal if and only if it is equivalent to*

- (1) *a table tiled by congruent triangles of angles  $(\pi/2, \pi/3, \pi/6)$  or  $(\pi/2, \pi/4, \pi/4)$ ;*
- (2) *an  $L$ -shaped table  $L(a, b)$  for some  $a, b \in \mathbb{Z}$ ; or*
- (3) *the triangle  $(\pi/2, 2\pi/5, \pi/10)$ .*

In fact, McMullen gave a complete description of the orbit-closures and invariant measures for the  $\mathrm{SL}_2\mathbb{R}$  action on the moduli space  $\Omega^1\mathcal{M}_2$  [Mc4].

## 5. THE COVERING CONSTRUCTION

In this final section, we discuss the basic idea which leads to Theorems 1.2 and 1.3. When a translation surface  $(Y, \eta)$  is a covering of another translation surface  $(X, \omega)$ , so that there is a covering map  $p : Y \rightarrow X$  such that  $\eta = p^*\omega$ , then the Veech groups are commensurable [GJ]; that is, after conjugating in  $\mathrm{SL}_2\mathbb{R}$ , the two groups share a finite-index subgroup. Gutkin and Judge used this relation of the Veech groups to characterize the translation surfaces (and billiard tables) with Veech groups commensurable to  $\mathrm{SL}_2\mathbb{Z}$ . When two translation surfaces are related only by a *branched covering*, any holomorphic map of finite degree  $Y \rightarrow X$  (so it may have critical points) which pulls  $\omega$  back to  $\eta$ , then the Veech group structure is not necessarily preserved.

**5.1. The branched covers of Hubert-Schmidt.** Hubert and Schmidt considered branched covers of lattice surfaces  $(X, \omega)$  branched over points of a special type. A point  $p$  in  $X$  is *periodic* if its orbit under the Veech group  $\Gamma(X, \omega)$  is finite in  $X$ . (Note that the Veech group can be viewed as the group of diffeomorphisms from  $(X, \omega)$  to itself which are linear with respect to the local Euclidean structure.) A point  $p$  is a *connection point* if any geodesic from a zero of  $\omega$  through the point  $p$  again encounters a zero of  $\omega$ . Hubert and Schmidt showed [HS]:

- (1) If  $p$  is a non-periodic connection point on a lattice surface  $(X, \omega)$ , then the subgroup  $\Gamma(X, \omega, p) := \{\gamma \in \Gamma(X, \omega) : \gamma(p) = p\}$  is infinitely generated; and
- (2) If a branched cover  $(Y, \eta) \rightarrow (X, \omega)$  over a lattice surface  $(X, \omega)$  is branched only over a connection point  $p$ , then  $\Gamma(Y, \eta)$  is commensurable with  $\Gamma(X, \omega, p)$ , and the surface  $(Y, \eta)$  satisfies the topological dichotomy.

They show further that the second statement holds when branching over more than one connection point if the base surface  $(X, \omega)$  has a property they call *strong holonomy type*.

**5.2. Examples of Smillie-Weiss and Cheung-Hubert-Masur.** Smillie and Weiss made use of the Hubert-Schmidt construction to prove Theorem 1.2, concentrating on branched covers with a single ramification point in the base translation space of genus  $g \geq 2$ . They cleverly combine two facts: one is the simple observation that the forgetful map from  $\mathcal{M}_{g,1}$  down to  $\mathcal{M}_g$  has compact fibers. The second is Masur's theorem that a minimal non-uniquely ergodic direction gives rise to a particular  $\mathrm{SL}_2\mathbb{R}$ -deformation which tends to infinity in  $\Omega^1\mathcal{M}_g$  [Ma1].

In [CHM], the authors again use one of the Hubert-Schmidt branched covers, but branched over *two* points. The special example of Theorem 1.3 arises in the following way. Begin with the translation surface of the triangle with angles  $(\pi/2, \pi/5, 3\pi/10)$ . This triangle unfolds into two reflected copies of the regular pentagon, forming a surface of genus 2. It is dynamically optimal [Ve], and it is equivalent to one of the *L*-shaped tables in McMullen's classification Theorem 4.3 [Mc1, §9]. The centers of

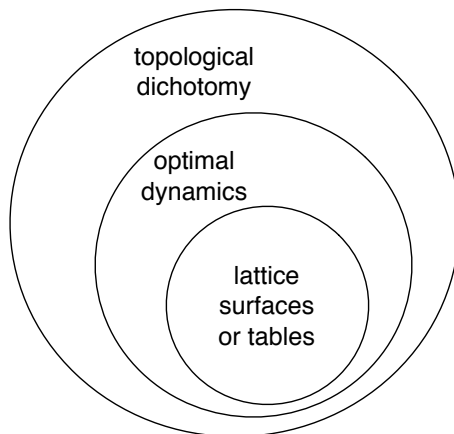


FIGURE 5.1. A diagram of inclusions for billiard tables or translation surfaces.

the two pentagons are connection points (see definition in §5.1). Taking a double cover of this genus 2 surface, branched over the two centers, produces the translation surface of genus 4 which is the unfolding of the triangle  $(2\pi/5, 3\pi/10, 3\pi/10)$ . The new surface satisfies the topological dichotomy by the arguments of [HS], but Cheung, Hubert, and Masur show that it has non-uniformly distributed, dense geodesics.

**5.3. Lattice tables and a dynamical characterization.** The schematic of Figure 5.1 indicates the relative inclusions of tables which have lattice Veech group, those which are dynamically optimal, and those satisfying topological dichotomy. By McMullen’s theorem (Corollary 4.2), the sets coincide for translation surfaces of genus 2. The examples of [SW] and [CHM] show that the containments are strict in the setting of translation surfaces of arbitrary genus. Billiard tables have not yet been found which are dynamically optimal without lattice Veech group, but a search is under way. Further investigations are also in progress about possible characterizations of the lattice condition, because it is the geometry of the  $SL_2\mathbb{R}$  action on the moduli space of translation surfaces which drives most of the interest in this class of billiards.

## REFERENCES

- [Ca] K. Calta. Veech surfaces and complete periodicity in genus two. *J. Amer. Math. Soc.* **17**(2004), 871–908.
- [Ch] Y. Cheung. Hausdorff dimension of the set of nonergodic directions. *Ann. of Math. (2)* **158**(2003), 661–678. With an appendix by M. Boshernitzan.
- [CHM] Y. Cheung, P. Hubert, and H. Masur. Topological dichotomy and strict ergodicity for translation surfaces. *Ergodic Theory Dynam. Systems* **28**(2008), 1729–1748.
- [GJ] E. Gutkin and C. Judge. Affine mappings of translation surfaces: geometry and arithmetic. *Duke Math. J.* **103**(2000), 191–213.

- [HS] P. Hubert and T. A. Schmidt. Infinitely generated Veech groups. *Duke Math. J.* **123**(2004), 49–69.
- [KZ] A. B. Katok and A. N. Zemljakov. Topological transitivity of billiards in polygons. *Mat. Zametki* **18**(1975), 291–300.
- [KMS] S. Kerckhoff, H. Masur, and J. Smillie. Ergodicity of billiard flows and quadratic differentials. *Ann. of Math. (2)* **124**(1986), 293–311.
- [Ma1] H. Masur. Hausdorff dimension of the set of nonergodic foliations of a quadratic differential. *Duke Math. J.* **66**(1992), 387–442.
- [Ma2] H. Masur. Ergodic theory of translation surfaces. In *Handbook of dynamical systems. Vol. 1B*, pages 527–547. Elsevier B. V., Amsterdam, 2006.
- [MS] H. Masur and J. Smillie. Hausdorff dimension of sets of nonergodic measured foliations. *Ann. of Math. (2)* **134**(1991), 455–543.
- [MT] H. Masur and S. Tabachnikov. Rational billiards and flat structures. In *Handbook of dynamical systems, Vol. 1A*, pages 1015–1089. North-Holland, Amsterdam, 2002.
- [Mc1] C. T. McMullen. Billiards and Teichmüller curves on Hilbert modular surfaces. *J. Amer. Math. Soc.* **16**(2003), 857–885.
- [Mc2] C. T. McMullen. Teichmüller curves in genus two: the decagon and beyond. *J. Reine Angew. Math.* **582**(2005), 173–199.
- [Mc3] C. T. McMullen. Teichmüller curves in genus two: torsion divisors and ratios of sines. *Invent. Math.* **165**(2006), 651–672.
- [Mc4] C. T. McMullen. Dynamics of  $SL_2(\mathbb{R})$  over moduli space in genus two. *Ann. of Math. (2)* **165**(2007), 397–456.
- [SW] J. Smillie and B. Weiss. Veech’s dichotomy and the lattice property. *Ergodic Theory Dynam. Systems* **28**(2008), 1959–1972.
- [Ve] W. A. Veech. Teichmüller curves in moduli space, Eisenstein series and an application to triangular billiards. *Invent. Math.* **97**(1989), 553–583.

DEPARTMENT OF MATHEMATICS, STATISTICS, AND COMPUTER SCIENCE, UNIVERSITY OF ILLINOIS AT CHICAGO, DEMARCO@MATH.UIC.EDU

## History of the Kervaire Invariant Problem

It always an important event when a famous old mathematical problem is solved. But with the solution, the history of the problem—why people thought it was so important in the first place, what were the twists and turns on the path to its solution—is often not recorded.

Both technical and expository treatments of the solution to the “Kervaire Invariant One” problem are being prepared, and will appear in time. I hope they’ll repeat some of the rich history of the problem. But in this issue of the Current Events Bulletin book we’ve chosen to focus on that history: at Mike Hopkins suggestion, I solicited reports from some of the major player in homotopy theory about the history of the problem and the mathematical thoughts and ideas that were related to it and gave it such importance in the history of the field. William Browder and Mark Mahowald—now joined by Paul Goerss—graciously agreed to provide something. You have before you, therefore, reports from great masters who struggled with this problem, and created an immense amount of beautiful mathematics in the process—without, however, reaching the eventual solution. We’ll hear about the solution from Mike Hopkins in person. . .

David Eisenbud





# HISTORY OF THE KERVAIRE INVARIANT PROBLEM

WILLIAM BROWDER

The history of this invariant could very well be considered to start with the paper of Pontryagin in 1938, where he introduced Framed Bordism (as it is now known) as a tool to calculate homotopy groups of spheres, using smooth manifolds. He proved that the second stable homotopy group of the  $n$ -sphere was zero, but this was soon shown to be incorrect by algebraic methods. The problem was the absence of the Kervaire invariant.

For oriented closed manifolds of dimension  $4k$  the middle dimensional intersection pairing defines a nonsingular symmetric bilinear form over the integers, whose signature gives a famous algebraic invariant often called the index of the manifold. For dimensions  $4k+2$  the intersection bilinear form is skew symmetric, and thus can be put in canonical form, so no apparent such invariants exist. But if, in the mod 2 version, we can enrich the intersection form to be associated to a quadratic form, that quadratic form has an invariant, called the Arf invariant after its discoverer.

I like to call this invariant the democratic invariant as it can be defined as follows:

Let  $V$  be a vector space of finite dimension over the integers mod 2, and let  $q : V \rightarrow Z/2$  be a quadratic form with associated bilinear form  $f$  (i.e.,  $q(a+b) = q(a) + q(b) + f(a,b)$ ). Consider  $q$  to be a vote between 0 and 1 (the candidates) among the elements of  $V$  (voters) and the Arf invariant of  $q$  is defined to be the winner of the election. If the bilinear form  $f$  is nonsingular the election is decisive. However, in the general case the election is a tie if and only if there is some element  $r$  of  $V$  such that  $f(r,x) = 0$  for all  $x$  in  $V$ , but  $q(r) = 1$ . (The election reaches a clear result unless some radical element votes positively).

Pontryagin had failed to note that an underlying obstruction to the process he was carrying out in dimension 2 was quadratic rather than linear, so that its Arf invariant was an obstruction for his argument, but he corrected this mistake in a later paper in 1955.

This might be considered the prehistory of the topological invariant, and in my view the history properly begins with the paper of Kervaire in 1960 where he constructed a PL 10-manifold which was not of the homotopy type of a smooth manifold. In it he constructed a cohomology operation from dimension 5 to 10, for a 4-connected closed 10-manifold which could be framed (stable tangent bundle trivial) on the complement of a point, and this operation was quadratic. In his example the Arf invariant was non trivial, while for any smooth manifold of that type, he proved it would be trivial because of the vanishing of some homotopy.

Kervaire's operation is defined in analogous circumstances for all dimensions and in the famous 1962 paper of Kervaire and Milnor defined the middle dimensional surgery obstruction in dimensions of the form  $4k+2$ . It followed from the surgery theory developed there that this defined a Framed Bordism invariant. You could

do surgery on the framed manifold to make it  $2k$ -connected to define the invariant and do surgery on a framed bordism to make the bordism similarly connected to prove it well defined.

If a PL manifold  $M$  has a trivial tangent bundle (or more properly microbundle) over the complement of a point it is in fact smoothable away from that point, from the theory of smoothing of PL manifolds of Mazur and Hirsch-Mazur. If  $M$  is of dimension  $2n$  and  $(n-1)$ -connected, then its  $n^{\text{th}}$  homology has a basis of embedded spheres, and the normal bundle of each of these spheres is stably trivial. When  $n = 1, 3$  or  $7$  it is trivial, but in other dimensions there are nontrivial possibilities, and if  $n$  is odd there is a single nontrivial possibility, namely the tangent bundle to the  $n$ -sphere. The quadratic form in these cases is simply given by whether this normal bundle is trivial or not. (For  $n = 1, 3$  or  $7$ , the definition of a quadratic form is related to the framing and is not homotopy invariant). The cohomology operation of Kervaire detects this non-triviality, and its model is actually the Thom complex of this tangent bundle.

If the Arf invariant of this form is zero, one can find enough embedded products  $S^n \times R^n$  representing a basis of the middle cohomology to carry out surgery to make  $M$  into a homotopy sphere, and otherwise you cannot.

Thus the question of whether or not a framed manifold could have a nonzero Kervaire invariant then became a central question for differential topology, equivalent to the calculation of the subgroup of homotopy spheres which bounded framed manifolds in dimension  $4k+1$ . The answer was yes for dimensions  $2, 6$  and  $14$  because of the parallelizability of the spheres of dimensions  $1, 3$  and  $7$ , but remained open for other dimensions of the form  $4k+2$ .

E. H. Brown in 1965 showed that for dimension  $8k+2$  Spin manifolds, the cohomology operation could be defined on the middle dimension without assuming  $4k$ -connectivity, so that the Kervaire Invariant could be made a Spin bordism invariant, using surgery only to make the manifold simply connected. Subsequently, he and F. P. Peterson in 1966 used this to show the Kervaire Invariant vanished on Framed Bordism in dimension  $8k+2$ .

Then in 1968, I proved that the Kervaire Invariant vanished on Framed Bordism in dimensions different from  $2^n - 2$ , and related possible nonvanishing in those dimensions to the existence of elements in the homotopy of spheres related to certain elements in the Adams spectral sequence. It turned out that this element had already been constructed by Mahowald and Tangora in dimension  $30$  and such an element was later constructed in dimension  $62$  by Barratt, Jones and Mahowald.

My method was to define the quadratic form by means of a functional Steenrod operation on a subgroup of the middle cohomology mod  $2$  which allowed me to define the form on a manifold  $M$  which had been "oriented" in a theory in which the  $2k+2$  Wu class was zero, a condition satisfied by any  $4k+2$  manifold. (The Wu classes are defined using the Steenrod operations in  $M$  and are directly related to the Stiefel-Whitney classes). A subtlety was that everything depended on how you made this Wu class vanish, how you chose the "orientation". This definition allowed one to define the Kervaire invariant immediately on the framed (or otherwise "oriented") manifold without doing any surgery or other geometrical operation, and so gave

a definition in a purely homotopy theoretical context of spaces satisfying Poincaré duality.

For a smooth manifold  $M$ , my definition can be translated into a condition on extending vector fields on submanifolds of  $M$  representing the middle dimensional mod 2 cohomology.

A simpler proof of my theorem on Framed Bordism was given by Jones and Rees, and Jones gave a beautiful construction of the 30 dimensional manifold representing the Mahowald-Tangora homotopy element.

After the results in dimensions 30 and 62, attention turned to dimension 126, the first open case, but this has resisted concerted attempts by many strong homotopy theorists and still is unknown. Many had tried to prove that all of these possible elements (or manifolds) existed but, conscious of the Hopf invariant 1 results (only three possible dimensions 1, 3 and 7), some began to try to prove that they did not exist beyond some dimension. Now this has been carried out by Hill, Hopkins and Ravenel for dimensions greater than 126.

Much other work has been done on the Kervaire Invariant because of its importance in surgery theory, e.g., Sullivan (product formula), Ranicki (algebraic surgery), and others.

#### REFERENCES

- [1] Pontryagin, L. S., *Gladkie mnogoobraziya i ih primeneniya v teorii gomotopii*. (Russian) [Smooth manifolds and their applications in homotopy theory.], Trudy Mat. Inst. im. Steklov. no. 45., Izdat. Akad. Nauk SSSR, Moscow, 1955. 139 pp.
- [2] Kervaire, Michel A., *A manifold which does not admit any differentiable structure*, Comment. Math. Helv. 34 1960, 257–270.
- [3] Kervaire, Michel A.; Milnor, John W., Groups of homotopy spheres. I, Ann. of Math. (2) 77 1963, 504–537.
- [4] Brown, Edgar H., Jr., *Note on an invariant of Kervaire*, Michigan Math. J. 12 1965, 23–24.
- [5] Brown, Edgar H., Jr.; Peterson, Franklin P., *The Kervaire invariant of  $(8k + 2)$ -manifolds*, Bull. Amer. Math. Soc. 71 1965, 190–193.
- [6] Browder, William, *The Kervaire invariant of framed manifolds and its generalization*, Ann. of Math. (2) 90 1969, 157–186.
- [7] Browder, William, *Surgery on simply-connected manifolds*, Ergebnisse der Mathematik und ihrer Grenzgebiete, Band 65. Springer-Verlag, New York-Heidelberg, 1972. ix+132 pp.
- [8] Mahowald, Mark; Tangora, Martin, *Some differentials in the Adams spectral sequence*, Topology 6 1967, 349–369.
- [9] Barratt, M. G.(1-NW); Jones, J. D. S.(4-WARW); Mahowald, M. E.(1-NW), *Relations amongst Toda brackets and the Kervaire invariant in dimension 62*, J. London Math. Soc. (2) 30 (1984), no. 3, 533–550.
- [10] Jones, John; Rees, Elmer, *Kervaire’s invariant for framed manifolds. Algebraic and geometric topology* (Proc. Sympos. Pure Math., Stanford Univ., Stanford, Calif., 1976), Part 1, pp. 141–147, Proc. Sympos. Pure Math., XXXII, Amer. Math. Soc., Providence, R.I., 1978.
- [11] Jones, John D. S., *The Kervaire invariant of extended power manifolds*, Topology 17 (1978), no. 3, 249–266.
- [12] Jones, John; Rees, Elmer, *Kervaire’s invariant for framed manifolds. Algebraic and geometric topology* (Proc. Sympos. Pure Math., Stanford Univ., Stanford, Calif., 1976), Part 1, pp. 141–147, Proc. Sympos. Pure Math., XXXII, Amer. Math. Soc., Providence, R.I., 1978.
- [13] Brown, Edgar H., Jr., *The Kervaire invariant and surgery theory. (English summary) Surveys on surgery theory*, Vol. 1, 105–120, Ann. of Math. Stud., 145, Princeton Univ. Press, Princeton, NJ, 2000.

- [14] Barratt, M. G.; Jones, J. D. S.; Mahowald, M. E., *The Kervaire invariant and the Hopf invariant*, Algebraic topology (Seattle, Wash., 1985), 135–173, Lecture Notes in Math., 1286, Springer, Berlin, 1987.
- [15] Barratt, M. G.; Jones, J. D. S.; Mahowald, M. E., *The Kervaire invariant problem. Proceedings of the Northwestern Homotopy Theory Conference*, (Evanston, Ill., 1982), 9–22, Contemp. Math., 19, Amer. Math. Soc., Providence, RI, 1983.
- [16] Brown, Edgar H., Jr., *The Arf invariant of a manifold. 1969 Conf. on Algebraic Topology*, (Univ. of Illinois at Chicago Circle, Chicago, Ill., 1968) pp. 9–18 Univ. of Illinois at Chicago Circle, Chicago, Ill.

# The Kervaire invariant in homotopy theory

Mark Mahowald and Paul Goerss\*

November 16, 2009

## Abstract

In this note we discuss how the first author came upon the Kervaire invariant question while analyzing the image of the  $J$ -homomorphism in the EHP sequence.

One of the central projects of algebraic topology is to calculate the homotopy classes of maps between two finite CW complexes. Even in the case of spheres – the smallest non-trivial CW complexes – this project has a long and rich history.

Let  $S^n$  denote the  $n$ -sphere. If  $k < n$ , then all continuous maps  $S^k \rightarrow S^n$  are null-homotopic, and if  $k = n$ , the homotopy class of a map  $S^n \rightarrow S^n$  is detected by its degree. Even these basic facts require relatively deep results: if  $k = n = 1$ , we need covering space theory, and if  $n > 1$ , we need the Hurewicz theorem, which says that the first non-trivial homotopy group of a simply-connected space is isomorphic to the first non-vanishing homology group of positive degree. The classical proof of the Hurewicz theorem as found in, for example, [28] is quite delicate; more conceptual proofs use the Serre spectral sequence.

Let us write  $\pi_i S^n$  for the  $i$ th homotopy group of the  $n$ -sphere; we may also write  $\pi_{k+n} S^n$  to emphasize that the complexity of the problem grows with  $k$ . Thus we have  $\pi_{n+k} S^n = 0$  if  $k < 0$  and  $\pi_n S^n \cong \mathbb{Z}$ . Given the Hurewicz theorem and knowledge of the homology of Eilenberg-MacLane spaces it is relatively simple to compute that

$$\pi_{n+1} S^n \cong \begin{cases} 0, & n = 1; \\ \mathbb{Z}, & n = 2; \\ \mathbb{Z}/2\mathbb{Z}, & n \geq 3. \end{cases}$$

The generator of  $\pi_3 S^2$  is the Hopf map; the generator in  $\pi_{n+1} S^n$ ,  $n > 2$  is the *suspension* of the Hopf map. If  $X$  has a basepoint  $y$ , the suspension  $\Sigma X$  is given by

$$\Sigma X = S^1 \times X / (S^1 \times y \cup 1 \times X)$$

where  $1 \in S^1 \subseteq \mathbb{C}$ . Then  $\Sigma S^n \cong S^{n+1}$  and we get a suspension homomorphism

$$E : \pi_{n+k} S^n \rightarrow \pi_{(n+1)+k} S^{n+1}.$$

---

\*The second author was partially supported by the National Science Foundation.

By the Freudenthal Suspension Theorem, this map is onto if  $k \leq n - 1$ , and if  $k < n - 1$  it is an isomorphism. The common value of this group for large  $n$  is the  $k$ th stable homotopy group of spheres, written  $\pi_k^s S^0$ . For short, we may write

$$\operatorname{colim}_k \pi_{n+k} S^n = \pi_k^s S^0.$$

Note that this formula makes sense even if  $k < 0$ .

There has been a great deal of computation in the stable homotopy groups of spheres; see, for example, Appendix 3 of [26]. The answer is fairly complete for  $k$  up to about 60; if we divide out by 2- and 3-torsion, this can be improved to about  $k = 1000$ . However, we are a long way from any sort of complete calculation. Research since the mid-1970s has shifted to the investigation of large-scale phenomena, especially after the paper by Miller, Ravenel, and Wilson [24] on periodic phenomena and the proofs of Ravenel's nilpotence conjectures by Devinatz, Hopkins, and Smith [10, 13].

Historically, the Kervaire invariant arose in Pontryagin's calculation of  $\pi_2^s S^0$ . He noted that  $\pi_k^s S^0$  is isomorphic to the group of cobordism classes of *framed*  $k$ -manifolds; that is, differentiable manifolds with a chosen trivialization of the stable normal bundle. Let  $\mathbb{F}_2$  be the field with two elements. and let  $M$  be a connected framed manifold of dimension  $k = 4m - 2$ . By collapsing all but the top cell of  $M$  we obtain a map

$$M \longrightarrow S^{4m-2}$$

which is an isomorphism of  $H^{4m-2}(-, \mathbb{F}_2)$ . Using surgery [6, 29] we can try to build a cobordism from  $M$  to the sphere. This may not be possible, but we do find that the non-singular pairing

$$\lambda : H^{2m-1}(M, \mathbb{F}_2) \times H^{2m-1}(M, \mathbb{F}_2) \rightarrow H^{4m-2}(M, \mathbb{F}_2) \cong \mathbb{F}_2$$

given by Poincaré duality has a quadratic refinement  $\mu$ ; that is, there is a function  $\mu : H^{2m-1}(M, \mathbb{F}_2) \rightarrow \mathbb{F}_2$  so that

$$\mu(x + y) + \mu(x) + \mu(y) = \lambda(x, y).$$

Up to isomorphism, the pair  $(H^{2m-1}(M, \mathbb{F}_2), \mu)$  is completely determined by the *Arf Invariant*. This invariant is 1 if  $\mu(x) = 1$  for the majority of the elements in  $H^{2m-1}(M, \mathbb{F}_2)$ ; otherwise it is 0.<sup>1</sup> The *Kervaire invariant* of  $M$  is the Arf invariant of this quadratic refinement.

After first getting the computation wrong, Pontryagin [25] noted that for a particular framing of  $S^1 \times S^1$ , the Kervaire invariant was non-zero, giving a non-trivial cobordism class. Then  $\pi_2^s S^0 \cong \mathbb{Z}/2\mathbb{Z}$  generated by this element.

To study the higher homotopy groups of spheres, we must consider more sophisticated methods. One such is the Adams spectral sequence

$$\operatorname{Ext}_A^s(\mathbb{F}_2, \Sigma^t \mathbb{F}_2) \Longrightarrow \pi_{t-s}^s S^0 \otimes \mathbb{Z}_2.$$

---

<sup>1</sup>For this reason, Browder has called the Arf invariant the “democratic invariant”.

Here  $A$  is the Steenrod algebra,  $\mathbb{Z}_2$  is the 2-adic integers, and  $\Sigma^t \mathbb{F}_2 = \tilde{H}^*(S^t, \mathbb{F}_2)$ . The Kervaire invariant elements are then classes

$$h_j^2 \in \text{Ext}_A^2(\mathbb{F}_2, \Sigma^{2^{j+1}} \mathbb{F}_2)$$

which could detect elements in  $\pi_{2^{j+1}-2}^s S^0$ . If  $j = 1$ , this element detects Pontryagin's class.

In his work in smoothing theory, Kervaire [16] constructed a topological manifold of dimension  $4m - 2$  for  $m \neq 1, 2, 4$  which had Kervaire invariant one and which was smooth if a point was removed. The question then became "Is the boundary sphere smoothable?" Browder [5], proved that it was smoothable if and only if  $m = 2^{j-1}$ ,  $j \geq 4$  and if the elements  $h_j^2$  detected a homotopy class. The homotopy class was constructed for  $j = 4$  before Browder's work by Peter May in his thesis [22]; the finer properties of this element were uncovered in [4]. The class in dimension 62 (that is,  $j = 5$ ) was constructed later in [3].

Hill, Hopkins, and Ravenel [11] have shown that for  $j \geq 7$  the class  $h_j^2$  is not a permanent cycle in the Adams spectral sequence and cannot detect a stable homotopy class. This settles Browder's question in all but one case. Their proof is a precise and elegant application of equivariant stable homotopy theory. It is also very economical: they develop the minimum amount needed to settle exactly the question at hand. The very economy of this solution leaves behind numerous questions for students of  $\pi_*^s S^0$ . One immediate problem is to find the differential on  $h_j^2$  in the Adams spectral sequence. The target would be an important element which we as yet have no name for.

The Kervaire invariant and Arf invariant have appeared in other places and guises in geometry and topology. For example, it is possible to formulate the Kervaire invariant question not for framed manifolds, but for oriented manifolds whose structure group reduces to  $SO(1) = S^1$ . In this formulation, Ralph Cohen, John Jones, and the first author showed that the problem had a positive solution [8]. The relevant homotopy classes are in the stable homotopy of the Thom spectrum  $MSO(1)$ ; they are constructed using a variant of the methods of [19], which certainly don't extend to the sphere. Note that  $MSO(1) = \Sigma^{-2} \mathbb{C}P^\infty$  is an infinite  $CW$  spectrum, but still relatively small. This may be the best of all positive worlds for this problem.

In [7], Brown found a way to extend the Kervaire invariant to another more general class of manifolds. And, by contemplating work of Witten, Hopkins and Singer found an application of the Arf invariant in dimension 6 to some problems in mathematical physics. See [12].

Parallel to this geometric story, the Kervaire invariant problem also arose in an entirely different line of research in homotopy theory, and here the negative solution of [11] leaves as many questions as it answers. This line of inquiry, long studied by the first author, asks just how the stable homotopy groups of spheres are born. To make this question precise, we must introduce the EHP sequence. This was discovered by James [15] in the mid-1950s and related techniques were exploited by Toda to great effect in his landmark book [27].

If  $X$  is a based space, let  $\Omega X$  denote the space of based loops in  $X$ . In his work on loop spaces [14], James produced a small  $CW$  complex with the



homotopy type of  $\Omega S^{n+1}$  and later he noticed that this gave a splitting

$$\Sigma \Omega S^{n+1} \simeq \bigvee_{t>0} S^{nt+1}.$$

Here  $\bigvee$  is the one-point union or *wedge*. By collapsing all factors of the wedge except for  $t = 2$  and then taking the adjoint, we obtain the first *Hopf invariant*

$$H : \Omega S^{n+1} \rightarrow \Omega S^{2n+1}$$

There is also the map  $E : S^n \rightarrow \Omega S^{n+1}$  adjoint to the identity; it induces the suspension homomorphism on homotopy groups. A calculation with the Serre spectral sequence shows that

$$S^n \xrightarrow{E} \Omega S^{n+1} \xrightarrow{H} \Omega S^{2n+1}$$

is a fiber sequence when localized at 2. As a consequence there is a long exact sequence in homotopy groups, once we divide out by the odd torsion:

$$(1) \quad \dots \rightarrow \pi_{i+2} S^{2n+1} \xrightarrow{P} \pi_i S^n \xrightarrow{E} \pi_{i+1} S^{n+1} \xrightarrow{H} \pi_{i+1} S^{2n+1} \rightarrow \dots$$

This is the EHP sequence. As mentioned,  $E$  is the suspension map and  $H$  is the Hopf invariant. The map  $P$  is more difficult to describe; however we do have that if  $\alpha \in \pi_* S^{n-1}$ , then (up to sign)

$$P(E^{n+2}\alpha) = [\iota_n, E\alpha]$$

where  $[-, -]$  is the Whitehead product and  $\iota_n \in \pi_n S^n$  is the identity. Thus, for example,  $P(\iota_{2n+1}) = [\iota_n, \iota_n]$ .

**From this point forward in this note, we will implicitly localize all groups at the prime 2.**

The EHP sequence gives an inductive method for calculating the homotopy groups of spheres; the key is to do double induction on  $n$  and  $k$  in  $\pi_{n+k} S^n$ . To this end we reindex the subscripts in Equation (1) and write a triangle

$$(2) \quad \begin{array}{ccc} \pi_{n+k} S^n & \xrightarrow{E} & \pi_{(n+1)+k} S^{n+1} \\ & \swarrow P \text{ (dotted)} & \nwarrow H \\ & \pi_{(2n+1)+(k-n)} S^{2n+1} & \end{array}$$

for the EHP sequence. The dotted arrow indicates a map of degree  $-1$ . Then, assuming we know  $\pi_{m+i} S^m$  for all  $m \leq n$  and for all  $i < k$ , we can try to calculate  $\pi_{(n+1)+k} S^{n+1}$ . Coupled with the unstable Adams Spectral Sequence, it is possible to do low dimensional calculations very quickly – but, as with all algebraic approximations to the homotopy groups of spheres, it gets difficult

fairly soon.<sup>2</sup> Tables for this computation can be found in a number of places; see, for example, [23] or §I.5 of [26].

**Question 1.1.** Suppose  $\alpha \in \pi_k^s S^0$  is a stable element.

1. What is the smallest  $n$  so that  $\alpha$  is in the image of  $\pi_{n+k} S^n \rightarrow \pi_k^s S^0$ ? Then  $S^n$  is the *sphere of origin*.
2. Suppose  $S^n$  is the sphere of origin of  $\alpha$  and  $a$  is a class in  $\pi_{n+k} S^n$  which suspends to  $\alpha$ . What is  $H(a)$ ? This is “the” *Hopf invariant* of  $\alpha$ .

**Technical Warning 1.2.** As phrased, the second question is not precise, as there maybe more than one  $a$  which suspends to  $\alpha$ . There are several ways out of this difficulty. One is to ignore it. In practice, this works well. Another is to note that the *EHP* sequences, as in Equation (2) assemble into an exact couple which gives a spectral sequence

$$E_{k,n}^1 = \pi_{n+k} S^{2n-1} \implies \pi_k^s S^0.$$

Then questions (1) and (2) can be rephrased by asking for the non-zero element in  $E^\infty$  which detects  $\alpha$ .

It is a feature of this spectral sequence that the  $E^2$ -page is an  $\mathbb{F}_2$ -vector space. This means, for example, that elements of high order must have high sphere of origin. Charts for this spectral sequence can be developed from [17] and can be found in explicit form in [23], which is based on work of the first author.

**Example 1.3.** As a simple test case, the sphere of origin the generator  $\eta \in \pi_1^s S^0 \cong \mathbb{Z}/2\mathbb{Z}$  is  $S^2$  and  $a$  can be taken to the Hopf map  $S^3 \rightarrow S^2$ . The Hopf invariant of this map is (up to sign) the identity  $\iota_3 \in \pi_3 S^3$ . We can ask when  $\iota_{2n-1} \in \pi_{2n-1} S^{2n-1}$  can be the Hopf invariant of some stable class. This is the *Hopf invariant one* problem, settled by Adams in [1]: it only happens when  $n$  is 2, 4, or 8; the resulting stable classes are  $\eta \in \pi_1^s S^0$ ,  $\nu \in \pi_3^s S^0$ , and  $\sigma \in \pi_7^s S^0$ .

There are instructive reformulations of the Hopf invariant one problem. First, by the *EHP* sequence,  $\iota_{2n-1}$  is the Hopf invariant of a stable class if and only if

$$[\iota_{n-1}, \iota_{n-1}] = 0 \in \pi_{2n-1} S^{n-1}.$$

Thus we are asking about the behavior of the Whitehead product.

Second, an argument with Steenrod operations shows there is an element of Hopf invariant one if and only if the element

$$h_j \in \text{Ext}_A^1(\mathbb{F}_2, \Sigma^{2^j} \mathbb{F}_2)$$

survives to  $E_\infty$  in the Adams spectral sequence. It is this last question that Adams settled showing

$$d_2 h_j = h_0 h_{j-1}^2$$

---

<sup>2</sup>Doug Ravenel has dubbed this general observation “The Mahowald Uncertainty Principle”.

if  $j \geq 4$ .

This example, while now part of our basic tool kit, remains very instructive for the interplay of stable and unstable information, and the role of the Adams spectral sequence. Notice also that we changed our question in the middle of the discussion.

**Question 1.4.** Let  $\alpha \in \pi_k^s S^0$  be a stable element, then  $\alpha$  desuspends uniquely to  $\pi_{n+k} S^n$  if  $n > k + 1$ . Suppose  $2n - 1 > k + 1$ . Is

$$\alpha \in \pi_{2n-1+k} S^{2n-1}$$

the Hopf invariant of a stable element in  $\pi_{2n-1+k} S^n$ ?

The solution of the Hopf invariant one problem, completely answers this question for a generator of  $\pi_0^s S^0$ . We will have other examples below.

It is exactly in thinking about Question (1.4) that the first author came to the Kervaire invariant problem.

In the middle 1960s, Adams [2] (with an addendum by the first author at the prime 2 [18]) wrote down an infinite family of non-zero elements in the homotopy groups of spheres. These elements we now call the *image of J*, and they were the first example of “periodic” families. They are easy to define, although less easy to show they are non-trivial.

Let  $SO(n)$  be the special orthogonal group. Then  $SO(n)$  acts on  $S^n$  by regarding  $S^n$  as the one-point compactification of  $\mathbb{R}^n$ . This action defines a map

$$SO(n) \rightarrow \text{map}_*(S^n, S^n)$$

from  $SO(n)$  to the space of pointed maps. Taking the adjoint, assembling all  $n$ , and applying homotopy yields map

$$(3) \quad J : \pi_k SO \rightarrow \pi_k^s S^0.$$

By Bott periodicity, we know the homotopy groups of  $SO$ . What Adams did was compute the image. To state the result, let  $k > 0$  be an integer, let  $\nu_2(-)$  denote 2-adic valuation, and define

$$\lambda(k) = \nu_2(k + 1) + 1.$$

Thus  $\lambda(7) = 4$  and  $\lambda(11) = 3$ . The image of the  $J$ -homomorphism lies in a split summand  $\text{Im}(J)_* \subseteq \pi_*^s S^0$  with  $\text{Im}(J)_1 \cong \mathbb{Z}/2\mathbb{Z}$  generated by  $\eta$  and for  $k \geq 2$ ,

$$\text{Im}(J)_k = \begin{cases} \mathbb{Z}/2^{\lambda(k)}\mathbb{Z} & k = 8t - 1, 8t + 3; \\ \mathbb{Z}/2\mathbb{Z} & k = 8t, 8t + 2; \\ \mathbb{Z}/2\mathbb{Z} \times \mathbb{Z}/2\mathbb{Z} & k = 8t + 1; \\ 0 & k = 8t + 4, 8t + 5, 8t + 6. \end{cases}$$

Let's write  $\rho_{8t-1}$  and  $\zeta_{8t+3}$  for the generators of the groups in degrees  $8t - 1$  and  $8t + 3$  respectively. Some of these elements are familiar; for example,  $\nu = \zeta_3$  and  $\sigma = \rho_7$ . The elements  $\eta\rho_{8t-1}$  and  $\eta^2\rho_{8t-1}$  are non-zero in  $\text{Im}(J)_*$ . There

is another generator  $\mu_{8t+1}$  in degree  $8t + 1$ ;  $\eta\mu_{8t+1} \neq 0$  and  $\eta^2\mu_{8t+1} = 4\zeta_{8t+3}$ . Despite the name, the  $J$ -homomorphism of Equation (3) is not onto  $\text{Im}(J)_*$ , as the elements  $\mu_{8t}$  and  $\eta\mu_{8t}$  are not in the image, although we see that they are intimately connected to that image. In fact, we can think of  $\mu_1$  as  $\eta$ ; then

$$\eta^2\mu_1 = \eta^3 = 4\nu = \zeta_3$$

and the equation  $\eta^2\mu_{8t+1} = 4\zeta_{8t+3}$  is then forced by the periodic behavior of these elements.

There were several revealing new features to this family. One was that it was infinite: this was the first systematic collection of elements produced in the stable homotopy groups of spheres and took us beyond the era of stem-by-stem calculations. Another feature was that this was the first of what we now call periodic families of stable homotopy classes. The attempt to understand stable homotopy theory in terms of periodic families led to a reorganization of the field, including the work of Miller, Ravenel, and Wilson [24], and the Ravenel's nilpotence conjectures, proved by Devinatz, Hopkins, and Smith [10, 13].

We now say that  $\text{Im}(J)_*$  is the  $v_1$ -periodic homotopy groups of spheres. We won't dwell on this point, but in modern language (a language not available in the 1960s), we say the the composite

$$\text{Im}(J)_* \xrightarrow{\subseteq} \pi_*^s S^0 \longrightarrow \pi_* L_{K(1)} S^0$$

is an isomorphism in degrees greater than 1 and an injection in degree 1. Here  $L_{K(1)} S^0$  is the localization of the sphere spectrum at the  $K$ -theory with coefficients in  $\mathbb{Z}/2\mathbb{Z}$ .

In the mid-1960s, the first author began an extensive study of the image of  $J$  in the EHP sequence; the first results appeared in [17] and there was followup paper [20] almost fifteen years later.

The sphere of origin and the Hopf invariants of the elements in the image of  $J$  are all known. For example, the sphere of origin of  $\zeta_{8k+3}$  is  $S^5$  and its Hopf invariant is an element on  $S^9$  which suspends to  $2^{\lambda(8k-1)-1}\rho_{8k-1}$ ; in particular, the Hopf invariant of  $\zeta_{8k+3}$  is another element in the image of  $J$ . The remarkable fact is that this is (almost) true in general; the exception is  $\eta\rho_{8k-1}$  which has sphere of origin  $S^3$  and Hopf invariant  $\nu\zeta_{8k-5}$  on the 5-sphere.<sup>3</sup> In this sense the image  $J$  is very nearly a closed family. The detailed answers are nicely laid out in [23].

But there are exceptions. Once the the sphere-of-origin and Hopf invariant calculations have been answered for the elements in the image of  $J$ , there are still a few elements left that could be Hopf invariants of new elements in the stable homotopy groups of spheres. There are some sporadic examples (see the table below) and there are also two infinite families. The first is

$$\nu = \zeta_3 \in \pi_{2^j+1-2} S^{2^{j+1}-5}, \quad j \geq 4.$$

---

<sup>3</sup>Since  $\nu = \zeta_3$ , the element  $\nu\zeta_{8k-5}$  could be regarded as an honorary element of the image of  $J$ , failing to attain full membership because it is unstable.

This turns out to be the Hopf invariant of an element

$$\eta_j \in \pi_{2j}^s S^0$$

detected by the element

$$h_1 h_j \in \text{Ext}_A^2(\mathbb{F}_2, \Sigma^{2^j+2}\mathbb{F}_2).$$

This element was constructed by the first author in [19].

The Kervaire invariant elements arose as part of conjectural solution to what happens for the second infinite family. To describe this conjecture we need some notation.

Let  $j \geq 2$  and define integers  $a$  and  $b$  by the equation  $j = 4a + b$  for  $0 \leq b \leq 3$ . Define  $\phi(j) = 8a + 2^b$ . Notice that if  $i \geq 2$ , then  $\pi_i SO \neq 0$  if and only if  $i = \phi(j) - 1$  for some  $j$ . Let  $\beta_j$  be a generator of the image of  $J$  in degree  $\phi(j) - 1$ ; thus, for example, we have

$$\beta_2 = \nu \quad \beta_3 = \sigma \quad \beta_4 = \eta\sigma \quad \beta_5 = \eta^2\sigma \quad \beta_6 = \zeta_{11}.$$

Notice we are excluding the generators  $\mu_{8k+1}$  and  $\eta\mu_{8k+1}$ .

The remaining classes available to be Hopf invariants were the infinite family

$$\beta_j \in \pi_{2n+\phi(j)} S^{2n+1}, \quad n + \phi(j) + 1 = 2^{j+1}.$$

The first author made the following conjecture in 1967 [17].

**Conjecture 1.5.** Let  $n + \phi(j) + 1 = 2^{j+1}$ . The Whitehead product  $[\iota_n, \beta_j] = 0$  if and only if  $h_j^2$  detects a non-zero homotopy class.

To paraphrase the conjecture we have: if  $h_j^2$  detects a non-zero homotopy class  $\Theta_j$ , then  $\Theta_j$  has sphere of origin  $S^{2^{j+1}-\phi(j)}$  and Hopf invariant  $\beta_j$ .

This conjecture has been proved in all aspects by Crabb and Knapp [9], but, of course, the negative solution of [11] leaves only the case  $j = 6$  of interest. Indeed, we now see that for  $j > 6$

$$[\iota_n, \beta_j] \neq 0 \in \pi_* S^{2^{j+1}-\phi(j)-1}.$$

So, somewhat surprisingly, the image of  $J$  has led us into unknown territory. What else can we say about this class?

**Open Problem 1.6.** There is another, richer, question left as well. If the Kervaire invariant class  $\Theta_j$  had existed, it would have had Hopf invariant  $\beta_j$ . A likely consequence of this was that for all odd  $k$ , the element

$$P(\beta_j) \in \pi_* S^{k2^{j+1}-\phi(j)-1}$$

would have had  $\Theta_j$  as its Hopf invariant. Now we have no idea what the Hopf invariant of this family of elements could be, but they are presumably a new and very interesting collection of elements in the stable homotopy groups of spheres. For example, they should play a key role in the iterated root invariant [21] of  $2\iota \in \pi_0^s S^0$ . The elements in this family should depend only on  $m$ , and not on  $k$ .

Here is a table showing the generators in the stable homotopy groups of spheres which are not in the image of  $J$ , yet which have Hopf invariants in the image of  $J$ . Listed also are their spheres of origin, and their Hopf invariants. There are five (or six) sporadic elements and one infinite family. The element  $\Theta_6$  is the unsettled case of the Kervaire invariant problem. It may or may not exist. The element  $\nu^*$  is the Toda bracket  $\langle \sigma, 2\sigma, \nu \rangle$  in  $\pi_{18}^s S^0$ . It is detected by  $h_2 h_4$  in the Adams Spectral Sequence. In this context the  $\eta_j$  family looks quite curious. Why does it have this privileged role?

Element	Sphere of Origin	Hopf Invariant
$\nu^2$	4	$\nu$
$\sigma^2$	8	$\sigma$
$\nu^*$	12	$\sigma$
$\Theta_4$	23	$\eta\sigma = \beta_4$
$\Theta_5$	54	$\eta^2\sigma = \beta_5$
$\Theta_6(?)$	116	$\zeta_{11} = \beta_6$
$\eta_j$	$2^j - 2$	$\nu$

## References

- [1] Adams, J. F., “On the non-existence of elements of Hopf invariant one”, *Ann. of Math. (2)*, 72 (1960), 20–104.
- [2] Adams, J. F., “On the groups  $J(X)$ , IV”, *Topology* 5 (1966), 21-71.
- [3] Barratt M., Jones, J., and Mahowald, M., “Relations amongst Toda brackets and the Kervaire invariant in dimension 62,” *J. London Math. Soc. (2)*, 30 (1984), 533-550.
- [4] Barratt, M. G. and Mahowald, M. E. and Tangora, M. C., “Some differentials in the Adams spectral sequence. II”, *Topology*, 9 (1970), 309–316.
- [5] Browder, W., “The Kervaire invariant of framed manifolds and its generalization”, *Ann. of Math. (2)* 90 (1969), 157–186.
- [6] Browder, W., *Surgery on simply-connected manifolds*, Springer-Verlag, Berlin 1972.
- [7] Brown, Jr., Edgar H., “Generalizations of the Kervaire invariant”, *Ann. of Math. (2)*, 95 (1972), 368–383.
- [8] Cohen, R. L. and Jones, J. D. S. and Mahowald, M. E., “The Kervaire invariant of immersions”, *Invent. Math.*, 79 (1985) no. 1, 95–123.

- [9] Crabb, M. C. and Knapp, K., “Applications of nonconnective  $\text{Im}(J)$ -theory”, *Handbook of algebraic topology*, 463–503, North-Holland, Amsterdam, 1995.
- [10] Devinatz, Ethan S. and Hopkins, Michael J. and Smith, Jeffrey H., “Nilpotence and stable homotopy theory. I”, *Ann. of Math. (2)*, 128 (1988), no. 2, 207–241.
- [11] Hill, M.A, Hopkins, M.J., and Ravenel, D.C, “On the non-existence of elements of Kervaire invariant one,” [arxiv.org/abs/0908.3724](https://arxiv.org/abs/0908.3724).
- [12] Hopkins, M. and Singer, I., “Quadratic functions in geometry, topology, and M-theory,” *J. Differential Geom.* 70 (2005), no. 3, 329–452.
- [13] Hopkins, Michael J. and Smith, Jeffrey H., “Nilpotence and stable homotopy theory. II”, *Ann. of Math. (2)*, 148 (1998), no. 1, 1–49.
- [14] James, I. M., “Reduced product spaces”, *Ann. of Math. (2)*, 62 (1955), 170–197.
- [15] James, I. M., “The suspension triad of a sphere”, *Ann. of Math. (2)*, 63 (1956), 407–429.
- [16] Kervaire, M, “A manifold which does not admit any differential structure,” *Comment. Math. Helv.*, 34 (1960), 256-270.
- [17] Mahowald, M., *The metastable homotopy of  $S^n$* , Memoirs of AMS 72 (1967).
- [18] Mahowald, M, “On the order of the image of  $J$ ”, *Topology* 6 (1967), 371-378.
- [19] Mahowald, Mark, “A new infinite family in  ${}_2\pi_*^s$ ”, *Topology*, 16 (1977), no. 3, 249–256.
- [20] Mahowald, Mark, “The image of  $J$  in the *EHP* sequence”, *Ann. of Math. (2)*, 116 (1982) no.1, 65–112.
- [21] Mahowald, Mark E. and Ravenel, Douglas C., “The root invariant in homotopy theory”, *Topology*, 32 (1993) no. 4, 865–898.
- [22] May, J.P., “The cohomology of restricted Lie algebras and of Hopf algebras; application to the Steenrod algebra”, Thesis, Princeton University, 1964.
- [23] Miller, H. R. and Ravenel, D. C., “Mark Mahowald’s work on the homotopy groups of spheres”, *Algebraic topology (Oaxtepec, 1991)*, Contemp. Math., 146, 1–30, Amer. Math. Soc., Providence, RI 1993.
- [24] Miller, Haynes R. and Ravenel, Douglas C. and Wilson, W. Stephen, “Periodic phenomena in the Adams-Novikov spectral sequence”, *Ann. Math. (2)*, 106 (1977), no. 3, 469–516.

- [25] Pontrjagin, C.R.. *Acad. Sci. URSS* 19 (1938) 147-149, 361-363
- [26] Ravenel, Douglas C., *Complex cobordism and stable homotopy groups of spheres*, Pure and Applied Mathematics, 121, Academic Press Inc., Orlando, FL, 1986.
- [27] Toda, Hirosi, *Composition methods in homotopy groups of spheres*, *Annals of Mathematics Studies*, No. 49, Princeton University Press, Princeton, N.J., 1962.
- [28] Spanier, E.H., *Algebraic Topology*, McGraw-Hill, New York, 1966.
- [29] Wall, C. T. C., *Surgery on compact manifolds*, Mathematical Surveys and Monographs, 69 (2nd ed.), A.M.S., Providence, R.I., 1999.

Department of Mathematics  
Northwestern University  
2033 Sheridan Road  
Evanston, IL, 60208

`mark@math.northwestern.edu`

`pgoerss@math.northwestern.edu`





**CURRENT EVENTS BULLETIN**  
**Previous speakers and titles**

For PDF files of talks, and links to *Bulletin of the AMS* articles, see  
<http://www.ams.org/ams/current-events-bulletin.html>.

**January 7, 2009 (Washington, DC)**

Matthew James Emerton, Northwestern University  
*Topology, representation theory and arithmetic: Three-manifolds and the Langlands program*

Olga Holtz, University of California, Berkeley  
*Compressive sensing: A paradigm shift in signal processing*

Michael Hutchings, University of California, Berkeley  
*From Seiberg-Witten theory to closed orbits of vector fields: Taubes's proof of the Weinstein conjecture*

Frank Sottile, Texas A & M University  
*Frontiers of reality in Schubert calculus*

**January 8, 2008 (San Diego, California)**

Günther Uhlmann, University of Washington  
*Invisibility*

Antonella Grassi, University of Pennsylvania  
*Birational Geometry: Old and New*

Gregory F. Lawler, University of Chicago  
*Conformal Invariance and 2-d Statistical Physics*

Terence C. Tao, University of California, Los Angeles  
*Why are Solitons Stable?*

**January 7, 2007 (New Orleans, Louisiana)**

Robert Ghrist, University of Illinois, Urbana-Champaign  
*Barcodes: The persistent topology of data*

Akshay Venkatesh, Courant Institute, New York University  
*Flows on the space of lattices: work of Einsiedler, Katok and Lindenstrauss*

Izabella Laba, University of British Columbia  
*From harmonic analysis to arithmetic combinatorics*

Barry Mazur, Harvard University  
*The structure of error terms in number theory and an introduction to the Sato-Tate Conjecture*

### **January 14, 2006 (San Antonio, Texas)**

Lauren Ancel Myers, University of Texas at Austin  
*Contact network epidemiology: Bond percolation applied to infectious disease prediction and control*

Kannan Soundararajan, University of Michigan, Ann Arbor  
*Small gaps between prime numbers*

Madhu Sudan, MIT  
*Probabilistically checkable proofs*

Martin Golubitsky, University of Houston  
*Symmetry in neuroscience*

### **January 7, 2005 (Atlanta, Georgia)**

Bryna Kra, Northwestern University  
*The Green-Tao Theorem on primes in arithmetic progression: A dynamical point of view*

Robert McEliece, California Institute of Technology  
*Achieving the Shannon Limit: A progress report*

Dusa McDuff, SUNY at Stony Brook  
*Floer theory and low dimensional topology*

Jerrold Marsden, Shane Ross, California Institute of Technology  
*New methods in celestial mechanics and mission design*

László Lovász, Microsoft Corporation  
*Graph minors and the proof of Wagner's Conjecture*

### **January 9, 2004 (Phoenix, Arizona)**

Margaret H. Wright, Courant Institute of Mathematical Sciences, New York University  
*The interior-point revolution in optimization: History, recent developments and lasting consequences*

Thomas C. Hales, University of Pittsburgh  
*What is motivic integration?*

Andrew Granville, Université de Montréal  
*It is easy to determine whether or not a given integer is prime*

John W. Morgan, Columbia University  
*Perelman's recent work on the classification of 3-manifolds*

### **January 17, 2003 (Baltimore, Maryland)**

Michael J. Hopkins, MIT  
*Homotopy theory of schemes*

Ingrid Daubechies, Princeton University  
*Sublinear algorithms for sparse approximations with excellent odds*

Edward Frenkel, University of California, Berkeley  
*Recent advances in the Langlands Program*

Daniel Tataru, University of California, Berkeley  
*The wave maps equation*



# 2010 CURRENT EVENTS BULLETIN

## Committee

**David Eisenbud**, *University of California, Berkeley, Chair*

**David Aldous**, *University of California, Berkeley*

**Lauren Ancel Meyers**, *University of Texas at Austin*

**Helene Barcelo**, *Mathematics Sciences Research Institute*

**Robert Bryant**, *Mathematics Sciences Research Institute*

**James Demmel**, *University of California, Berkeley*

**Matthew Emerton**, *Northwestern University*

**Susan Friedlander**, *University of Southern California*

**William Fulton**, *University of Michigan*

**Robert Ghrist**, *University of Pennsylvania*

**Andrew Granville**, *University of Montreal*

**Helmut Hofer**, *Courant Institute, New York University*

**Olga Holtz**, *University of California, Berkeley*

**Michael Hutchings**, *University of California, Berkeley*

**Linda Keen**, *Herbert H. Lehman College (CUNY)*

**Carlos Kenig**, *University of Chicago*

**Izabella Laba**, *University of British Columbia*

**László Lovász**, *Eotvos Lorand University*

**John Morgan**, *State University of New York, Stony Brook*

**George Papanicolaou**, *Stanford University*

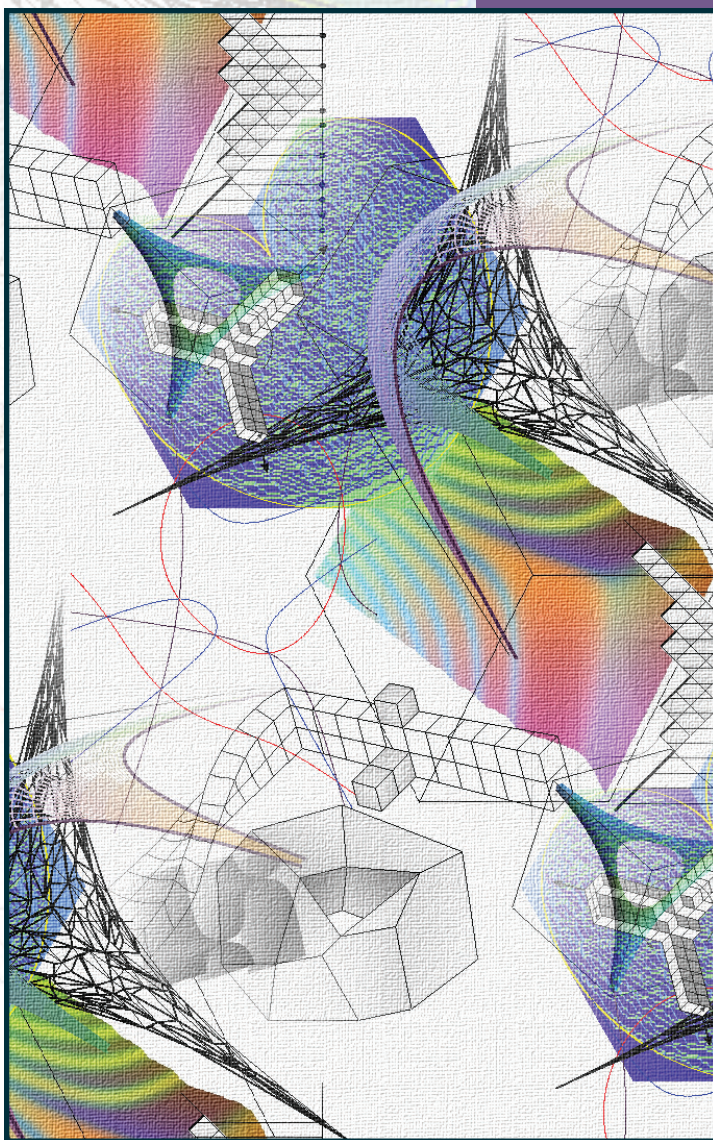
**Richard Schoen**, *Stanford University*

**Frank Sottile**, *Texas A & M University*

**Akshay Venkatesh**, *New York University-Courant Institute*

**Karen Vogtmann**, *Cornell University*

**Andrei Zelevinsky**, *Northeastern University*



---

Cover graphic associated with Green's talk courtesy of Cliff Reiter, originally published in *Vector* 22 No. 4 (2006) 106-116.

Cover graphic associated with Wagner's talk created by David Wagner.

Cover graphic associated with DeMarco's talk courtesy of Curtis T. McMullen, Department of Mathematics, Harvard University.

Cover graphic associated with Hopkins' talk created by Michael Hopkins.

The back cover graphic is reprinted courtesy of Andrei Okounkov.

