

Rational billiards and flat structures

Howard Masur and Serge Tabachnikov

July 3, 2006

Contents

1	Polygonal billiards, rational billiards	3
1.1	Polygonal billiards	3
1.2	Examples: a pair of elastic point-masses on a segment and a triple of point-masses on a circle	4
1.3	Unfolding billiard trajectories, rational polygons	5
1.4	Example: billiard in the unit square	6
1.5	Rational billiard determines a flat surface	8
1.6	Minimality of the billiard flow in rational polygons	10
1.7	Rational billiards and interval exchange maps	13
1.8	Flat metrics and quadratic differentials	14
2	Teichmuller space, strata of quadratic differentials and $SL(2, \mathbf{R})$ actions	18
2.1	Teichmuller space and mapping class group	18
2.2	Compactifications	19
2.3	Strata of quadratic differentials	20
2.4	$SL(2, \mathbf{R})$ action on quadratic differentials	20
3	Ergodicity	22
3.1	Veech nonergodic example	22
3.2	Ergodicity in almost every direction	28
3.3	A combinatorial construction	31
3.4	Further ergodicity results	36
3.5	Ergodicity of general polygonal billiards	37

3.6	Constructive approach to polygons with ergodic billiard flow	41
4	Periodic orbits	41
4.1	Periodic directions are dense	41
4.2	Counting saddle connections and maximal cylinders	48
5	Veech groups and Veech surfaces	49
5.1	Definition and examples of Veech groups	49
5.2	Veech dichotomy	53
5.3	Asymptotics	55
5.4	Covers	56
6	Interval exchange transformations	58
6.1	Topological structure of orbits	58
6.2	Number of invariant measures. Lack of mixing	60
6.3	Ergodicity of interval exchange transformations	64
6.4	Asymptotic flag of an interval exchange transformation	65
7	Miscellaneous results	68
7.1	Stable periodic trajectories	68
7.2	Encoding billiard trajectories. Polygonal billiards have zero entropy .	69
7.3	Complexity of billiard trajectories in rational polygons	73
7.4	Periodic trajectories in some irrational billiards	76
7.5	A non-periodic trajectory that is not dense in the configuration space	78
	References	79

Acknowledgements

Both authors were supported by the National Science Foundation grants. The second author is grateful to the Max-Planck-Institut in Bonn for its hospitality. We are grateful to E. Gutkin, A. Katok, A. Zorich, C. Judge, and Y. Cheung for their criticism and help.

1 Polygonal billiards, rational billiards

1.1 Polygonal billiards

Informally speaking, the theory of mathematical billiards can be partitioned into three areas: convex billiards with smooth boundaries, billiards in polygons (and polyhedra) and dispersing and semi-dispersing billiards (similarly to differential geometry in which the cases of positive, zero and negative curvature are significantly different).

These areas differ by the types of results and the methods of study: in the former a prominent role is played by the KAM theory and the theory of area preserving twist maps; the latter concerns hyperbolic dynamics and has much in common with the study of the geodesic flow on negatively curved manifolds. The recent progress in the study of polygonal billiards is mostly due to applications of the theory of flat structures on surfaces (a.k.a. quadratic differentials) and the study of the action of the Lie group $\mathbf{SL}(2, \mathbf{R})$ on the space of quadratic differentials which is the main topic of this paper. We refer to [39], [40] and [65] for a general survey of mathematical billiards.

We will be considering a plane polygonal billiard Q , that is, a compact polygon (in general, not necessarily convex or simply connected). To fix ideas, the billiard flow is a flow in the unit tangent bundle to Q with discontinuities corresponding to reflections in the boundary ∂Q . These reflections are described by the familiar law of geometrical optics: *the angle of incidence equals the angle of reflection*. The billiard map T is a section of the billiard flow; it acts on unit tangent vectors $(x, v) \in TQ$ whose foot point x is an interior point of a side of Q and the vector v has the inward direction. The map T has an invariant measure μ . Let t be the length parameter along the perimeter of Q and let θ be the angle made by v with the respective side. Then the invariant measure is given by the next formula:

$$\mu = \sin \theta \, d\theta \, dt.$$

One often needs to consider parallel beams of billiard trajectories in polygons. For example, every even-periodic trajectory belongs to a 1-parameter family of parallel periodic trajectories of the same period and length, and these trajectories make a strip. An odd-periodic trajectory also includes into a strip, this time consisting of trajectories whose period and length is twice as great. Note that the billiard reflection in a side of Q transforms a parallel beam to another family of parallel trajectories, and that the width of the beam remains the same.

Fig. 1

Part of the interest in polygonal billiards comes from the fact that they are closely related to problems of mechanics. We discuss one such simple example below.

1.2 Examples: a pair of elastic point-masses on a segment and a triple of point-masses on a circle

Consider two points with masses m_1 and m_2 on the unit segment $[0, 1]$. The points may elastically collide and reflect from the end points of the segment ("walls"). The configuration space of the system is given by the inequalities $0 \leq x_1 \leq x_2 \leq 1$ where x_1 and x_2 are the coordinates of the points. Rescale the variables: $\bar{x}_i = \sqrt{m_i}x_i$, $i = 1, 2$. The configuration space is now a right triangle in the $(\bar{x}_1\bar{x}_2)$ -plane with the acute angle $\tan^{-1}(\sqrt{m_1/m_2})$.

Consider a collision of the points. Let v_1, v_2 be the velocities before, and u_1, u_2 – after the collision. The conservation of momentum and energy laws read:

$$m_1u_1 + m_2u_2 = m_1v_1 + m_2v_2,$$

$$m_1u_1^2/2 + m_2u_2^2/2 = m_1v_1^2/2 + m_2v_2^2/2.$$

In the rescaled coordinates the velocities are rescaled by the same factors, therefore

$$\sqrt{m_1}\bar{u}_1 + \sqrt{m_2}\bar{u}_2 = \sqrt{m_1}\bar{v}_1 + \sqrt{m_2}\bar{v}_2,$$

$$\bar{u}_1^2 + \bar{u}_2^2 = \bar{v}_1^2 + \bar{v}_2^2.$$

The latter equation says that the magnitude of the vector (\bar{u}_1, \bar{u}_2) does not change after the collision; the former one says that the scalar product of this vector with the vector $(\sqrt{m_1}, \sqrt{m_2})$ is preserved as well. The vector $(\sqrt{m_1}, \sqrt{m_2})$ is tangent to the side of the configuration triangle given by the equation $\bar{x}_1/\sqrt{m_1} = \bar{x}_2/\sqrt{m_2}$.

Thus the configuration trajectory reflects in this side according to the billiard reflection law. Likewise one considers collisions with the walls $x = 0$ and $x = 1$: they correspond to the billiard reflections in the other two sides of the configuration triangle. One concludes that the dynamical system of two elastic particles on a segment is isomorphic to the billiard in a right triangle whose shape depends on the ratio of the masses of the particles.

Similar arguments show that the system of three elastic point-masses on the circle with a fixed center of mass is isomorphic to the billiard in an acute triangle (see [24]). Let the masses be m_1, m_2, m_3 ; then the angles α_i of the triangle are given by

$$\tan \alpha_i = m_i \sqrt{\frac{m_1 + m_2 + m_3}{m_1 m_2 m_3}}, \quad i = 1, 2, 3.$$

In the limit $m_3 \rightarrow \infty$ one obtains the previous example of two point-masses on a segment.

Many other mechanical models reduce to billiards, in particular, the model of gas in a closed vessel as a collection of elastic balls in a compact domain. Such models are of great importance in statistical physics.

1.3 Unfolding billiard trajectories, rational polygons

Instead of reflecting a billiard trajectory in a side of the billiard polygon Q one may reflect Q in this side and unfold the trajectory to a straight line. This process is iterated at each reflection: each successive copy of the billiard polygon is obtained from the previous one by the reflection in the side, met by the straightened trajectory. This unfolding method has many applications in the study of polygonal billiards.

Fig. 2

Let $A(Q)$ be the group of motions of the plane generated by the reflections in the sides of Q . Denote the reflection in the side s by σ_s . Notice that for every two sides s and t of Q one has:

$$\sigma_{\sigma_s(t)} = \sigma_s \sigma_t \sigma_s.$$

It follows that every copy of Q involved in the unfolding is the image of Q under an element of the group $A(Q)$. The product of an even number of elements of this group preserves orientation while an odd number reverses it.

To keep track of the directions of billiard trajectories in Q consider the group $G(Q)$ that consists of the linear parts of the motions from $A(Q)$. This subgroup of the orthogonal group is generated by the reflections in the lines through the origin which are parallel to the sides of the polygon Q .

The group $G(Q)$ acts on the unit circle. When a billiard trajectory in Q reflects in a side s its direction is changed by the action of the element of $G(Q)$ which is the projection of σ_s to $G(Q)$.

Definition 1.1 *A billiard polygon Q is called rational if the group $G(Q)$ is finite and irrational otherwise.*

If Q is a rational polygon then the group $G(Q)$ is the dihedral group of symmetries of a regular polygon. A necessary condition for Q to be rational is that all its angles are rational multiples of π . It is also sufficient if the boundary is connected; that is, the polygon is simply connected. The set of plane n -gons can be considered as a subset in \mathbf{R}^{2n} (each vertex has two degrees of freedom). We give the space of n -gons the subspace topology. Rational polygons are dense in the space of polygons.

A given billiard trajectory in a rational billiard table will have only finitely many different directions; this finite collection of directions plays the role of an integral of motion. More precisely, let p be the composition of the projection of the unit tangent bundle $Q \times \mathbf{S}^1$ on \mathbf{S}^1 and the projection of \mathbf{S}^1 to the quotient space $\mathbf{S}^1/G(Q)$ (which is an interval). Then the function p is constant along every billiard trajectory in a rational polygon.

We illustrate the unfolding procedure in the simplest example of a rational polygon, the square.

1.4 Example: billiard in the unit square

Unfolding a trajectory one obtains a line in the plane which is tiled by the unit squares, the images of the original square Q under the action of the group $A(Q)$. Two lines in the plane correspond to the same billiard trajectory in Q if they differ by a translation through a vector from the lattice $2\mathbf{Z} + 2\mathbf{Z}$.

Consider the fundamental domain of the group $2\mathbf{Z} + 2\mathbf{Z}$ which is the square made of four copies of Q . Identify the opposite sides to obtain a flat torus. A billiard trajectory becomes a geodesic line on this torus; every geodesic line has a constant slope λ . The unit tangent bundle of the torus is represented as the union of the tori, parameterized by the slopes λ ; the geodesic flow on each torus is a constant flow.

The dynamics on an individual invariant torus depends on whether λ is rational or irrational: in the former case the geodesic flow is periodic, and in the latter it is ergodic, and in fact, uniquely ergodic. In particular, a billiard trajectory with a rational slope is periodic, while the one with an irrational slope is dense in the square.

One can also analyze periodic trajectories. The unfolding of such a trajectory is a segment in the plane whose end-points differ by a vector from $2\mathbf{Z} + 2\mathbf{Z}$. Every parallel trajectory is also periodic with the same period and length.

Assume that an unfolded periodic trajectory goes from the origin to the lattice point $(2p, 2q)$. If p and q are coprime this is a prime periodic trajectory, and if p and q have a common multiplier then the periodic trajectory is multiple. The length of the trajectory is $2\sqrt{p^2 + q^2}$, and to a choice of p and q there correspond two orientations of the trajectory. Thus the number of (strips of) parallel trajectories of length less than L equals the number of pairs of integers, satisfying the inequality $p^2 + q^2 < L^2/2$.

This is the number of lattice points inside the circle of radius $L/\sqrt{2}$, centered at the origin. In the first approximation, this number equals the area $\pi L^2/2$, and to take only prime periodic trajectories into account (that is, only coprime (p, q)) one divides this by $\pi^2/6$. Thus one obtains a quadratic asymptotic estimate on the number of periodic trajectories of the length not exceeding a fixed number. We will see in a later section that this result holds for general rational polygons.

Integrable billiards. What works so well in the above example is the fact that the images of Q under the group $A(Q)$ tile the plane. A similar consideration applies to rectangles, equilateral triangles and right triangles with an acute angle $\pi/4$ or $\pi/6$. These polygons are called *integrable*, and the billiard flow reduces to a constant flow on a torus. Note that in all these cases one may define the extension of a billiard trajectory through a vertex of the billiard polygon.

Although integrable polygons are exceptional, some of the features of the billiard dynamics in the integrable case extend to general rational polygons. These results will be discussed in succeeding sections.

Almost integrable billiards. One class of rational billiards for which one can make such precise statements is the class of *almost integrable* billiards, intermediate between integrable and general rational polygons, studied by Gutkin ([27]). A polygon Q is called almost integrable if the group $A(Q)$ is a discrete subgroup of the group of motions of the plane. There are exactly four such groups generated by the reflections in the sides of the four integrable polygons. An almost integrable polygon can be drawn on the corresponding lattice.

Fig. 3

Given an almost integrable polygon, the billiard flow decomposes into directional flows F_θ (just as in the case of a square – see next section for a detailed discussion). Choose a basis e_1, e_2 of the respective lattice. A direction is called rational if it is given by a vector $a_1 e_1 + a_2 e_2$ with $a_1/a_2 \in \mathbf{Q}$. Gutkin proved (see also [9]) that, similarly to the square case, the following conditions are equivalent:

- (i) θ is an irrational direction;

- (ii) F_θ is minimal;
- (iii) F_θ is ergodic;
- (iv) F_θ is aperiodic, that is, $F_\theta^t \neq id$ for all $t \neq 0$.

1.5 Rational billiard determines a flat surface

The following construction of a flat surface from a rational billiard table plays a central role in the present study (see [19], [41], [42], [61] and [47]; the latter paper was the first to relate billiards in polygons and quadratic differentials).

Let Q be a rational polygon. The group $G(Q)$ is the dihedral group D_N generated by the reflections in the lines through the origin that meet at angles π/N where N is a positive integer. This group has $2N$ elements, and the orbit of a generic point $\theta \neq k\pi/N$ on the unit circle consists of $2N$ points. Let the angles of Q be $\pi m_i/n_i$ where m_i and n_i are coprime integers. If Q is simply connected (the assumption we make throughout this section) then N is the least common multiple of the denominators n_i .

Consider the unit tangent bundle $Q \times \mathbf{S}^1$, the phase space of the billiard flow, and let M_θ be the subset of points whose projection to \mathbf{S}^1 belongs to the orbit of θ under D_N . Then, M_θ is an invariant surface of the billiard flow in Q . This invariant surface is a level surface of the above mentioned function p , "the integral of motion".

Assume that $\theta \neq k\pi/N$ and enumerate the angles in the D_N -orbit of θ on the unit circle counterclockwise: $\theta = \theta_1, \theta_2, \dots, \theta_{2N}$. The surface M_θ is obtained from $2N$ copies of Q , namely, $Q \times \theta_i \subset Q \times \mathbf{S}^1$, $i = 1, \dots, 2N$ by gluing their sides according to the action of D_N .

Consider $2N$ disjoint and parallel copies of Q in the plane. Call them Q_1, \dots, Q_{2N} and orient the even ones clockwise and the odd ones counterclockwise. Choose an index $i = 1, \dots, 2N$ and a side s of Q_i ; reflect the direction θ_i in this side. The resulting direction is θ_j for some $j = 1, \dots, 2N$. Glue the side s of Q_i to the identical side of Q_j . After these gluings are made for all values of i and all choices of the side s of Q_i , the sides of all the polygons Q_i are pasted pairwise, and the gluings agrees with the orientation. The result is an oriented compact surface that depends only on the polygon Q , but not on the choice of θ , and we denote it by M . The directional billiard flows F_θ on M in directions θ are obtained one from another by rotations.

For example, if Q is a square then $N = 2$ and the result is a torus made of four identical squares. If Q is a right triangle with an acute angle equal to $\pi/8$ then the surface M is obtained from a regular octagon, the result of gluing 16 copies of the triangle, by pairwise gluing its opposite sides; this surface has genus 2.

Fig. 4

The genus of M is given in the next lemma.

Lemma 1.2 *Let the angles of a billiard k -gon be $\pi m_i/n_i$, $i = 1, \dots, k$ where m_i and n_i are coprime, and N be the least common multiple of n_i 's. Then*

$$\text{genus } M = 1 + \frac{N}{2} \left(k - 2 - \sum \frac{1}{n_i} \right).$$

Proof. We need to analyze how the pastings are made around a vertex of Q . Consider the i -th vertex V with the angle $\pi m_i/n_i$. Let G_i be the group of linear transformations of the plane generated by the reflections in the sides of Q , adjacent to V . Then G_i consists of $2n_i$ elements.

According to the construction of M the number of copies of the polygons Q_j that are glued together at V equals the cardinality of the orbit of the test angle θ under the group G_i , that is, equals $2n_i$. Originally we had $2N$ copies of the polygon Q , and therefore, $2N$ copies of the vertex V ; after the gluings we have N/n_i copies of this vertex on the surface M .

It follows that the total number of vertices in M is $N(\sum 1/n_i)$. The total number of edges is Nk , and the number of faces is $2N$. Therefore the Euler characteristic of M equals

$$N \sum \frac{1}{n_i} - Nk + 2N = 2 - 2g$$

where g is the genus, and the result follows. \square

The billiard flow on M is obtained from the constant flows in the directions θ_i in the polygons Q_i . The result is a (unit) vector field on M with singularities at the vertices. The above proof shows that the i -th vertex of M is the result of gluing $2n_i$ copies of the angle $\pi m_i/n_i$ which sums up to an angle of $2\pi m_i$. One may realize such a singularity geometrically as follows. Take m_i copies of a Euclidean upper half plane H_j and m_i copies of a lower half plane L_j ; $j = 1, \dots, m_i$. Then glue the positive real axis of H_j to the positive real axis of L_j and glue the negative real axis of L_j to the negative real axis of H_{j+1} ($m_i + 1 = 1$). The result is a singularity with a total angle of $2\pi m_i$. We will call it a cone angle $2\pi m_i$ singularity. From this description one can also see that for any direction θ , there are $2m_i$ separatrices (m_i of them incoming and m_i outgoing) emanating from the singularity in direction θ .

It is easy to describe the set of polygons for which all the angles are 2π , that is, $m_i = 1$; then the singularities of M are removable. The sum of interior angles of a k -gon is $\pi(k - 2)$. Thus if $m_i = 1$ for all i then

$$\frac{1}{n_1} + \dots + \frac{1}{n_k} = k - 2.$$

This equation has only four solutions with $n_i \geq 2$, considered up to permutations:

$$\left(\frac{1}{3}, \frac{1}{3}, \frac{1}{3}\right), \left(\frac{1}{2}, \frac{1}{4}, \frac{1}{4}\right), \left(\frac{1}{2}, \frac{1}{3}, \frac{1}{6}\right), \left(\frac{1}{2}, \frac{1}{2}, \frac{1}{2}, \frac{1}{2}\right).$$

These solutions correspond to the already mentioned integrable polygons, and in each case the invariant surface M is a torus with a constant flow.

To summarize the construction, given a rational billiard polygon Q , one constructs a compact surface M whose genus is given by the above lemma. This surface inherits a flat metric from Q with a finite number of cone-type singularities, corresponding to the vertices of Q , with cone angles multiples of 2π . The directional billiard flow F_θ on M is a constant flow in a fixed direction with singularities at the cone points.

1.6 Minimality of the billiard flow in rational polygons

The next result on the minimality of the billiard flow in a rational polygon is much easier than the stronger results on ergodicity described later in the paper. We discuss it because it serves a model for these harder theorems.

Definition 1.3 *A flow is called minimal if any of its orbits is dense in the phase space.*

Definition 1.4 *A directed saddle connection in direction θ is an orbit of F_θ that goes from a singularity to a singularity (possibly, the same one) and has no interior singularities.*

A saddle connection is also called a generalized diagonal since it corresponds to a billiard orbit in direction θ that goes from a vertex to a vertex.

Since the saddle connection is represented by a geodesic in the Euclidean structure it determines a vector in \mathbf{R}^2 . For a saddle connection β we will denote by h_β and v_β the horizontal and vertical components of this vector.

Definition 1.5 *A metric cylinder in direction θ is the isometric image of an open right angled flat Euclidean cylinder consisting of closed geodesics in direction θ . A metric cylinder is maximal if it cannot be enlarged.*

We already noted that a periodic orbit in direction θ for the billiard flow included into a strip of parallel periodic orbits thus giving rise to a metric cylinder in direction θ . The following lemmas and theorem are standard (see [64]), and the proofs do not use the group structure coming from billiards. Thus they will hold verbatim for the general situation of flat structures (that will be formally defined shortly).

Lemma 1.6 *If M has singularities, then the boundary of a maximal metric cylinder consists of a finite number of saddle connections in direction θ .*

Proof. If a metric cylinder fills the surface, then M is the flat torus which has no singularities. Thus we can assume the maximal cylinder has a boundary. The obstruction to enlarging a maximal cylinder are saddle connections on the boundary. \square

Lemma 1.7 *Suppose α^+ is a trajectory, infinite in the positive direction, and starting at point P_0 . Let β be an interval perpendicular to α^+ with P_0 as one of its endpoints. Then α^+ returns to β .*

Proof. Since there are a finite number of singularities, there are a finite number of trajectories starting at points of β that hit a singularity before crossing β again. By shortening β to a subinterval β' with one endpoint P_0 and some other endpoint Q_0 , we can assume that no trajectory leaving β' hits a singularity before returning to β . Now flow the interval β' in the positive direction. The interval sweeps out rectangles of increasing area. Since the area of the surface is finite, the interval β' must return and overlap β . If α^+ itself does, we are done. Otherwise the trajectory leaving Q_0 returns to β' and some trajectory leaving a point $Q_1 \in \beta'$ returns to P_0 . We now consider the interval $\beta'' \subset \beta'$ with endpoints P_0 and Q_1 and apply the previous analysis to it. Flowing in the forward direction it must return to β and now α^+ itself must return to β . \square

Theorem 1.8 *For all but countably many directions θ the flow F_θ is minimal on the surface M .*

Proof. Since there are only finitely many singularities and countably many homotopy classes of arcs with fixed endpoints joining the singularities, it is clear that there are only countably many saddle connections (the same statement follows from the fact that the group $A(Q)$ is countable.) We claim that the flow F_θ is minimal if there is no saddle connection in direction θ .

By Lemma 1.6 there cannot be a metric cylinder for then there would be a saddle connection. Suppose there was an infinite trajectory l_θ in direction θ which is not dense. Let $A \neq M$ be the set of ω limit points of l_θ . Then A is invariant under the flow F_θ . Since $A \neq M$ one can choose a trajectory γ on ∂A ; let P_0 be its initial point. We will show that γ is a saddle connection. If not then γ is infinite in at least one of the two directions. We show that there is an open neighborhood of P_0 contained in A , a contradiction to P_0 being a boundary point.

Let β be a perpendicular arc with P_0 an endpoint. It is enough to show that there exists a segment $[P_0, Q] \subset \beta$ which is contained in A . For then doing this on both sides we would have our open neighborhood. Now Lemma 1.7 implies that γ hits β again at some P_1 . If the interval $[P_0, P_1] \subset A$ we are done. Suppose not. Then there exists $Q_1 \in [P_0, P_1]$ which is not in A . Since A is closed, there is a largest open subinterval $I_1 \subset [P_0, P_1]$ containing Q_1 which is in the complement of A . Let P_2 be the endpoint of I_1 closest to P_0 . Then $P_2 \in A$ and the trajectory through P_2 must be a saddle connection. For if it were infinite in either direction, it would intersect I_1 . Since A is invariant under the flow, this contradicts that I_1 misses A . \square

Returning back to the rational billiard polygon Q we see that for all but countably many directions every billiard trajectory is dense in Q . Approximating a general billiard by a rational one we obtain the next result on topological transitivity of polygonal billiards ([41]).

Definition 1.9 *A flow is called topologically transitive if it has a dense orbit.*

Theorem 1.10 *For every $k \geq 3$ there exists a dense G_δ subset in the space of simply connected k -gons that consists of polygons with topologically transitive billiard flows.*

Proof. Identify the phase space of the billiard flow in each k -gon with $\mathbf{D}^2 \times \mathbf{S}^1$, and assume that this identification depends continuously on the polygon. Let B_i be a countable basis for the topology of $\mathbf{D}^2 \times \mathbf{S}^1$. Denote by X_n the set of k -gons Q such that for each open set U in the phase space, there exists a billiard trajectory starting in U that visits all the images of the sets B_1, \dots, B_n in the phase space of the billiard flow in Q . Each set X_n is open and their intersection is a G_δ set.

Let us show that this intersection is dense. Let Y_m be the set of rational k -gons with the angles $\pi m_i/n_i$, m_i and n_i coprime, such that the least common multiple of n_i 's is at least m . For every polygon $Q \in Y_m$ the invariant surface M_θ is $1/m$ -dense in the phase space. Therefore for every n there exists m such that for every $Q \in Y_m$ the surface M_θ intersects all the images of the sets B_1, \dots, B_n in the phase space of the billiard flow in Q . Since M_θ has a dense trajectory for all but countably many θ the space Y_m lies in X_n . Finally, Y_m is dense in the space of k -gons for every m , and so X_n is dense for every n as well. It follows that $\cap X_n$ is dense by Baire's theorem.

Let Q be a billiard polygon in $\cap X_n$. We claim that there is a dense billiard trajectory in the phase space of Q . Let U be a compact domain in the phase space $\mathbf{D}^2 \times \mathbf{S}^1$. Suppose inductively we have chosen a compact neighborhood $U_{n-1} \subset U$. Since $Q \in \cap X_n$, one can find a billiard trajectory starting in U_{n-1} in Q that visits B_1, \dots, B_n . By continuity there is a neighborhood $U_n \subset U_{n-1}$ such that each trajectory in U_n visits B_1, \dots, B_n within time T_n . Choose a phase point $v \in \cap U_n$. Then the trajectory of the vector v is dense in the phase space of Q . \square

1.7 Rational billiards and interval exchange maps

The billiard flow on the invariant surface can be reduced to a one-dimensional transformation; this reduction is a particular case of the reduction of the billiard flow to the billiard transformation.

Definition 1.11 *Let (I_1, \dots, I_n) be a partition of the interval $[0, 1)$ into nonintersecting semiclosed intervals, enumerated from left to right, and let σ be a permutation of n elements. The respective interval exchange transformation $T : [0, 1) \rightarrow [0, 1)$ is a transformation whose restriction to every I_i is a parallel translation and such that the intervals $T(I_1), \dots, T(I_n)$ follow from left to right in the order $\sigma(1), \dots, \sigma(n)$.*

Clearly, an interval exchange transformation preserves the Lebesgue measure on the unit interval.

Example. The exchange of two intervals $[0, a)$ and $[a, 1)$ is identified with the rotation of the circle \mathbf{R}/\mathbf{Z} through $1 - a$.

Remark. We will often need to consider an interval exchange transformation defined on an interval I , different from $[0, 1)$. In such a case it is understood that I is rescaled to $[0, 1)$ by an affine transformation.

The reduction of the billiard flow in a fixed direction θ to an interval exchange goes as follows. Recall that the invariant surface M was constructed by pasting pairs of

equal and parallel sides of $2N$ copies of the billiard polygon Q . Choose one side from each such pair and call these segments s_1, \dots, s_m . Arrange these segments along a line and let I be their union. Give each segment the measure determined by the orthogonal projection on the direction, perpendicular to θ . Since "the width of a beam" is an invariant transversal measure of a constant flow, the billiard transformation induces a piecewise isometry T of the segment I . If the segments are appropriately oriented then T is orientation preserving. It remains to rescale I to the unit interval and to define T at its points of discontinuity so that it is continuous on the left. The result is an interval exchange transformation.

The reduction of the billiard flow in a fixed direction to an interval exchange transformation is by no means unique. For example, assume that the flow F_θ is topologically transitive. Choose an interval I in M , transverse to the flow, and give it the measure determined by the orthogonal projection on the direction, perpendicular to θ . Then the first return map to I along the trajectories of F_θ is orientation and measure preserving, that is, an interval exchange transformation. Note that since F_θ is topologically transitive, the first return map is defined no matter how small I is. The points of discontinuity of the first return map correspond to the orbits of the flow through the singular points of the flow F_θ .

1.8 Flat metrics and quadratic differentials

As we have seen a rational billiard defines a closed surface formed by gluing Euclidean polygons isometrically along their edges. The vertices of the polygon correspond to points with cone angle singularities of the metric. In this section we generalize this notion to what we will call flat structures with parallel line fields, or just flat structures for short, so that the set of rational billiards of a given genus is a subset of the space of flat structures. The main reason for this generalization is that the group $SL(2, \mathbf{R})$ acts on the space of flat structures while not preserving the space of rational billiards. The study of this action allows to make conclusions about flat structures in general, which then yield results about rational billiards in particular.

Let M be a compact C^∞ surface and Σ a finite set of points in M . On $M \setminus \Sigma$ we require coordinate charts $v = (x, y)$ such that the transition functions on the overlaps are of the form

$$v \rightarrow v + c \quad \text{or} \quad v \rightarrow -v + c.$$

That is, the transition functions are translations and reflections in the origin followed by translations. Since they preserve the Euclidean metric, this allows one to define

a locally Euclidean metric on M . In addition, these transition functions preserve families of parallel lines in the plane so that for each direction θ , there is a well-defined foliation F_θ of the surface consisting of lines in direction θ .

Note that if the transition functions are all translations, then the line field defines a vector field and in each direction θ we have a flow F_θ . This is exactly the situation for flat structures defined by rational billiards.

Theorem 1.8 that holds for rational billiards still holds in the more general case of a flat structure. Moreover we still have the notions of saddle connections and metric cylinders. The form of the transition functions implies that in a neighborhood of a point in Σ there are polar coordinates (r, θ) such that the metric can be written as

$$ds^2 = dr^2 + (crd\theta)^2$$

where c is a half integer. We say that the metric has cone type singularity with cone angle $2\pi c$. The curvature κ at a singular point is defined by the formula

$$\kappa = 2\pi - 2\pi c.$$

We may concretely describe the metric in a neighborhood of a singular point by gluing together half planes as described above. Notice that our previous discussion shows that a rational polygon determines such a flat structure with transition functions all of the first type (parallel translations) and c an integer.

One may visualize a flat structure as a finite union of polygons glued along parallel (or antiparallel) sides, such that each side is glued to exactly one other, and with special vertices v_1, \dots, v_n such that at v_i there is a conical singularity of excess angle πk_i .

Another way to define the same structure would be to begin by requiring that M have Gaussian curvature zero, or equivalently, it is locally isometric to \mathbf{R}^2 with the Euclidean metric, away from a finite set of points Σ and cone type singularities near points of Σ , where c is an arbitrary real. This in itself is not enough, for it does not rule out the possibility of rotations in the transition functions, which would not allow for parallel line fields. One way to guarantee the correct structure in this language is to consider parallel translation with respect to the connection determined by the metric. The global obstruction to parallel translation being well defined is the holonomy group which is a subgroup of $O(2)$. If the holonomy group is either trivial or $\{I, -I\}$ then parallel translation is well-defined and M possesses a parallel line field. The existence of such a line field implies that the cone angles are multiples of π . Pick one such line field and call it the vertical line field. In the neighborhood

of each point of $M \setminus \Sigma$ we can construct a chart which is an orientation preserving isometry and takes the vertical line field to the vertical line field in \mathbf{R}^2 . The change of coordinate functions between these charts then have the desired form.

Quadratic differentials. The same structure can also be defined in terms of complex analysis. Recall that a Riemann surface or a complex analytic structure on a surface M consists of an atlas of charts (U_α, z_α) where $\{U_\alpha\}$ is a covering of M by open sets, $z_\alpha : U_\alpha \rightarrow \mathbf{C}$ is a homeomorphism and if $U_\alpha \cap U_\beta \neq \emptyset$, then $z_\alpha \circ z_\beta^{-1} : z_\beta(U_\alpha \cap U_\beta) \rightarrow \mathbf{C}$ is complex analytic. The maps z_α are called uniformizing parameters.

A meromorphic quadratic differential ϕ assigns to each uniformizing parameter z_α a meromorphic function $\phi_{z_\alpha}(z_\alpha)$ with the property that

$$\phi_{z_\beta}(z_\beta) \left(\frac{dz_\beta}{dz_\alpha} \right)^2 = \phi_{z_\alpha}(z_\alpha)$$

in $U_\alpha \cap U_\beta$. Then the quadratic differential $\phi(z)dz^2$ is invariantly defined on M .

Although the value of ϕ is not well-defined, the set of zeroes and poles of ϕ and their orders are. It is a classic result in Riemann surface theory that ϕ has a finite number of zeroes with orders k_i and poles of order l_i satisfying

$$\sum k_i - \sum l_i = 4g - 4.$$

In this paper we will assume that if there are poles then they are simple. This implies that the norm or area of ϕ defined by

$$\|\phi\| = \int_M |\phi(z)| |dz|^2$$

is finite. Notice that the area element $|\phi(z)| |dz|^2$ is well-defined independently of choice of coordinates. A quadratic differential determines a metric with the length element $|\phi(z)|^{1/2} |dz|$. If the quadratic differential has a simple pole then the metric is not complete on the punctured surface.

The vertical trajectories of ϕ are the arcs along which

$$\phi(z)dz^2 < 0$$

and the horizontal trajectories the arcs along which

$$\phi(z)dz^2 > 0.$$

Equivalence of definitions. We now indicate why a flat structure (with a parallel line field) and holomorphic quadratic differential define the same object. First suppose ϕ is a quadratic differential on M and $p \in M$ is a point which is neither a pole nor zero. We may choose a uniformizing parameter ζ in a neighborhood of p with p corresponding to $\zeta = 0$. We may then choose a branch of $\phi^{1/2}(\zeta)$ near $\zeta = 0$ and define a new uniformizing parameter w by

$$w = w(\zeta) = \int_0^\zeta \phi^{1/2}(\tau) d\tau.$$

Then, in the w coordinates, the quadratic differential is given by $\phi_w(w) \equiv 1$. Such coordinates are called natural. If w and w' define natural coordinates in overlapping neighborhoods, then

$$w' = \pm w + c$$

and so one has a flat structure. If the holomorphic quadratic differential is the square of an Abelian differential, then the transition functions are translations. Again we remark that this is the case with rational billiards. The Riemannian metric is the flat metric $|dw|$.

At a zero of ϕ of order k (at a simple pole take $k = -1$) the quadratic differential can be written as

$$\left(\frac{k+2}{2}\right)^2 \zeta^k d\zeta^2$$

for a choice of coordinate ζ . Then

$$w(\zeta) = \zeta^{\frac{k+2}{2}},$$

where the uniformizing parameter w is defined as above. The full angle $0 \leq \arg \zeta \leq 2\pi$ is subdivided into $k+2$ equal sectors

$$\frac{2\pi}{k+2}j \leq \arg \zeta \leq \frac{2\pi}{k+2}(j+1), \quad j = 0, \dots, k+1.$$

Every sector is mapped by $w(\zeta)$ onto the upper or lower halfplane. Thus a zero of ϕ of order k corresponds to a singular point of a flat structure with the cone angle $\pi(k+2)$. The pre-images of the horizontal lines are the horizontal trajectories of the quadratic differential ϕ , and the trajectory structure has the form of a $k+2$ *pronged* singularity.

Fig. 5

Conversely, a flat structure with a parallel line field defines a quadratic differential. Let $v = (x, y)$ and $v' = (x', y')$ be coordinates in overlapping charts so that

$$v' = \pm v + c.$$

Setting $z = x + iy$ and $z' = x' + iy'$, the transition functions are complex analytic, so define a Riemann surface structure on $M \setminus \Sigma$, where Σ is the set of singularities. We may then define a quadratic differential ϕ on this Riemann surface by assigning the constant function 1 to the parameter z . The $2\pi c$ cone angle singularity corresponds to a zero of order $2c - 2$ (recall that c is a half integer).

There is an additional description of these structures via the theory of measured foliations, developed by Thurston (see [66], [7], [17]). This description will not play any role here and we will not dwell on it.

We will go back and forth between the terminology of quadratic differentials and flat structures with parallel line fields.

Example. We finish the section with an example taken from [78]: the Riemann surface corresponding to the billiard in a regular n -gon with n odd is conformally equivalent to the Fermat curve $x^n + y^n = 1$, and the respective quadratic differential is dx^2/y^4 ; see [3] for a similar description of Riemann surfaces corresponding to rational triangles.

2 Teichmuller space, strata of quadratic differentials and $SL(2, \mathbf{R})$ actions

2.1 Teichmuller space and mapping class group

A general reference for Teichmuller spaces, and compactifications is the paper of Bers [6]. Let M be a surface of genus g with n punctures. We assume $3g - 3 + n \geq 0$.

Definition 2.1 *Teichmuller space $T_{g,n}$ is the space of equivalence classes of complex structures or Riemann surface structures X on M where $X_1 \sim X_2$ if there is a biholomorphic map from X_1 to X_2 which is isotopic to the identity on M .*

One may define the Teichmuller distance function $d_T(\cdot, \cdot)$ on T_g by

$$d_T(X_1, X_2) = 1/2 \inf \log K(f),$$

where the infimum is taken over all quasiconformal maps f isotopic to the identity on M and $K(f)$ is the maximal dilation of f as measured by the complex structures X_1, X_2 .

In the special case of $g = 1, n = 0$; that is, elliptic curves, $T_{1,0}$ is well known to be the upper half plane and the Teichmüller metric the Poincaré or non-euclidean metric.

Definition 2.2 *The mapping class group $Mod(g, n)$ is the group $Diff^+(M)/Diff_0(M)$; the group of orientation preserving diffeomorphisms modulo those isotopic to the identity.*

The mapping class group acts on $T_{g,n}$ by pull-back; given a complex structure X defined by coordinate charts (U_α, z_α) and $f \in Mod(g)$, we find a new complex structure $f \cdot X$ defined by the atlas of coordinate charts $(f(U_\alpha), z_\alpha \circ f^{-1})$. The quotient space $R_{g,n} = T_{g,n}/Mod(g)$ is the moduli space of Riemann surfaces of genus g with n punctures. It is well-known that $Mod(1, 0)$ is the group $SL(2, Z)$.

2.2 Compactifications

The moduli space $R_{g,n}$ is well-known not to be compact. It is possible to deform a Riemann surface by pinching along one or more disjoint simple closed curves, by letting the hyperbolic length of the curves go to 0. The resulting surface then has nodes or punctures and may not be connected. For example if a closed surface of genus g is pinched along a single closed curve γ and γ does not disconnect M , the resulting surface has genus $g - 1$ with 2 punctures, while if γ is dividing, the resulting surface has two components each of which has one puncture and the sum of their genera is g . We can compactify $R_{g,n}$ by adjoining the moduli spaces of surfaces obtained in this fashion.

Denote the compactification by $\bar{R}_{g,n}$. The topology has the following property. Let X_0 be any surface in the compactification and $X_n \rightarrow X_0$. Remove any neighborhood U of the punctures of X_0 . Then for large enough n , there is a conformal embedding of $X_0 \setminus U \rightarrow X_n$.

The compactification by these moduli spaces is well-behaved with respect to quadratic differentials. For suppose ϕ_n are unit norm quadratic differentials on X_n which converge to $X_0 \in \bar{R}_{g,n}$. Then there is a subsequence of ϕ_n which converges uniformly on compact sets of X_0 via the conformal embedding to an integrable quadratic differential ϕ_0 . However it may be the case that $\phi_0 \equiv 0$ on one or more components of X_0 . This issue will be discussed further in the section on periodic orbits.

2.3 Strata of quadratic differentials

A meromorphic quadratic differential ϕ on a surface of genus g with n punctures with at most simple poles at the punctures defines for some j , a $j + 1$ tuple $\sigma(\phi) = (k_1, \dots, k_j, \epsilon)$, where k_i are the orders of the zeroes and poles and $\epsilon = +$ if the quadratic differential is the square of an Abelian differential and $-$ if it is not.

We say that two quadratic differentials defining the same σ are equivalent if there is a homeomorphism of the surface, isotopic to the identity taking singular points to singular points of the same order and which at other points has the same local form as the change-of-coordinate functions. The set of equivalence classes is denoted by $Q(\sigma)$ and is called a stratum. The resulting space is denoted by $Q(\sigma)$ and one can prove that it is a manifold. In the case of a compact surface ($n = 0$) and all $k_i = 1$ (and $\epsilon = -1$) $Q(\sigma)$ is called the principle stratum for compact Riemann surfaces. For fixed g, n , the union of the $Q(\sigma)$ fit together to form the space $Q_{g,n}$ of all meromorphic quadratic differentials with at most simple poles over $T_{g,n}$. It is a well-known part of Teichmuller theory that $T_{g,n}$ is a complex manifold and this space $Q_{g,n}$ is the cotangent bundle, although we will not make use of this structure in this paper.

A stratum need not be connected – [70] and [49]. See also [57] that give necessary and sufficient conditions on a given $j + 1$ tuple σ so that there exists a differential ϕ such that $\sigma = \sigma(\phi)$.

If we consider compact Riemann surfaces ($n = 0$), then the strata $Q(\sigma)$ are not closed subsets of $Q_{g,0}$, unless there is a single zero ($j = 1$). This is because a sequence in $Q(\sigma)$ may collapse a pair of lower order zeroes into a higher order one.

The group $Mod(g, n)$ acts on each $Q(\sigma)$ by pull-back of structures. If $\{u_\alpha\}$ is an atlas of natural coordinates defining a quadratic differential ϕ , and f is a diffeomorphism, then $\{u_\alpha \circ f^{-1}\}$ defines a new family of natural coordinates. The quotient is denoted by $QD(\sigma)$ and will play a crucial role in this paper. Since $Mod(g, n)$ does not act freely, the space $QD(\sigma)$ has the structure of an orbifold.

2.4 $SL(2, \mathbf{R})$ action on quadratic differentials

The group $SL(2, \mathbf{R})$ acts on each stratum by linear transformations of local coordinates. If $\{u_\alpha\}$ is an atlas of natural coordinates defining a quadratic differential ϕ , and $A \in SL(2, \mathbf{R})$ then $\{Au_\alpha\}$ defines a new family of natural coordinates for a quadratic differential $A\phi$. It follows that A takes singularities of ϕ to singularities of $A\phi$ of the same order.

If one visualizes a flat structure S as a union of polygons $R_1 \cup \dots \cup R_n$, then $gS = gR_1 \cup \dots \cup gR_n$.

The following one parameter subgroups of $SL(2, \mathbf{R})$ are of special interest:

$$g_t = \begin{pmatrix} e^{t/2} & 0 \\ 0 & e^{-t/2} \end{pmatrix}, \quad r_\theta = \begin{pmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{pmatrix}, \quad h_s = \begin{pmatrix} 1 & s \\ 0 & 1 \end{pmatrix},$$

referred to as the geodesic or Teichmüller, circular and horocyclic flows, respectively.

The effect of the flow $\phi \rightarrow g_t \phi$ is to stretch along the horizontal trajectories of ϕ by a factor of e^t and contract along the vertical trajectories by e^t . Associated to g_t is the map f_t of the underlying Riemann surfaces $X \rightarrow X_t$, called the Teichmüller map. The map f_t takes zeroes of ϕ to zeroes of ϕ_t . If $w = u + iv$ are natural coordinates for ϕ away from the zeroes, and $w_t = u_t + iv_t$ are natural coordinates for ϕ_t away from its zeroes, then in these coordinates, f_t is given by

$$u_t = e^{t/2}u, \quad v_t = e^{-t/2}v.$$

The famous theorem of Teichmüller asserts that *given any homeomorphism $f : X \rightarrow Y$ of Riemann surfaces of finite type, there is a unique Teichmüller map from X to Y which minimizes the maximal dilation among all quasiconformal homeomorphisms in the homotopy class*

The action of r_θ on a quadratic differential ϕ is the same as multiplying ϕ by $e^{2i\theta}$. This multiplication defines an action of the circle $\mathbf{R}/\pi\mathbf{R}$ on the set of quadratic differentials. Note that the action of r_θ leaves the flat metric invariant, but changes the vertical line field.

The action on saddle connections is as follows. Let

$$A = \begin{pmatrix} a & b \\ c & d \end{pmatrix}$$

be a matrix in $SL(2, \mathbf{R})$. Let β be a directed saddle connection, and let h_1 and v_1 be the horizontal and vertical components with respect to ϕ of the associated vector. Then the horizontal and vertical components h_2 and v_2 of the vector associated to β with respect to $A\phi$ are given by matrix multiplication; namely

$$\begin{pmatrix} a & b \\ c & d \end{pmatrix} \cdot \begin{pmatrix} h_1 \\ v_1 \end{pmatrix} = \begin{pmatrix} h_2 \\ v_2 \end{pmatrix}.$$

The orbit of a quadratic differential under $SL(2, \mathbf{R})$ is the unit tangent bundle to an isometrically embedded hyperbolic disc inside Teichmüller space. This disc is

called a Teichmüller disc and has been considered by numerous authors – see [50], for example. In a later section we will examine certain examples of Teichmüller discs that arise in the so-called Veech billiards and their generalizations.

It is immediate from the definitions that the action of $SL(2, \mathbf{R})$ commutes with the action of $Mod(g)$ and thus descends to an action on $QD(\sigma)$. It is this action we will study in some detail.

Invariant measure on the strata. Let $QD^1(\sigma) \subset QD(\sigma)$ consist of the flat surfaces with unit area. It is possible to define an $SL(2, \mathbf{R})$ invariant measure μ_0 on each $QD^1(\sigma)$ which is absolutely continuous with respect to the orbifold structure. We show how to do this for the principle stratum $QD^1(\sigma)$. The measure is first defined locally on $Q(\sigma)$. A saddle connection of ϕ_0 persists under small perturbation of ϕ_0 . There are a finite number of such saddle connections β_1, \dots, β_n whose horizontal and vertical components $(h_{\beta_1}, v_{\beta_1}, \dots, h_{\beta_n}, v_{\beta_n})$ serve as local coordinates for $Q(\sigma)$ in a neighborhood of ϕ_0 . Specifically, pass to a double cover $\pi : \tilde{M} \rightarrow M$, ramified over the zeroes of ϕ_0 . There is an involution τ of \tilde{M} which interchanges the sheets. We choose the β_i such that their lifts to \tilde{M} form a basis for the first homology of \tilde{M} , odd with respect to the involution τ . The measure μ on $Q(\sigma)$ is then Lebesgue measure on \mathbf{R}^{2n} pulled-back to $Q(\sigma)$ via these local coordinates. One sees that μ does not depend on choice of basis β_i , so actually defines a measure. One then defines the μ_0 measure of $E \subset Q^1(\sigma)$ to be the μ measure of $\{r\phi : 0 \leq r \leq 1, \phi \in E\}$. It is easily seen to descend to $QD(\sigma)$ and to be $SL(2, \mathbf{R})$ invariant. The measure μ_0 was shown to be finite and ergodic on the principle stratum in [51] and then on each component of general stratum in [69], [70]. The following problem is quite interesting: *to classify all ergodic $SL(2, \mathbf{R})$ invariant measures on $QD(\sigma)$.*

3 Ergodicity

3.1 Veech nonergodic example

Definition 3.1 *A foliation F is uniquely ergodic if it is minimal and the transverse measure is unique up to scalar multiplication.*

It is well-known that if F is uniquely ergodic, then it is ergodic with respect to the transverse measure. Unique ergodicity is equivalent to the following condition. Let G be any foliation with transverse measure, transverse to F . For any point x_0 , the transverse foliation G defines an arc length $l(t)$ on the leaf of F through x_0 in

either direction. The pair (F, G) defines a measure μ on M . For every continuous function $f(x)$ on M and every point x_0 we have

$$\lim_{T \rightarrow \infty} \frac{1}{T} \int_0^T f(l(t)) dt = \int_M f(x) d\mu(x).$$

The classical example of a uniquely ergodic foliation is the irrational flow on the flat torus – see [39] for a detailed discussion.

Veech’s example. In [68] Veech constructed examples of minimal and not uniquely ergodic dynamical systems. An example of a minimal flow on a surface which is not uniquely ergodic was constructed by Sataev [62]. Examples of minimal nonuniquely ergodic interval exchanges were constructed by Keane [44].

We will look at Veech’s example in some detail. Take two copies of the unit circle and mark off a segment J of length $2\pi\alpha$ in the counterclockwise direction on each with one endpoint at 0. Now take θ irrational and consider the following dynamical system. Start with a point p , say in the first circle. Rotate counterclockwise by $2\pi\theta$ until the first time the orbit lands in J ; then switch to the corresponding point in the second circle, rotate by $2\pi\theta$ until the first time the point lands in J ; switch back to the first circle and so forth. Veech showed that if θ is an irrational number with unbounded partial quotients, then there are irrational numbers α such that this system is minimal and such that the Lebesgue measure is not ergodic so that the system is not uniquely ergodic (see also [39], section 14.5 e). We may describe the same dynamical system by a flow on a flat surface that comes from a rational billiard.

We actually give a slight generalization of the Veech example. Start with a square torus T ; that is, the unit square with lower left vertex at $(0, 0)$ whose opposite sides are identified. Equivalently, it is \mathbf{R}^2 modulo the integer lattice. Fix a point $(x_0, y_0) \in \mathbf{R}^2$, called a slit. Let w be a segment in \mathbf{R}^2 joining $(0, 0)$ to (x_0, y_0) , projected to T . Take two copies of T , each slit along w , and identify the positive side of w on one copy to the negative side in another. It is easy to see that this results in a quadratic differential ϕ on a surface of genus 2. Since the total angle around each of the points $(0, 0)$ and (x_0, y_0) is 4π , they correspond to zeroes of order 2 of ϕ . The surface of genus 2 is partitioned into 2 sheets M_w^+ and M_w^- separated from each other by the union of the two slits.

In the special case that $x_0 = 0$ and $y_0 = \alpha$, so that the slit is along the vertical axis, the resulting surface has two metric cylinders in the vertical direction, one for each sheet. Choose the core curve of each metric cylinder. The first return map to this pair of curves for the flow in direction θ is precisely the dynamical system studied by Veech.

This same dynamical system can be described in terms of rational billiards. Take a rectangle of length 1 and width $1/2$, with an interior vertical barrier of length $(1 - \alpha)/2$ from the midpoint of a horizontal side. It is not hard to see that the corresponding flat surface, formed by gluing four copies of this table, is exactly the flat surface described in the special case.

Fig. 6

Returning to more general construction, we say that the slit is irrational if either $x_0, y_0 \neq 0$ and y_0/x_0 irrational, or one coordinate is 0 and the other is irrational.

Theorem 3.2 *Suppose the slit (x_0, y_0) is irrational and S the corresponding flat surface. Then there exists uncountably many directions θ such that the flow on S in direction θ is minimal and not ergodic.*

Geometric criterion for nonergodicity. The proof of Theorem 3.2 represents unpublished work of John Smillie and the first author. In order to give the proof we first give a geometric criterion of how to find nonergodic foliations. This criterion was given in [56]. To set the notation let $P_n = [A_n, B_n]$ be a sequence of partitions of a flat surface S such that their common boundary is a union of saddle connections all in the same direction θ_n . Assume that the directions θ_n converge to θ_∞ . Let F_{θ_∞} be the foliation in direction θ_∞ . Rotate the coordinate system so that θ_∞ is the vertical direction. Let h_n be the sum of the horizontal components of the vectors associated to the saddle connections separating A_n from B_n . Let μ be the measure defined by the flat structure.

Theorem 3.3 *Suppose that*

- (i) $\lim_{n \rightarrow \infty} h_n = 0$
- (ii) for some $c, c', 0 < c \leq \mu(A_n) \leq c' < 1$
- (iii) $\sum_{n=1}^{\infty} \mu(A_n \Delta A_{n+1}) < \infty$

Then the vertical foliation F_{θ_∞} is nonergodic.

Proof. Let

$$A_\infty = \liminf A_n = \{x : \exists N \text{ such that for } n \geq N, x \in A_n\}.$$

$$B_\infty = \liminf B_n = \{x : \exists N \text{ such that for } n \geq N, x \in B_n\}.$$

We first show that A_∞ and B_∞ satisfy

- (1) $\mu(M \setminus (A_\infty \cup B_\infty)) = 0$.
- (2) $A_\infty \cap B_\infty = \emptyset$
- (3) $\mu(A_\infty \Delta A_n) \rightarrow 0$ as $n \rightarrow \infty$.
- (4) $0 < \mu(A_\infty) < 1$.

By (iii) and the Borel-Cantelli Lemma the set of x' which are in infinitely many $A_n \Delta A_{n+1}$ has μ measure 0. From this we have (1). Statement (2) is immediate. To see (3) note that

$$A_\infty \Delta A_n \subset \bigcup_{i=n}^{\infty} A_i \Delta A_{i+1}$$

so that

$$\mu(A_\infty \Delta A_n) \leq \sum_{i=n}^{\infty} \mu(A_i \Delta A_{i+1}).$$

Hypothesis (iii) implies that the right hand side goes to 0 as n goes to infinity, proving (3). Statement (4) follows from (3) and (ii).

Returning to the proof of Theorem 3.3, we can assume that the vertical foliation F_{θ_∞} is orientable: if not, we may replace M by an orientable double cover and replace A_∞ and B_∞ by their lifts. Now let f_t be the flow along vertical leaves.

We claim that for any t the set A_∞ is μ a.e. invariant; that is,

$$\mu(f_t(A_\infty) \Delta A_\infty) = 0.$$

Suppose on the contrary, that for some t_0 , we have $\mu(f_{t_0}(A_\infty) \Delta A_\infty) = \delta > 0$. By (3) and (i) we may choose n large enough so that

$$\mu(A_\infty \Delta A_n) < \delta/8 \quad \text{and} \quad e^{t_0} h_n < \delta/8.$$

Since f_{t_0} is μ preserving, the first inequality above and the assumption give

$$\mu(f_{t_0}(A_n) \Delta A_n) \geq \delta - 2\delta/8 = 3\delta/4.$$

Thus at time t_0 , $3\delta/8$ of the measure of A_n flows to its complement. However at most $e^{t_0} h_n < \delta/8$ of measure can cross the boundary of A_n , a contradiction, proving the claim. We would like to conclude that A_∞ is flow invariant, although by the claim it is only a.e. invariant for each time. The theorem is a consequence of the following general lemma.

Lemma 3.4 *Let f_t be a flow on a space X preserving a probability measure μ . Suppose there is a set A such that for every t , $\mu(f_t(A)\Delta A) = 0$. Then there is a set A' invariant under f_t with $\mu(A\Delta A') = 0$.*

Proof. Let λ be Lebesgue measure on \mathbf{R} . Let

$$A' = \{x : f_t(x) \in A \text{ for } \lambda \text{ a.e. } t\}.$$

It is clear that A' is f_t invariant. We will show that $\mu(A\Delta A') = 0$. Let

$$C_0 = \{(x, t) : x \in A\} \text{ and } C_1 = \{(x, t) : f_t(x) \in A\}.$$

For every t we have

$$\mu(\{x : (x, t) \in C_0\Delta C_1\}) = \mu(f_t(A)\Delta A) = 0.$$

This implies that $(\mu \times \lambda)(C_0\Delta C_1) = 0$. By Fubini there is a set X' of full μ measure so that for all $x \in X'$,

$$\lambda(\{t : (x, t) \in C_0\Delta C_1\}) = 0.$$

If $x \in X' \cap A$ then $(x, t) \in C_0$ so that $\lambda(\{t : (x, t) \notin C_1\}) = 0$. Therefore the set $\{t : f_t(x) \in A\}$ has full λ measure so $x \in A'$. If $x \in X' \setminus A$ then the set $\{t : (x, t) \notin C_1\}$ has full λ measure so that $x \notin A'$. Thus $A\Delta A'$ is contained in the complement of X' so $\mu(A\Delta A') = 0$. This finishes the proof of the Lemma and therefore Theorem 3.3. \square

Proof of Theorem 3.2 Let w be the segment joining $(0, 0)$ to (x_0, y_0) projected to T . Suppose w' is another segment on T with the same endpoints as w . Then we can take w' as the projection of a segment joining $(0, 0)$ to $(x_0 + p, y_0 + q)$ where $p, q \in \mathbf{Z}$ and w' is also thought of a slit. Suppose w and w' intersect an odd number of times in their interior so they divide each other into an even number of pieces. This is equivalent to saying that w and w' are homologous mod 2 on T . The sheet interchange measured by $c = (M_w^+ \cap M_{w'}^-) \cup (M_w^- \cap M_{w'}^+)$ is a union of an even number of parallelograms with sides on w and w' (here w and w' are thought of as vectors). Thus the area of c is at most $|w \times w'|$.

Fig. 7

Now fix a sequence of positive numbers ρ_j with $\sum \rho_j < \infty$. We will build an infinite directed tree with each vertex leading to 2 further vertices. At level j there

will be 2^{j-1} vertices. Each vertex will correspond to a pair of integers (p, q) which will determine a slit with endpoints $(0, 0)$ and $(x_0 + p, y_0 + q)$.

For any pair (p, q) form the quotient $(p+x_0)/(q+y_0)$, the slope of the corresponding slit. Let δ_j be the minimum distance between any two slopes as (p, q) varies over the vertices at level j . Now inductively, suppose we have determined the tree up to level j . For each vertex (p, q) at level j we will choose two vertices (p', q') and (p'', q'') at level $j + 1$, and call these the two *children* of (p, q) . To find such children begin by choosing an even integer d such that

$$\frac{\rho_j}{(q + y_0)(q + y_0 + d)} < \delta_j/4.$$

Then consider the inequality

$$d|(p + x_0)n - (q + y_0)m| < \rho_j$$

The assumption on the irrationality of the original slit (x_0, y_0) implies that $\frac{x_0+p}{y_0+q}$ is irrational so there are infinitely many coprime solutions (m, n) of the above inequality. Choose any two of them and set

$$p' = p + dm \quad \text{and} \quad q' = q + dn,$$

calling the two resulting pairs the children. A direct calculation shows that

$$\delta'_j = \left| \frac{p + x_0}{q + y_0} - \frac{p' + x_0}{q' + y_0} \right| < \delta_j/4.$$

That is to say, the distance between the slope of a slit and either child is bounded by $\delta_j/4$. Then the distance between the slopes of the two children of (p, q) is at most $\delta_j/2$, which implies that

$$\delta_{j+1} < \delta_j/2.$$

Let w be the segment corresponding to the parent and w' the segment of a child. Since d is even, w' is homologous to $w \pmod{2}$. In fact $w' - w$ represents d times the primitive class of (m, n) . As we have seen the area of the sheet interchange c is bounded by $|w' \times w|$, and an easy calculation shows that this is bounded by

$$|(p + x_0)n - (q + y_0)m| < \rho_j.$$

Thus, in the situation of Theorem 3.3, we have constructed a sequence of partitions satisfying (iii). Now for any geodesic in the tree, the sequence of ratios

$(p + x_0)/(q + y_0)$ is a Cauchy sequence hence converging to some θ_∞ . From this, condition (i) of Theorem 3.3 follows. Condition (ii) holds since the partition divides the surface into two pieces of equal areas. Thus by Theorem 3.3 the limiting foliation F_{θ_∞} is not ergodic.

We now show that, in fact, there are uncountably many limiting directions. Since there are only countably many directions that are not minimal, we then can conclude that there are uncountably many limiting directions which are minimal. Since each pair (p, q) has two children there are clearly uncountably many sequences. We therefore need to show that any two limits are distinct. Suppose $r_j = (p_j + x_0)/(q_j + y_0) \rightarrow \theta$ and a distinct sequence $r'_j = (p'_j + x_0)/(q'_j + y_0) \rightarrow \theta'$. Suppose the sequences differ for the first time at stage j , so that the slopes r_j and r'_j satisfy

$$|r_j - r'_j| \geq \delta_j. \tag{4.5}$$

Then by (4.3) and (4.4) we have

$$\left| \theta - \frac{p_j + x_0}{q_j + y_0} \right| \leq \sum_{k \geq j} \delta'_k < \sum_{k \geq j} \frac{\delta_k}{4} < \frac{\delta_j}{2},$$

and similarly $|\theta' - r'_j| < \delta_j/2$. Combined with (4.5), we have $\theta \neq \theta'$. □

The above theorem gives uncountably many minimal nonergodic directions. By requiring α to satisfy a Diophantine condition, it is even possible to show that there is a set of such examples of positive Hausdorff dimension (see Section 3.4 below). The main result of the next section however says that the set of nonergodic directions must have Lebesgue measure 0.

3.2 Ergodicity in almost every direction

Theorem 3.5 ([47]) *For any holomorphic quadratic differential ϕ (rational billiard) the set of $\theta \in [0, 2\pi]$ such that $r_\theta\phi$ has minimal but non-uniquely ergodic vertical foliation has Lebesgue measure 0.*

The proof of this theorem uses the $SL(2, \mathbf{R})$ action on quadratic differentials and some combinatorial constructions related to quadratic differentials. The quadratic differential ϕ belongs to some stratum $QD(\sigma)$. Let $QD_\epsilon(\sigma)$ be the set of quadratic differentials in $QD(\sigma)$ all of whose saddle connections have length at least ϵ . Recall that R_g is the moduli space of Riemann surfaces of genus g .

Proposition 3.6 $QD_\epsilon(\sigma)$ is compact in $QD(\sigma)$.

Proof. $QD_\epsilon(\sigma)$ is clearly closed in $QD(\sigma)$ since if $\phi_n \rightarrow \phi$ and ϕ has a geodesic segment of length less than ϵ so does ϕ_n for n large. It is therefore enough to show that $QD_\epsilon(\sigma)$ lies over a compact set in R_g . If not, there is a sequence $\phi_n \in QD_\epsilon(\sigma)$ of quadratic differentials lying on Riemann surfaces X_n going to infinity in R_g . Passing to subsequences we may assume that X_n converges to a Riemann surface X_∞ with nodes or punctures acquired by pinching along a set of disjoint simple closed curves $\alpha_1, \dots, \alpha_p$. By passing to a further subsequence we can assume that ϕ_n converges uniformly on compact sets to an integrable quadratic differential ϕ_∞ on X_∞ .

Thus ϕ_∞ has at most simple poles at the nodes. There is therefore a curve α_j homotopic to a puncture with ϕ_∞ length at most $\epsilon/2$. For large n the uniform convergence implies that α_j has ϕ_n length at most $\epsilon/2$, a contradiction. \square

Recall, that associated to a holomorphic quadratic differential ϕ is a structure on an associated Riemann surface. Recall also the the Teichmuller flow $g_t\phi$ from Section 2.4.

Definition 3.7 A quadratic differential ϕ is called *divergent* if the Riemann surfaces associated to $g_t\phi$ eventually leave every compact set of R_g as $t \rightarrow \infty$.

Theorem 3.8 Suppose the vertical foliation F_ϕ of ϕ is minimal but not uniquely ergodic. Then ϕ is divergent.

Proof. Suppose the theorem is not true so that for some sequence $t_n \rightarrow \infty$ the Riemann surfaces $X_n = X_{t_n}$ converge in R_g to some X_0 . By passing to a further subsequence we may assume that $\phi_n = g_{t_n}(\phi)$ converges to some ϕ_0 on X_0 . Let f_{t_n} be the corresponding Teichmuller map. Denote by Σ and Σ_0 the sets of zeroes of ϕ and ϕ_0 , respectively.

The normalized transverse measures on the topological foliation F_ϕ form a finite dimensional convex set for which the extreme points ν_i are mutually singular ergodic measures – see [77]. Let ν_0 be the transverse measure of F_ϕ . We have $\nu_0 = \sum_{i=1}^m c_i \nu_i$, for constants c_i . Pick a segment I of a horizontal trajectory of ϕ so that $\nu_i(I) \neq \nu_j(I)$ for $i \neq j$. Let μ the transverse measure to the horizontal foliation defined by ϕ . This means that the area element $|\phi(z)dz^2| = \nu_0 \times \mu$. Let $E_i \subset X$ be the disjoint sets of generic points of the measure $\nu_i \times \mu$ for the interval I . We have $(\nu_i \times \mu)(E_i) = 1$. Each E_i is a union of leaves of F_ϕ .

For each i let $A_i \subset X_0$ denote the set of accumulation points of $f_{t_n}(x) \in X_n$ for $x \in E_i$. Now let $U \subset X_0$ be any open set. We claim that there is a j such that

$$U \cap A_j \neq \emptyset.$$

To prove the claim, choose an open W so that $\bar{W} \subset U$. Since $X_n \rightarrow X_0$, for each n , W may be considered as a subset of X_n with area bounded below by $\delta > 0$. Consider $W_n = f_{t_n}^{-1}(W) \subset X$. Since f_{t_n} preserves area, $(\nu_0 \times \mu)(W_n) > \delta$.

Then for each n , there is a $j = j(n)$ such that

$$(\nu_j \times \mu)(E_j \cap W_n) = (\nu_j \times \mu)(W_n) \geq \frac{\delta}{mc_j}.$$

Choose a j so that this inequality holds for infinitely many n . This implies that there exists some $x \in E_j \cap W_n$ for infinitely many n . Since $f_{t_n}(x) \in W$, $f_{t_n}(x)$ has an accumulation point in \bar{W} , proving the claim.

Partition X_0 into a finite number of rectangles in the natural coordinates of ϕ_0 whose sides are horizontal and vertical segments. Adjacent rectangles meet along a common horizontal or vertical segment of a side. Each point of Σ_0 is required to be a vertex of a rectangle. For each pair of adjacent rectangles R_1, R_2 choose open sets $U_i \subset R_i$ with the property that if $y_i \in U_i$ there is a vertical segment l_i with one endpoint y_i such that l_1, l_2 are two sides of a coordinate rectangle contained in $R_1 \cup R_2$. For any two points y_1, y_2 in the same coordinate rectangle there are vertical segments l_i with endpoints y_i which are the vertical sides of a coordinate rectangle. Now by the claim, by passing to further subsequences of t_n , for each open set U_i we may choose $y_i \in U_i$ and $x_i \in \cup_{j=1}^m E_j$ such that $y_i = \lim f_{t_n}(x_i)$.

Fig. 8

Now suppose U_1, U_2 are such open sets contained in either adjacent or the same rectangle with corresponding y_i and x_i . We claim that x_1 and x_2 must belong to the same E_j . Assume otherwise. By renumbering assume $x_i \in E_i, i = 1, 2$. Let l_i be the vertical segments with endpoints y_i and which are vertical segments of a coordinate rectangle. They have the same length. Each l_i is the limit of vertical segments $l_{i,n}$ of equal length of $g_{t_n}\phi$ with one endpoint $f_{t_n}(x_i)$. Each $l_{i,n}$ is the image under f_{t_n} of a vertical segment $L_{i,n} \subset E_i$ passing through x_i ; $L_{i,n}$ have equal length which go to ∞ with n . Therefore the number of intersections of $L_{i,n}$ with I goes to ∞ and since $L_{i,n} \subset E_i$,

$$\lim_{n \rightarrow \infty} \frac{\text{card}(l_{i,n} \cap f_{t_n}(I))}{\text{card}(l_{2,n} \cap f_{t_n}(I))} = \lim_{n \rightarrow \infty} \frac{\text{card}(L_{1,n} \cap I)}{\text{card}(L_{2,n} \cap I)} = \frac{\nu_1(I)}{\nu_2(I)} \neq 1.$$

But since $l_{i,n}$ has limit l_i and every horizontal segment of ϕ_0 that intersects l_1 intersects l_2 and vice versa,

$$\frac{\text{card}(l_{1,n} \cap f_{t_n}(I))}{\text{card}(l_{2,n} \cap f_{t_n}(I))}$$

can be made arbitrarily close to 1 by taking n large enough. This is a contradiction, proving the claim.

Any two rectangles R_1, R_n can be connected by a chain of rectangles R_1, \dots, R_n where R_i and R_{i+1} are adjacent along a side. This fact and the last claim imply that there is a single index j such that $x_i \in E_j$ for all x_i .

Choose some index $i \neq j$ and $x \in E_i \setminus \Sigma$. Let $y \in X_0$ be an accumulation point of $f_{t_n}(x)$. By passing to further subsequences of t_n we can assume that there is a single rectangle R such that $f_{t_n}(x) \in R$ and $\lim_{n \rightarrow \infty} f_{t_n}(x) = y$. Let $U_k \subset R$ be one of the open sets found previously. We may find vertical segments $l \subset R$ with one endpoint y and $l' \subset R$ with endpoint $y_k = \lim f_{t_n}(x_k)$ such that l, l' are two vertical sides of a coordinate rectangle contained in R . We may find vertical segments l_n and l'_n of $g_{t_n}\phi$ of equal length, $l_n \rightarrow l$ and $l'_n \rightarrow l'$, l_n has endpoint $f_{t_n}(x)$ and l'_n has endpoint $f_{t_n}(x_k)$. As before

$$\frac{\text{card}(l_n \cap f_{t_n}(I))}{\text{card}(l'_n \cap f_{t_n}(I))}$$

can be made arbitrarily close to 1 by taking n large enough, and yet since x and x_k are points of E_i and E_j respectively, the ratio has limit

$$\frac{\nu_i(I)}{\nu_j(I)} \neq 1,$$

and we have our contradiction. □

The proof of Theorem 3.5 follows from the above theorem and the next

Theorem 3.9 *The set of $\theta \in [0, 2\pi]$ such that $r_\theta\phi$ is divergent, has Lebesgue measure zero.*

3.3 A combinatorial construction

In order to prove this theorem we will need to discuss a certain combinatorial construction. We fix a stratum $QD(\sigma)$.

Definition 3.10 *Two saddle connections are disjoint if their only common points are vertices.*

Denote by $p_0 = p_0(\sigma)$ the maximum number of disjoint saddle connections. A set of disjoint saddle connections Γ is called a *system*. By $|\Gamma|$ we mean the maximum length of any member of Γ . By a *complex* K we mean a subset of M consisting of disjoint saddle connections and triangles, such that a triangle is in K if and only if its sides are in K . An ϵ complex is a complex all of whose saddle connections have length at most ϵ . We denote the (topological) boundary of K by ∂K . Then ∂K is a system of saddle connections.

Note that if K is a complex, then $K \setminus \partial K$ (viewed as a subset of M) is a (possibly empty) domain. We say that the *complexity* of a complex is j if it has j saddle connections, including those on the boundary. The main construction is as follows.

Proposition 3.11 (Combining complexes) *Suppose K is a complex of complexity i with nonempty boundary and there is a saddle connection α which is either disjoint from K or crosses ∂K . Then there is a complex K' containing K of complexity $j > i$ with the property that $|\partial K'| \leq 2|\partial K| + |\alpha|$.*

Proof. Suppose first that K fails to be convex in the sense that there are segments AB and BC on ∂K such that the angle they make at B is smaller than π . This means that AB followed by BC is not a geodesic. Further suppose that the geodesic joining A to C does not lie in K . Then we may add this geodesic and the inequality on length holds.

Thus we may assume that K is convex in the sense that there are no such segments AB and BC . Now if α is disjoint from K , we may just add it to K and again the length inequality holds.

Thus we may assume that α crosses ∂K . Let Q be the point where α crosses ∂K on a segment ω with endpoints A, B . Suppose that in following α from Q into the exterior of K α hits a singularity P in the exterior of K before again crossing ∂K or hits a singularity on ∂K .

We may form a polygon with no interior singularities whose sides are ω , a possibly broken geodesic joining A to P and a possibly broken geodesic joining B to P . The edges and diagonals of this polygon have lengths bounded by $|\alpha| + |\omega|$. Since ω is a boundary edge of K not all of the edges and diagonals of the polygon can be in K . We may therefore add an edge that is not in K .

Fig. 9

Thus we are reduced to the case that before next hitting a singularity, α crosses ∂K at Q' , an interior point of a segment ω' with endpoints C, D . It may be the case

that $A = D$ or $B = C$ or both. Let α_1 be the subsegment of α joining Q to Q' . If $A = D$ or $B = C$, then α_1 is homotopically nontrivial relative ∂K ; otherwise ω and ω' would bound a triangle in the exterior of K , contrary to the convexity assumption. Now move α_1 parallel to itself with one endpoint moving along ω toward one endpoint, say A , and the other along ω' towards, say, D . We can choose the direction so that the lengths of the parallel segments do not increase in length. This can be done until for a first time a segment α_2 , parallel to α_1 , meets a singularity P , which may be either A or D . We proceed exactly as in the previous paragraph. We can use P to build a bigger complex whose boundary satisfies the same length estimate as before. \square

Definition 3.12 *Let $\epsilon < C$. A saddle connection γ is (ϵ, C) isolated if it has length less than ϵ and every saddle connection that crosses γ in its interior has length greater than C .*

Since any two (ϵ, C) isolated curves are disjoint, the number of such curves is bounded above by p_0 .

Definition 3.13 *For a quadratic differential ϕ let $n_\epsilon(\phi)$ be the maximum number of simplices in a connected ϵ complex.*

Definition 3.14 *For a stratum $QD(\sigma)$ let $N(\epsilon, C)$ be the set of quadratic differentials in the stratum which possess (ϵ, C) isolated saddle connections.*

Proposition 3.15 *Let μ denote the Lebesgue measure on the circle. Suppose there is a set S of angles θ of positive measure so that $r_\theta\phi$ is divergent for $\theta \in S$. Then there are*

- (1) a sequence of times $T_i \rightarrow \infty$
- (2) a sequence of sets $S_i \subset S$ and a number $\delta > 0$ such that $\mu(S_i) \geq \delta$
- (3) a sequence $\epsilon_i \rightarrow 0$
- (4) a positive constant C such that $g_{T_i} r_\theta \phi \in N(\epsilon_i, C)$ for $\theta \in S_i$.

The idea behind this Proposition is that since the Riemann surface of $g_t r_\theta \phi$ eventually leaves every compact set, there will be times T_i after which there will always be a segment of length less than ϵ_i in the corresponding metrics. For a fixed C , however, the sets $N(\epsilon_i, C)$ do not form a neighborhood basis of infinity so it may not be the case that for sufficiently large T_i that there will be always (ϵ_i, C) isolated segments. Nevertheless by passing to slightly smaller sets than S_i we will be able to make the statement for isolated segments.

Proof. Choose a sequence $\epsilon_i \rightarrow 0$. Since for each ϵ_i the set QD_{ϵ_i} is compact, for each divergent ray $r_\theta(q)$ there are times T_i such that $g_t r_\theta(q) \in QD \setminus QD_{\epsilon_i}$ for all $t \geq T_i$. Then T_i can be chosen so that this relationship holds for all $\theta \in S'$ where $S' \subset S$ and $\mu(S') > \mu(S)/2$.

Consider now triples of sequences (ϵ_i, T_i, S_i) such that $\epsilon_i \rightarrow 0$, $T_i \rightarrow \infty$, $\mu(S_i)$ bounded away from 0 and $g_{T_i} r_\theta \phi \in QD \setminus QD_{\epsilon_i}$ for $\theta \in S_i$. The set of triples is nonempty for we can take $S_i = S'$ for all i . Choose a sequence of triples that maximizes

$$\min_i \min_{\theta \in S_i} n_\epsilon(g_{T_i} r_\theta \phi).$$

For the corresponding sequence, for each i and $\theta \in S_i$ let $C_{i,\theta}$ be the length of the shortest saddle connection that crosses a boundary segment of the ϵ_i complex of $g_{T_i} r_\theta \phi$. For each i , let m_i be a number so that

$$\mu(\{\theta \in S_i : C_{i,\theta} \geq m_i\}) = \mu(S_i)/2.$$

We claim that the numbers m_i are bounded below. Suppose not. Choose a subsequence converging to 0 and for that subsequence let

$$S'_i = \{\theta \in S_i : C_{i,\theta} < m_i\},$$

so $\mu(S'_i) = \mu(S_i)/2$ is bounded below, and let

$$\epsilon'_i = 3\epsilon_i + m_i \rightarrow 0.$$

Fix $\theta \in S'_i$. Since there is an ϵ_i complex and a saddle connection of length at most m_i crossing it, by Proposition 3.11 we can find an ϵ'_i complex on $g_{T_i} r_\theta \phi$ with more simplices, contradicting the maximality of our sequence of triples. Thus the set of m_i is bounded below, by some $C > 0$. If we replace S_i with the set of θ such that $C_{i,\theta} \geq m_i \geq C$, every saddle connection crossing the ϵ_i complex has length at least C , proving the proposition. \square

Let α be a geodesic segment. For some θ_α , the segment α is vertical with respect to $r_{\theta_\alpha}\phi$. Rotate the coordinate system so that $\theta_\alpha = 0$. Denote by $l(t, \theta)$ the length of α with respect to the structure $g_t r_\theta \phi$, so that $|\alpha|_\phi = l(0, 0)$. In these coordinates

$$l(t, \theta) = l(0, 0)(e^t \sin^2 \theta + e^{-t} \cos^2 \theta)^{1/2}.$$

Fix $t \geq 0$ and choose $\epsilon > 0$. Define two intervals $I_\alpha \subset J_\alpha \subset (-\pi/2, \pi/2)$ with respect to the rotated coordinates:

$$I_\alpha = \{\theta : |\sin \theta| \leq \epsilon/(ve^{t/2})\},$$

$$J_\alpha = \{\theta : |\sin \theta| \leq C/(2ve^{t/2})\}.$$

Lemma 3.16 *For $\theta \notin I_\alpha$ one has: $l(t, \theta) \geq \epsilon$. There are constants ϵ_1, T, K independent of α such that for $\epsilon < \epsilon_1$ and $t > T$ we have $\mu(I_\alpha)/\mu(J_\alpha) < K\epsilon$. In addition if $l(t, \theta) < \epsilon$ for some θ , then $l(t, \theta) < C$ for all $\theta \in J_\alpha$.*

Proof. The first statement is immediate from the definition. If $l(t, \theta) < \epsilon$ for some θ then, since $l(t, \theta)$ attains its minimum at $\theta = 0$,

$$ve^{-t/2}|\cos \theta| \leq ve^{-t/2} \leq l(t, \theta) < \epsilon.$$

For $\theta \in J_\alpha$, we then have

$$l(t, \theta) \leq (C^2/4 + \epsilon^2 e^{-t})^{1/2} \leq C/2 + \epsilon.$$

If we choose $\epsilon_1 = C/2$ and $\epsilon < \epsilon_1$ then $l(t, \theta) < C$. This proves the last statement.

Let I'_α and J'_α be the images of I_α and J_α under the sine function. If $C/(2ve^{t/2}) < 1/2$ then J'_α is contained in the interval $(-1/2, 1/2)$ and $\mu(I'_\alpha)/\mu(J'_\alpha) = 2\epsilon/C$. Since the arcsine function restricted to the interval $(-1/2, 1/2)$ is Lipschitz and has a Lipschitz inverse, it changes the lengths of intervals by a bounded amount. Thus we can choose a constant K depending on C such that $\mu(I_\alpha)/\mu(J_\alpha) < K\epsilon$ as claimed.

To ensure that $C/(2ve^{t/2}) < 1/2$ choose T such that $C/(2me^{T/2}) < 1/2$ where m is the length of the shortest saddle connection on ϕ . Since $v \geq m$ and $t > T$, the inequality holds and the proof of the lemma is finished. \square

The last technical statement we need is as follows.

Proposition 3.17 *There are constants T, ϵ_1, K' such that for $t > T$ and $\epsilon < \epsilon_1$ one has:*

$$\mu(\{\theta : g_t r_\theta \phi \in N(\epsilon, C)\}) < K'\epsilon.$$

Proof. Fix T and $\epsilon < \epsilon_1$ as in the preceding lemma. Consider all saddle connections α which are (ϵ, C) isolated with respect to $g_t r_\theta \phi$ for some θ . For each such α let \hat{I}_α be the smallest interval of θ 's for which α is (ϵ, C) isolated. Then $\hat{I}_\alpha \subset I_\alpha \subset J_\alpha$.

Construct a new open interval \hat{J}_α as follows: the left-hand endpoint of \hat{J}_α will be halfway between the left-hand endpoints of \hat{I}_α and J_α , and the right-hand endpoint of \hat{J}_α will be halfway between the right-hand endpoints of \hat{I}_α and J_α . The point of this construction is the fact that if \hat{J}_α and \hat{J}_β intersect, then either J_α intersects \hat{I}_β or J_β intersects \hat{I}_α . Assume the first. We can find $\theta \in J_\alpha \cap \hat{I}_\beta$ such that β is (ϵ, C) isolated with respect to $g_t r_\theta q$. Now, with respect to the metric of $g_t r_\theta \phi$, the saddle connection α has length less than C , and β cannot cross any geodesic segment of length less than C . Thus α and β are disjoint. The maximum number of disjoint segments is p_0 . Therefore no θ can lie in more than p_0 segments \hat{J}_α . The sum of the lengths of the \hat{J}_α is at most $2\pi p_0$. The sum of the lengths of the J_α is at most $4\pi p_0$. The sum of the lengths of the intervals I_α is at most $4K\pi\epsilon p_0$, and the intervals I_α cover the set of θ for which $g_t r_\theta \phi \in N(\epsilon, C)$. Thus we can take $K' = 4K\pi p_0$. \square

Proof of Theorem 3.9. Assume the theorem is false. Then choose constants $\epsilon_i \rightarrow 0$, times $T_i \rightarrow \infty$, a constant C , and sets S_i so that Proposition 3.15 holds. In particular, $\mu(S_i) \geq \delta > 0$. Let K', T be as in Proposition 3.17; then, for ϵ_i small enough and $T_i > T$ large enough, we have: $\mu(S_i) < K'\epsilon_i$. However, for ϵ_i small enough, one has: $K'\epsilon_i < \delta$, a contradiction. \square

3.4 Further ergodicity results

In this section we mention, without proofs, some further results on ergodicity of vertical foliations of quadratic differentials.

The measure 0 result in Theorem 3.5 was improved in [55] to give a statement about Hausdorff dimension.

Theorem 3.18 *For any quadratic differential ϕ the set of $\theta \in [0, 2\pi]$ such that $r_\theta \phi$ has minimal but non-uniquely ergodic vertical foliation has Hausdorff dimension at most $1/2$.*

Recently Y.Cheung ([12]) has shown that the bound $1/2$ is sharp by revisiting the Veech examples described at the beginning of this chapter. Assume that the number α in the Veech example is not Liouville, that is, there exists $s \geq 2$ and $c > 0$ such that $|\alpha - p/q| > c/q^s$ for all p, q . Under this assumption Cheung proved that the set of

directions θ for which the system is minimal but not uniquely ergodic has Hausdorff dimension $1/2$. This is to be contrasted with an unpublished result of Boshernitzan who showed the same set has Hausdorff dimension 0 for a residual set of α .

On the other hand positive Hausdorff dimension turns out to be typical as was shown in [56]. Recall the stratum $QD^1(\sigma)$ and $SL(2, \mathbf{R})$ ergodic invariant measure μ_0 on each component of $QD^1(\sigma)$. Assume that $\sigma \neq (\emptyset; +), (1, -1; -), (-1, -1, -1, -1; -)$, that is, the quadratic differentials do not define a flat torus, once punctured torus, or four times punctured sphere.

Theorem 3.19 *For each component C of $QD^1(\sigma)$ there is a $\delta = \delta(C) > 0$ such that for μ_0 almost all $\phi \in C$, the set of $\theta \in [0, 2\pi]$ such that $r_\theta\phi$ has minimal but non-uniquely ergodic vertical foliation has Hausdorff dimension δ .*

3.5 Ergodicity of general polygonal billiards

Now we apply Theorem 3.5 to polygonal billiards. As before, the collection of plane polygonal regions with n vertices and a given combinatorial type is identified with a subset of \mathbf{R}^{2n} and is given a subspace topology. Recall that $G(Q)$ denotes the group generated by reflections in the sides of a polygon Q .

The next result is in the spirit of Theorem 1.10 but is stronger.

Theorem 3.20 ([47]) *Let X be a closed subset of the space of billiard tables with the property that for any number N the set of rational tables $Q \in X$ with $\text{card}(G(Q)) \geq N$ is dense. Then ergodic tables in X form a dense G_δ subset.*

Proof. We can assume that $X \subset \mathbf{R}^{2n}$ is compact. Each $Q \in X$ is a polygonal region in the plane. We can assume that the area of Q is 1. Let PX be the bundle whose base space is X and the fiber P_Q over $Q \in X$ is the phase space of the corresponding billiard table Q . Note that $PX \subset X \times \mathbf{R}^2 \times \mathbf{S}^1$. Let μ_Q be the product of the area measure on Q with unit Lebesgue measure on \mathbf{S}^1 , and let ϕ_t be the billiard flow on PX .

Choose a sequence of continuous functions f_1, f_2, \dots on PX which, when restricted to each P_Q , are dense in $L^2(P_Q)$. We make the further assumption that if v is an outward and v' is the corresponding inward vector on the boundary of Q , then $f_i(v) = f_i(v')$. Let $E(i, n, T)$ be the set of $Q \in X$ for which

$$\int_{z \in P_Q} \left(\frac{1}{T} \int_0^T f_i(\phi_t(z)) dt - \int_{P_Q} f_i d\mu \right)^2 d\mu < 1/n.$$

Let

$$E(i, n) = \cup_{T=1}^{\infty} E(i, n, T),$$

and let

$$E = \cap_{i=1}^{\infty} \cap_{n=1}^{\infty} E(i, n).$$

The set of $Q \in X$ for which ϕ_t restricted to P_Q is ergodic is precisely E – see [58]. We will prove that

- (1) The sets $E(i, n, T)$ are open, and
- (2) For a given i and n , there is an N such that $E(i, n)$ contains all rational tables Q for which $\text{card}(G(Q)) > N$.

Assuming these statements, the theorem is proved as follows. The first statement implies that the sets $E(i, n)$ are open. The second one implies that the $E(i, n)$ are dense. Then E is a G_δ , and it follows from the Baire category theorem that E is dense in X .

Statement (1) is a consequence of the next result.

Lemma 3.21 *Let $T > 0$ be fixed, and let f be a continuous function on PX respecting the boundary identifications. Then*

$$\int_{z \in P_Q} \left(\frac{1}{T} \int_0^T f(\phi_t(z)) dt - \int_{P_Q} f d\mu_Q \right)^2 d\mu_Q$$

depends continuously on Q .

Proof.

For $Q \in X$ let $a(Q) = \int_{P_y} f d\mu_Q$. Clearly, $a(Q)$ depends continuously on Q .

Replace the function f by the function \tilde{f} defined as follows: for $z \in P_Q$ let $\tilde{f}(z) = f(z) - a(Q)$. Then the proof of the lemma reduces to the proof of the continuity of the following function:

$$c(Q) = \int_{z \in P_Q} \left[\frac{1}{T} \int_0^T \tilde{f}(\phi_t(z)) dt \right]^2 d\mu_Q.$$

The difficulty is that, due to discontinuities of the billiard flow at the corners, the billiard orbits of close velocity vectors may diverge. Roughly speaking, this difficulty is overcome by deleting the set of velocity vectors whose trajectories hit a corner; this

set has zero measure, and the time T trajectory depends continuously on the velocity vector in its complement.

More specifically, for $z \in P_Q \subset PX$ introduce an auxiliary function $l(z)$ as the infimum of the distance from the time T trajectory of vector z to the set of vertices of the polygon Q . Clearly $l(z)$ depends continuously on z , and the zero level set of $l(z)$, that is, the set of velocity vectors whose time T trajectory hits a vertex, has zero measure.

Let $\epsilon > 0$ be fixed, and let $M = \sup \tilde{f}$. For a neighborhood N_1 of Q in the space of polygons denote by Q_1 the intersection of all polygons from N_1 (note that Q_1 is not a polygon). We can choose N_1 so small that the area of Q_1 is at least $1 - \epsilon/6M^2$.

Let $C_\delta \subset P_Q$ be the set of velocity vectors z for which $l(z) \geq \delta$. Choose δ small enough so that the measure of C_δ is at least $1 - \epsilon/6M^2$. Let $D_\delta \subset C_\delta$ consist of pairs $(x, v) \in \mathbf{R}^2 \times \mathbf{S}^1$ with $x \in Q_1$. Note that the measure of D_δ is at least $1 - \epsilon/3M^2$. Since $Q_1 \times \mathbf{S}^1$ is contained in P_R for all $R \in N_1$, we can identify $N_1 \times D_\delta$ with a subset of PX in a natural way. Let $\bar{l}(R)$ be the infimum of $l(z)$ for $z \in (P_R \cap (N_1 \times D_\delta))$. Now \bar{l} is continuous and $\bar{l}(Q) = \delta$. We can find a neighborhood $N_2 \subset N_1$ so that for all $R \in N_2$, $\bar{l}(R) > 0$. Let $d(R)$ denote

$$d(R) = \int_{z \in (P_R \cap (N_1 \times D_\delta))} \left[\frac{1}{T} \int_0^T \tilde{f}(\phi_t(z)) dt \right]^2 d\mu_R.$$

Then

$$|d(R) - c(R)| \leq M^2 \mu(P_R \setminus (P_R \cap (N_1 \times D_\delta))) \leq \epsilon/3$$

for each $R \in N_2$. For $R \in N_2$ the function $f(\phi_t(z))$ depends continuously on $t \in [0, T]$ and $z \in D_\delta$. Thus we can find a neighborhood N_3 of Q in which d varies by less than $\epsilon/3$. Then, for $R \in N_3$,

$$|c(Q) - c(R)| \leq |c(Q) - d(Q)| + |d(Q) - d(R)| + |d(R) - c(R)| \leq \epsilon.$$

This completes the proof of continuity of c and the proof of the Lemma. \square

To finish the proof of the Theorem one needs to establish statement (2); this statement follows from the next lemma.

Lemma 3.22 *Fix $n > 0$ and let f be a continuous function on PX . Choose $\delta > 0$ so that if $|\theta_1 - \theta_2| < \delta$ then $|f(\theta_1) - f(\theta_2)| < 1/2n$. Let $N \geq 2/\delta$, and let Q be a rational polygon with $|G(Q)| \geq N$. Then for sufficiently large T one has:*

$$\left[\int_{z \in P_Q} \left[\frac{1}{T} \int_0^T f(\phi_t(z)) dt - \int_{P_x} f d\mu_Q \right]^2 d\mu_Q \right]^{1/2} < 1/n.$$

Proof. Since Q is fixed we will drop the subscripts from P, G and μ . For $\theta \in \mathbf{S}^1$ let $u(\theta) = 1/|G| \int_{z \in M_\theta} f(z) dA$ where dA is the area measure on M_θ ; here, as before, M_θ is the invariant surface of the billiard flow in direction θ . For $z \in M_\theta$ let $u'(z) = u(\theta)$. For $z \in P$ let $v_T(z) = 1/T \int_0^T f(\phi_t(z)) dt$. The quantity which appears in the lemma is the norm in the space $L^2(P)$ of the function $v_T - f f d\mu$. We claim that

$$\lim_{T \rightarrow \infty} \|v_T - u'\| = 0.$$

To prove the claim, notice that the surfaces M_θ are parametrized by $\theta \in \mathbf{S}^1/\mathbf{G}$. We evaluate the norm by integrating first with respect to M_θ and then with respect to θ :

$$\|v_T - u'\| = [|G| \int_{\theta \in \mathbf{S}^1/\mathbf{G}} \frac{1}{|G|} \int_{z \in M_\theta} (v_T(z) - u'(z))^2 dA d\theta]^{1/2}.$$

Let

$$w_T(\theta) = [\frac{1}{|G|} \int_{z \in M_\theta} (v_T(z) - u'(z))^2 dA]^{1/2}.$$

Then

$$\|v_T - u'\| = [|G| \int_{\theta \in \mathbf{S}^1/\mathbf{G}} w_T(\theta)^2 d\theta]^{1/2}.$$

For a given θ the ergodicity of ϕ restricted to M_θ implies that $\lim_{T \rightarrow \infty} w_T(\theta) = 0$. Theorem 3.5 implies ergodicity for almost all θ . Since the functions w_T are bounded and converge pointwise almost everywhere to 0, they converge to 0 in norm. This completes the proof of the claim.

The second claim is that

$$\|u' - \int f d\mu\| \leq 1/2n.$$

Now $\|u' - \int f d\mu\|$ is equal to the norm of $u - \int f d\mu$ in $L^2(\mathbf{S}^1)$. Let θ_1 and θ_2 be points in the circle at which u assumes its minimum and maximum values m and M , respectively. Note that u is constant on the orbits of G . The distance between neighboring points in a G orbit is less than $2/|G| < 2/N < \delta$. By replacing θ_2 by $g(\theta_2)$, where $g \in G$, we may assume that $|\theta_1 - \theta_2| < \delta$. It follows from the continuity assumption on f that since $|\theta_1 - \theta_2| < \delta$, then $|u(\theta_1) - u(\theta_2)| < 1/2n$. Since u is defined by averaging f , $\int_P f = \int_{\mathbf{S}^1} u$. Thus $m \leq \int f \leq M$; hence $|u(\theta) - \int f| \leq 1/2n$ and $\|u - \int f\| \leq 1/2n$. This completes the proof of the claim.

We now complete the proof of the lemma. Choose T sufficiently large so that $\|v_T - u'\| \leq 1/2n$. Then

$$\|v_T - \int f\| \leq \|v_T - u'\| + \|u' - \int f\| \leq 1/2n + 1/2n = 1/n,$$

and we are done. □

3.6 Constructive approach to polygons with ergodic billiard flow

Ya. Vorobets gave a constructive description of a topologically massive set of polygons with ergodic billiard flow – see [75]. We describe his result without proof.

Definition 3.23 *Let $\phi(N)$ be a positive function on \mathbf{N} whose limit is zero as $N \rightarrow \infty$. Let Q be a simple k -gon with angles $\alpha_1, \dots, \alpha_k$ between the adjacent sides. We say that Q admits approximation by rational polygons at the rate $\phi(N)$ if for every $n > 0$ there is $N > n$ and positive integers n_1, \dots, n_k , each coprime with N , such that $|\alpha_i - \pi n_i/N| < \phi(N)$ for all i .*

The next theorem gives an explicit estimate of how well a polygon should be approximated by rational ones to guarantee ergodicity of the billiard flow.

Theorem 3.24 *Let Q be a polygon that admits approximation by rational polygons at the rate*

$$\phi(N) = \left(2^{2^{2^N}}\right)^{-1}.$$

Then the billiard flow in Q is ergodic.

The paper [75] also contains constructive proofs of other results on polygonal billiards and flat surfaces. In particular it gives a new proof of the quadratic upper bound on the number of saddle connections, a result we will discuss in more detail in the next section.

4 Periodic orbits

4.1 Periodic directions are dense

Let ϕ be a quadratic differential with the corresponding flat metric and β a free homotopy class of simple closed curves. It is a standard fact that there is a geodesic representative in the homotopy class. The geodesic will typically pass through zeroes of ϕ , making angles in excess of π at the zero. If however there is a geodesic representative that does not pass through a singularity, then the geodesic can be moved

parallel to itself and sweeps out a cylinder of homotopic geodesics; we called such a cylinder a *metric cylinder*.

The boundary of the metric cylinder is a union of parallel saddle connections. In the special case of the flat torus, there are closed geodesics in a dense set of directions, and, of course, for each such direction the closed geodesics fill the surface. The main objective of this section is to prove the following theorem.

Theorem 4.1 ([52]) *For any quadratic differential ϕ_0 there is a dense set of directions $\theta \in \mathbf{S}^1$ such that ϕ_0 has a metric cylinder in the direction θ .*

Corollary 4.2 *For any rational billiard table there is a dense set of directions with a periodic orbit in that direction.*

Theorem 4.1 was strengthened in the following result of Boshernitzan, Galperin, Kruger and Troubetzkoy – see [11].

Theorem 4.3 *For any quadratic differential ϕ_0 there is a dense set of vectors in the tangent space to the surface such that the orbit determined by that vector is closed.*

The theorem is proved in [11] for rational billiards. The proof holds verbatim for the general case of a quadratic differential.

The idea of the proof of Theorem 4.1 is to study limit points of the $SL(2, \mathbf{R})$ orbit of ϕ_0 . Recall first, that in the compactification $\bar{R}_{g,0}$ of the moduli space $R_{g,0}$ of compact Riemann surfaces, we allow pinching Riemann surfaces along $p \leq 3g - 3$ simple closed curves to produce Riemann surfaces X_0 with nodes or punctures. Also recall that if $X_n \rightarrow X_0$ in this compactification and ϕ_n is a sequence of unit norm quadratic differentials on X_n , then via a conformal embedding of $X_0 \setminus U$ into X_n , for U an arbitrary neighborhood of the punctures on X_0 , we can pass to a subsequence and assume that ϕ_n converges uniformly on $X_0 \setminus U$ to a finite norm quadratic differential ϕ_∞ on $X_0 \setminus U$. It may or may not be the case that $\phi_\infty \equiv 0$.

Definition 4.4 *A limiting quadratic differential ϕ_∞ on X_0 is called exceptional if*

- (a) *some component of X_0 is a torus and ϕ_∞ defines the flat metric on the torus or*
- (b) *some component of X_0 is a punctured sphere and ϕ_∞ has no zeroes or*
- (c) *$\phi_\infty \equiv 0$ on every component of X_0 .*

We note that in case (b), since ϕ_∞ has no zeroes, there must be exactly four punctures at which ϕ_∞ has a simple pole, and at any other puncture, ϕ_∞ is regular. The importance of this definition is illustrated by the following proposition. We are indebted to Yair Minsky for suggesting the proof. A different proof appeared in [52] based on [51].

Proposition 4.5 *There exists $M > 0$ with the property that if $X_n \rightarrow X_0$ in $\bar{R}_{g,0}$ and ϕ_n are unit norm holomorphic quadratic differentials on X_n which converge uniformly on compact sets of X_0 to an exceptional ϕ_∞ on X_∞ , then for large enough n , ϕ_n has a metric cylinder of length at most M .*

Proof. In cases (a) and (b) there is a dense set of directions such that ϕ_∞ has a metric cylinder in that direction. Choose one such metric cylinder and a neighborhood of the punctures such that the cylinder contains a closed curve missing that neighborhood. The uniform convergence of ϕ_n to ϕ_∞ on the complement of this neighborhood implies that for large n , ϕ_n has a closed regular geodesic with some uniformly bounded length.

Thus we can assume we are in case (c). For any $\epsilon > 0$ and for each loop C surrounding a puncture of X_0 we have $|C|_{\phi_n} < \epsilon$ for large n . The loops C divide X_n into p annuli B_i homotopic to the pinching curves and $\phi_n \rightarrow \phi_\infty$ uniformly on $X_0 \setminus (\cup B_i)$. Since $\phi_\infty \equiv 0$, the ϕ_n area of $X_0 \setminus (\cup B_i)$ goes to 0 and so the ϕ_n area of $\cup B_i$ goes to 1. For large n , for at least one annulus B , the ϕ_n area of B is at least $1/(3g - 3)$. This annulus has two boundary components C .

Fix such a large n . If there is no metric cylinder in the homotopy class of C then we may let α be the unique ϕ_n geodesic homotopic to the loop C . Then α passes through singularities, at some of which the angle is in excess of π . We will arrive at a contradiction. First assume that α is embedded. For each r let $N(r)$ be the r neighborhood of α ; the set $N(r)$ is convex. For small values of r , $N(r)$ is an annulus containing α in its interior (here is where we use that α is embedded). For all r it is a domain with nonpositive Euler characteristic. Let $\alpha_r = \partial N(r)$. For all except an isolated set of r , the curve α_r is smooth.

The Gauss-Bonnet formula says that

$$\int_{N(r)} K + \int_{\alpha_r} \kappa ds = 2\pi\chi(N(r)) \leq 0,$$

where K is the Gaussian curvature and κ is the geodesic curvature. Now $K = 0$ except at the interior points x of $N(r)$, where the curvature is negative and concentrated so that the contribution to the first term is $-m\pi$, for m a positive integer. The number

m may jump a bounded number of times as more zeroes are included in $N(r)$ and so is maximized by some M . Thus

$$\int_{\alpha_r} \kappa ds \leq m\pi \leq M\pi.$$

with equality in the first inequality as long as $N(r)$ remains an annulus. Let $A(r) = ||N(r)||$, the area of $N(r)$, and let $L(r)$ denote the length of α_r . Then

$$A'(r) = L(r) \quad \text{and} \quad L'(r) = \int_{\alpha_r} \kappa ds \leq m\pi \leq M\pi,$$

again with equality in the first inequality as long as $N(r)$ remains an annulus. In particular this gives $L(r) \geq \pi r$, as long as $N(r)$ remains an annulus. Now the fact that $A''(r) \leq M\pi$ together with $A(0) = 0$ and $A'(0) \leq \epsilon$, gives

$$A(r) \leq M\pi r^2 + \epsilon r.$$

Choose ϵ small enough so that

$$A(\epsilon) \leq \frac{1}{3g-3}.$$

We show that $B \subset N(\epsilon)$ and that this leads to the desired contradiction. Let $r_0 = \epsilon/\pi$. The first step is to notice that each component C of the boundary of B must intersect $N(r_0)$. If $N(r_0)$ is an annulus this follows from the fact that $L(r_0) \geq r_0\pi = \epsilon$ and $N(r_0)$ is convex, so any curve, namely C , of length smaller than ϵ cannot be homotopic to and exterior to $N(r_0)$. If $\chi(N(r_0)) < 0$, then C must intersect $N(r_0)$ since C is homotopic to α . Since C has length at most ϵ the fact that C intersects $N(r_0)$ implies that C must be contained in $N(\epsilon(1/\pi + 1/2)) \subset N(\epsilon)$. Then both boundary components C of B are contained in $N(\epsilon)$, so we must have $B \subset N(\epsilon)$. But now

$$||B|| \leq A(\epsilon) \leq \frac{1}{(3g-3)},$$

which is a contradiction.

In case α is not embedded, so that α is not in the interior of $N(r)$, we lift α to the annular cover so that it is embedded. We have the same lower bound growth estimates for the length of curves in a neighborhood of the lifted α as well as the upper bound for the area of the neighborhood. This neighborhood projects to $N(r)$. The rest of the proof then goes through. \square

Thus in order to prove Theorem 4.1 we will find in the $SL(2, \mathbf{R})$ orbit of ϕ_0 a sequence of quadratic differentials that converge to an exceptional limit. We will then apply Proposition 4.5.

The following lemma says that the Teichmüller flow g_t extends continuously to the compactification. We need to add to the definition that if $\phi_0 \equiv 0$ on a Riemann surface X_0 , then $g_t\phi_0 \equiv 0$ for every t .

Lemma 4.6 *Suppose ϕ_n converges uniformly to ϕ_0 as above. Then for every t , $g_t(\phi_n) \rightarrow g_t(\phi_0)$ uniformly on compact sets of the Riemann surface of $g_t\phi_0$ as $n \rightarrow \infty$.*

Proof. Let X_0 be the Riemann surface of ϕ_0 and X_n the Riemann surface of ϕ_n . For any component of X_0 on which ϕ_0 is not identically 0, choose natural coordinates $z_0 = x_0 + iy_0$ for ϕ_0 . Then, via the conformal embedding of compact subsets of X_0 in X_n , one may view z_0 as a local parameter for X_n , n large, and ϕ_n is a holomorphic function of z_0 . We have, by assumption, that $\phi_n(z_0) \rightarrow \phi_0(z_0) \equiv 1$. If we let $z_n = x_n + iy_n$ be the natural coordinates of ϕ_n , which are then functions of z_0 , we have

$$z_n(z_0) \rightarrow z_0.$$

Since the Teichmüller maps corresponding to $g_t\phi_n$ and $g_t\phi_0$ are given in the natural coordinates by

$$x_n \rightarrow e^{t/2}x_n, \quad y_n \rightarrow e^{-t/2}y_n,$$

and

$$x_0 \rightarrow e^{t/2}x_0, \quad y_0 \rightarrow e^{-t/2}y_0,$$

it is clear that $g_t\phi_n \rightarrow g_t\phi_0$.

If $\phi_0 \equiv 0$ on a component, then for any local parameter z_0 on a compact set, considered via the conformal embedding as a local parameter for X_n , n large, we have $\phi_n(z_0) \rightarrow 0$. Then $z_n(z_0) \rightarrow 0$ as well. This implies that $g_t\phi_n \rightarrow 0$ too. \square

As we have remarked before, for a given genus g and $n = 0$ say, a stratum $QD(\sigma)$ is not a closed subset of the entire space of holomorphic quadratic differentials on compact surfaces of that genus, unless the stratum is defined by every quadratic differential having a single zero. This is because a sequence in $QD(\sigma)$ may collapse a pair of lower order zeroes into a higher order zero. In addition, the degeneration of Riemann surfaces allows limits of quadratic differentials on which the topology of the surface has changed. Because a degenerated surface need not be connected, in the following definition we allow disconnected surfaces and corresponding strata.

Definition 4.7 A stratum $QD(\sigma')$ is a degeneration of a stratum $QD(\sigma)$ if $QD(\sigma') \neq QD(\sigma)$ and there exists a sequence $\phi_n \in QD(\sigma)$ on Riemann surfaces $X_n \in \bar{R}_{g,0}$ such that $X_n \rightarrow X_0$ in $\bar{R}_{g,0}$ and ϕ_n converges uniformly on compact sets to some $\phi_\infty \in QD(\sigma')$ where ϕ_∞ is not identically zero.

For any quadratic differential ϕ which is not identically zero, denote by $l(\phi)$ the length of the shortest saddle connection of ϕ . Notice that $l(\cdot)$ is not continuous under degeneration within a fixed genus since zeroes are collapsed to higher order zeroes. That is to say, suppose $QD(\sigma')$ is a degeneration of $QD(\sigma)$ within the same genus and suppose $\phi_n \in QD(\sigma)$ converges to a nonzero $\phi_\infty \in QD(\sigma')$. Then $l(\phi_n) \rightarrow 0$ while $l(\phi_\infty) > 0$.

Now suppose a nonzero quadratic differential ϕ_0 on a connected surface is given. For any closed interval $I \subset [0, 2\pi]$ let

$$E_I(\phi_0) = \{\phi : \exists \theta_n \rightarrow \theta \in I, \exists t_n \rightarrow \infty, \exists \psi_n : r_{\psi_n} g_{t_n} r_{\theta_n} \phi_0 \rightarrow \phi\}.$$

That is to say, E_I is the set of all ω limit points ϕ in the $SL(2, \mathbf{R})$ orbit of ϕ_0 where the limiting initial rotation angle lies in I . Here the limiting ϕ may belong to any stratum and may be identically zero. Note that θ_n need not lie in I . Now Theorem 4.1 will follow from the next proposition.

Proposition 4.8 For any ϕ_0 , $E_I(\phi_0)$ contains an exceptional ϕ .

Proof. Since I is closed, $E_I = E_I(\phi_0)$ is compact. Assume E_I does not contain any exceptional ϕ ; we will arrive at a contradiction. Let

$$E_I^T = \{\phi \in E_I : \forall \phi' \in E_I, QD(\sigma'(\phi')) \text{ is not a degeneration of } QD(\sigma(\phi))\}.$$

That is, E_I^T is the set of $\phi \in E_I$ which belong to a stratum that cannot be further degenerated within E_I . Let

$$l = \inf_{\phi \in E_I^T} l(\phi).$$

If $l = 0$ choose a sequence $\phi_n \in E_I^T$ with $\lim l(\phi_n) \rightarrow 0$. By passing to a subsequence we may assume that $\phi_n \rightarrow \phi_\infty \in E_I$. Since $l = 0$ then either the Riemann surface of ϕ_n has degenerated or some set of zeroes of ϕ_n have been collapsed to a higher order zero of ϕ_∞ . In either case ϕ_∞ belongs to a stratum which is a degeneration of the stratum containing ϕ_n , a contradiction to the definition of E_I^T .

Suppose on the other hand that $l > 0$. Choose $t_n \rightarrow \infty$, and a sequence θ_n which converges to some $\theta_1 \in I$ such that

$$g_{t_n} r_{\theta_n} \phi_0 \rightarrow \phi_1 \in E_I^T.$$

Then $l(\phi_1) \geq l$. Let γ be a saddle connection of ϕ_1 whose length is $l(\phi_1)$. Let θ_γ be such that the vector associated to γ is vertical with respect to the structure of $r_{\theta_\gamma} \phi_1$. As $s \rightarrow \infty$ the length of γ with respect to the metric of $g_s r_{\theta_\gamma} \phi_1$ goes to 0. Thus we can choose s such that the length of γ with respect to the metric of $g_s r_{\theta_\gamma} \phi_1$ is less than $l/2$. We will have a contradiction if we can show that $g_s r_{\theta_\gamma} \phi_1 \in E_I^T$. Since the stratum containing ϕ_1 cannot be degenerated within E_I we only need to show that $g_s r_{\theta_\gamma} \phi_1 \in E_I$. Now by Lemma 4.6

$$g_s r_{\theta_\gamma} \phi_1 = \lim_{n \rightarrow \infty} g_s r_{\theta_\gamma} g_{t_n} r_{\theta_n} \phi_0.$$

We claim that there exist sequences $\theta'_n \rightarrow 0$, $t'_n \rightarrow \infty$, and ψ_n such that

$$g_s r_{\theta_\gamma} g_{t_n} = r_{\psi_n} g_{t'_n} r_{\theta'_n}.$$

Canonically identify G with the unit tangent vectors to the upper half-plane with Id identified with the vertical vector at $i = \sqrt{-1}$. Then $SL(2, \mathbf{R})$ acts by Mobius transformations. Furthermore ∞ is an attracting fixed point of the geodesic flow g_t . Therefore for each n , there is a unique rotation $r_{-\psi_n}$ such that $r_{-\psi_n} g_s r_{\theta_\gamma} g_{t_n}$ takes the vertical vector at i to a vector based on the imaginary axis and then for some small θ'_n and large t'_n , $r_{-\theta'_n} g_{-t'_n}$ takes this vector to the vertical vector at i , proving the claim. Finally,

$$g_s r_{\theta_\gamma} \phi_1 = \lim_{n \rightarrow \infty} r_{\psi_n} g_{t'_n} r_{\theta_n + \theta'_n} \phi_0,$$

that is, $g_s r_{\theta_\gamma} \phi_1 \in E_I$, and we are done. \square

Proof of Theorem 4.1. Choose a closed subinterval I' of I contained in the interior of I . Then Proposition 4.8 says that there must be some exceptional $\phi \in E_{I'}$. By Proposition 4.5, for a sequence $t_n \rightarrow \infty$ and $\theta_n \rightarrow \theta_0 \in I'$, $g_{t_n} r_{\theta_n} \phi_0$ has a metric cylinder in the homotopy class of some β_n and the length is uniformly bounded above by some M . Since $t_n \rightarrow \infty$ and g_{t_n} expands horizontal lengths, the horizontal length of β_n in the metric of $r_{\theta_n} \phi_0$ goes to 0.

Thus for some $\theta'_n \rightarrow 0$, the quadratic differential $r_{\theta'_n} r_{\theta_n} \phi_0$ has a metric cylinder in the class of β_n and since $\theta_n + \theta'_n \in I$, for large n , we are done. \square

Proof of Theorem 4.3. Let $T_1(S) = \{(p, \theta)\}$ denote the unit tangent space to S ; here $p \in S$ and $\theta \in \mathbf{S}^1$. Let $\epsilon > 0$. We begin by fixing a uniquely ergodic direction $\theta \in \mathbf{S}^1$, provided by theorem 3.5. Unique ergodicity implies that we can choose N large enough so that for all $x \in S$ the leaf through x of length N is $\epsilon/2$ -dense in S .

There are only a finite number of saddle connections of length less than or equal to $2N$ and so their tangent vectors determine a finite set of directions. This means that we can choose δ sufficiently small so that if θ' satisfies $|\theta - \theta'| < \delta$, then any saddle connection in direction θ' has length at least $2N$. We may also choose $\delta < \epsilon/2N$. From Theorem 4.1 we know that there is a point (x_0, θ_0) on a closed leaf L_{θ_0} in a metric cylinder such that $|\theta_0 - \theta| < \delta$. Let L_θ be the leaf through x_0 in direction θ extended to length N in each direction, or to a zero. Now we claim that in at least one of the two possible directions, moving distance N , the endpoint of L_θ and the endpoint of L_{θ_0} are within $\epsilon/2$ of each other. We prove the claim. For each direction there is a segment joining the endpoints of L_θ and L_{θ_0} such that the segment, L_θ , and L_{θ_0} form a simply connected domain. If both simply connected domains contained a zero in their closure, there would be a saddle connection γ of length at most $2N$ in direction θ' satisfying $|\theta' - \theta| < \delta$, contradicting the choice of δ . Thus in one direction the simply connected domain is actually a metric triangle with no interior zeroes. The choice of $\delta < \epsilon/2N$ implies now that the endpoints are within $\epsilon/2$ of each other, proving the claim.

Since the segment of L_θ is $\epsilon/2$ -dense in S , the leaf L_{θ_0} is ϵ -dense. This can be done for every uniquely ergodic direction θ . The set of vectors determining uniquely ergodic directions is dense in the whole phase space T_1S . Since ϵ is arbitrary this completes the proof. \square

4.2 Counting saddle connections and maximal cylinders

In this section we discuss some results on the asymptotics of the number of saddle connections and periodic orbits. No proofs will be given. As noted in Section 1.4, the number of (parallel families of) periodic orbits of length at most L grows asymptotically as $\pi L^2/2\zeta(2)$. For a general rational billiard or a flat structure, the asymptotics are not known. However we do have quadratic upper and lower bounds. To be specific, given a flat structure ϕ and $L > 0$, let $N_1(\phi, L)$ be the number of saddle connections of length at most L on ϕ and $N_2(\phi, L)$ the number of maximal metric cylinders of length at most L . The following theorem was proved in [53] and [54].

Theorem 4.9 *For each ϕ there exist positive constants $c_1 = c_1(\phi)$ and $c_2 = c_2(\phi)$ such that*

$$c_2 L^2 \leq N_2(\phi, L) \leq N_1(\phi, L) \leq c_1 L^2.$$

Other proofs of the upper estimate are given in [75] and [16]. It is possible to prove precise asymptotics for generic flat structures.

Theorem 4.10 ([16]) *For each component of $QD(\sigma)$ there exist constants c and s such that for μ_0 a.e. $\phi \in QD(\sigma)$ one has:*

$$\lim_{L \rightarrow \infty} \frac{N_1(\phi, L)}{L^2} = c$$

and

$$\lim_{L \rightarrow \infty} \frac{N_2(\phi, L)}{L^2} = s.$$

In the next section we discuss some examples in which the quadratic asymptotics are established and the constants are computable.

5 Veech groups and Veech surfaces

5.1 Definition and examples of Veech groups

In this section we discuss the groups of affine transformations associated to flat structures. These groups give rise to subgroups in $SL(2, \mathbf{R})$. If this subgroup is a lattice in $SL(2, \mathbf{R})$, then the flat structure has particularly nice properties.

Let ϕ be a quadratic differential on the Riemann surface X . Let Σ be a finite subset of X which contains the zeroes of ϕ (but may be larger). Let $X' = X \setminus \Sigma$. Following Veech we give the next definition.

Definition 5.1 *The affine group $Aff^+(\phi)$ is the group of orientation preserving self homeomorphisms of X that map X' to itself and are affine with respect to the natural coordinates of ϕ .*

This means that for each point $p \in X'$ we choose coordinates z near p so that $\phi = dz^2$ and coordinates w near $f(p)$ so that $\phi = dw^2$, and in these coordinates, the map f is affine. We call the derivative (i.e., linear part) of that homeomorphism the

derivative of f . It is well-defined independently of p , z and w up to a factor of $\pm I$. Let $a(f) = \pm A$ be the derivative. Then by an observation of Veech [71] we have

$$|f^*\phi| = \det(A)|\phi|$$

and

$$\int |\phi| = \int |f^*(\phi)|,$$

which implies that the determinant is 1 and so A is an element of $G = SL(2, \mathbf{R})$.

Definition 5.2 *The Veech group $\Gamma(\phi)$ is the image of $Aff^+(\phi)$ under the derivative map.*

The case of most interest is when Γ is a lattice in G . Recall that a lattice Γ in G is a discrete subgroup such that G/Γ has finite volume; a lattice is nonuniform if the quotient is not compact (see, e.g., [5]).

We present several examples.

Thurston example. The following example is essentially due to Thurston [66]. Take the unit square with lower left vertex at the origin and identify opposite vertical sides. Mark off 4 points on both the top and bottom of the resulting cylinder C , dividing the top and bottom into equally spaced intervals and so that the points on top lie directly above the points on the bottom. Arrange the intervals on top into two pairs and glue each pair isometrically as indicated in the figure. Do the same for the bottom.

Fig. 10

The Euclidean metric on the rectangle extends to a quadratic differential $\phi = dz^2$ on the glued surface X , which has genus 2. The 4 marked points on the top identify to a single zero of order 2; similarly for the bottom. Let Σ consist of these two points.

A complete proof that $Aff^+(\phi)$ is a lattice (for a very similar example) can be found in [15]. Here we will just give an outline of that proof. The meromorphic function $F(z) = \wp(z)/\wp(1)$, where $\wp(z)$ is the Weierstrass function, has fundamental periods of 2 and $2i$ and gives a double covering of the sphere branched over 1, 0, -1 and ∞ . Because $F(z) = F(-z)$, it induces a meromorphic covering map f of the identification space X onto the sphere, branched over the same points.

The quadratic differential ϕ is a real multiple of the pull-back under f of the quadratic differential $\psi = d\zeta^2/\zeta(\zeta^2 - 1)$ on the sphere. Now ψ has closed horizontal

trajectories. The affine group $Aff^+(\psi)$ contains $PSL(2, \mathbf{R})$ as a finite index subgroup. A finite index subgroup H_0 of $PSL(2, \mathbf{R})$ acting on the sphere has lifts under f to maps of X and this set of lifts \tilde{H}_0 forms a finite index subgroup of $Aff^+(\phi)$. Thus $Aff^+(\phi)$ is a lattice.

We can describe two elements of $Aff^+(\phi)$. Recall that a (right) Dehn twist about a simple closed curve α on a surface is a homeomorphism of the surface which fixes α and twists any curve β crossing α once to the right about α . The map $(x, y) \rightarrow (x + y, y)$ is a Dehn twist about the core curve of C and has derivative

$$A_1 = \begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix}.$$

To find a second affine map notice that the gluings force all trajectories in the vertical direction to be closed as well. Indeed, an upward vertical trajectory hits the top horizontal side of the square at some point S and continues downwards from the point T of the same horizontal side which is identified with S . The surface X decomposes into vertical metric cylinders C_1, C_2 of circumference 2 and height $1/4$, each composed of 2 vertical strips glued together. This allows us to define a right Dehn twist on each C_i with derivative

$$A_2 = \begin{pmatrix} 1 & 0 \\ -8 & 1 \end{pmatrix}.$$

The elements A_1 and A_2 generate an infinite index subgroup of $Aff^+(\phi)$. Thurston showed that the subgroup generated by them contains pseudo-Anosov elements.

Definition 5.3 *Two subgroups Γ_1, Γ_2 of G are commensurable if there is some $g \in G$ such that $g\Gamma_1g^{-1} \cap \Gamma_2$ has finite index in both $g\Gamma_1g^{-1}$ and Γ_2 .*

The group constructed above is commensurable to $SL(2, Z)$. Gutkin and Judge [31] showed that the Veech group is commensurable with $SL(2, Z)$ if and only if the surface can be tiled by a Euclidean parallelogram.

Veech examples. The next set of examples are due to Veech – see [71], [72]. Part of the motivation was to construct an example of a non-arithmetic Veech group, that is, a lattice, not commensurable to $SL(2, Z)$. These examples are the flat structures ϕ_n associated with billiards in right triangles with angles $\pi/n, n \geq 5$ and the set Σ consists of the zeroes of ϕ_n . Veech proved the following theorem.

Recall that the (p, q, r) triangle group is the index 2 subgroup of the group generated by reflections in the sides of a triangle with angles $\pi/p, \pi/q, \pi/r$ in the hyperbolic

plane; the triangle group is a subgroup of G , the group of isometries of the hyperbolic plane. If an angle is 0, then the vertex of the triangle is at infinity (see, e.g., [5]).

Theorem 5.4 *For n odd, $\Gamma(\phi_n)$ is a $(2, n, \infty)$ triangle group. For even $n = 2m$, $\Gamma(\phi_n)$ is the (m, ∞, ∞) triangle group. In either case the Riemann surface X is the surface associated to $y^2 + x^n = 1$.*

One can check that for n odd the invariant surfaces for the above mentioned triangle and the isosceles triangle with equal angles of π/n are isomorphic so that they determine the same Veech group.

We will not give a proof of the theorem but illustrate the main idea by an example. It is convenient to discuss $\Gamma(\phi_8)$ as this invariant surface was already described in figure 4. The flat surface built out of the triangle is the regular octagon P with opposite sides identified. The surface X has genus 2 and the quadratic differential ϕ_8 has a single zero of order 4 corresponding to the identification of all vertices of P . The counterclockwise rotation f_1 of order 8 about the center of P is obviously affine with respect to ϕ_8 and is holomorphic with respect to X . It fixes the zero. Its derivative is

$$A_1 = \begin{pmatrix} \sqrt{2}/2 & -\sqrt{2}/2 \\ \sqrt{2}/2 & \sqrt{2}/2 \end{pmatrix}.$$

To find a second affine map, rotate the octagon so that two of the opposite sides are vertical. Then every horizontal trajectory is closed and that the surface decomposes into two metric cylinders. The first cylinder C_1 consists of horizontal segments joining the opposite vertical sides. Its boundary consists of the two horizontal segments joining opposite vertices. Next notice that because of the identifications made, each horizontal segment leaving one of the nonvertical sides is also closed; it first hits the other nonvertical side before closing. The family of these closed trajectories determine the second metric cylinder C_2 . Its boundary includes the first boundary together with an extra closed saddle connection which is the top and bottom horizontal segments of the octagon. If we normalize so that the sides of P are of length 1, then C_1 has circumference $\cot \pi/8 = 1 + \sqrt{2}$ and height 1, while C_2 has circumference $1 + \cot \pi/8 = 2 + \sqrt{2}$ and height $\sqrt{2}/2$. Notice that the ratio of height to circumference for the first cylinder is twice that for the second one. The ratio of height to circumference is called the *modulus* of a metric cylinder.

Fig. 11

Then there exists an affine map which preserves each cylinder fixing the boundary pointwise and whose derivative is

$$A_2 = \begin{pmatrix} 1 & 2 \cot \pi/8 \\ 0 & 1 \end{pmatrix}.$$

The map is the square of the Dehn twist on C_1 and the Dehn twist on C_2 . Now A_1, A_2 generate a $(8, \infty, \infty)$ triangle group which is not commensurable with $SL(2, Z)$ ([5]).

The following theorem combines results from [15],[76] and [78], and complements Theorem 5.4.

Theorem 5.5 *The Veech groups of the flat surfaces associated with the three series of triangles with angles*

$$(\pi/2n, \pi/n, (2n-3)\pi/2n), \quad ((2k-1)\pi/4k, (2k-1)\pi/4k, \pi/2k)$$

and

$$(k\pi/(2k+1), k\pi/(2k+1), \pi/(2k+1)), \quad n \geq 4, k \geq 2,$$

are, respectively,

$$(3, n, \infty), \quad (2k, \infty, \infty) \quad \text{and} \quad (2k+1, \infty, \infty)$$

triangle groups.

5.2 Veech dichotomy

We now discuss some properties of $\Gamma(\phi)$. These results are all due to Veech – [71], [72].

Lemma 5.6 *Let $f \in Aff^+(\phi)$ be such that $A = a(f) \neq \pm Id$, and let $\zeta \in \mathbf{R}^2 \setminus \{0\}$ be such that $A\zeta = \pm\zeta$. Then every leaf in the direction corresponding to ζ is either a saddle connection or closed so that the surface decomposes into the cylinders C_1, \dots, C_p in direction ζ . Moreover there exists a positive integer k such that for each i the map f^k preserves C_i , fixing the boundary, and is a power of the Dehn twist of C_i .*

Proof. Clearly f fixes the foliation in direction ζ . For some k , that can be assumed even, f^k maps each separatrix to itself. Since $A\zeta = \pm\zeta$, f^k must be the identity on

each separatrix. Each separatrix is either a saddle connection or dense in an open set. In the latter case, since f^k is real analytic, f^k would be the identity in the open set, and so $A^k = \pm Id$. However A is a unipotent, therefore $A^k \neq \pm Id$. Thus every separatrix is a saddle connection which means that the surface decomposes into cylinders such that f^k is the identity on the boundary of each cylinder. As it is affine on each cylinder, it must be a power of the Dehn twist. \square

Veech also shows that the moduli of the cylinders are rationally related.

Lemma 5.7 *If $a(f) \neq \pm Id$ then f is not isotopic to the identity.*

Proof. Since f is affine in the natural coordinates determined by ϕ , it is in fact a Teichmuller map determined by ϕ and some dilation. Teichmuller maps are extremal quasiconformal maps within their homotopy class. This means that f cannot be homotopic to the identity, since the identity is conformal. \square

Corollary 5.8 *$Aff^+(\phi)$ is a subgroup of the mapping class group of the Riemann surface.*

Lemma 5.9 *$\Gamma(\phi)$ is a discrete subgroup of G .*

Proof. Suppose $f_n \in Aff^+(\phi)$ is such that $a(f_n) \rightarrow \pm Id$. By passing to a subsequence and using the Arzela-Ascoli theorem we can assume that f_n converges uniformly to a diffeomorphism f . For large m , $f_m \circ f_{m+1}^{-1}$ is isotopic to the identity, which by Lemma 5.7 gives $a(f_{m+1}) = a(f_m)$ and so $a(f_m) = \pm Id$ for large m . \square

We now consider consequences of the assumption that $\Gamma(\phi)$ is a lattice. Begin again by canonically identifying G with the set of unit tangent vectors to the upper half-plane with e , the identity in G identified with the vertical vector at $i = \sqrt{-1}$. Then $SL(2, \mathbf{R})$ acts on itself by Mobius transformations. Recall that $SL(2, \mathbf{R})$ also acts on quadratic differentials. If we associate the identity in G to ϕ , then we canonically identify the $SL(2, \mathbf{R})$ orbit of ϕ with $SL(2, \mathbf{R})$.

Theorem 5.10 *(Veech dichotomy). Suppose $\Gamma(\phi)$ is a lattice. For any direction either the foliation in that direction is minimal or its every leaf is closed or a saddle connection. In the latter case the surface decomposes into cylinders of closed leaves and there is $f \in Aff^+(\phi)$ with $a(f) \neq I$ which preserves these cylinders and acts as Dehn twists on them. Furthermore Γ is a nonuniform lattice.*

Proof. If the foliation is not minimal there is a saddle connection in that direction. By rotating, we can assume the direction is vertical. Then $g_t\phi$ has a saddle connection which approaches 0 in length, which, under the identification, says that $g_t(e)$ leaves every compact set of G/Γ , where $e \in G$ is the identity element. Since Γ is a lattice this can only happen if the point at infinity is a fixed point for some parabolic element of Γ . In particular, the lattice is not uniform. Such a parabolic element fixes the vertical vector $\zeta = [0, 1]$. It remains to apply Lemma 5.6. \square

A different proof of the above theorem can be found in [76].

Theorem 5.10 says that, just as in the case of the flat torus, there is a dichotomy: either the flow is minimal or every orbit is closed. Since the triangle groups are lattices, the Veech dichotomy applies, in view of Theorem 5.4, to billiards in the isosceles triangles with equal angle π/n ; it applies to regular polygons as well.

In fact, the analogy with the flat torus extends to ergodicity as well.

Theorem 5.11 *Suppose $\Gamma(\phi)$ is a lattice. If the foliation in a direction is minimal, then it is uniquely ergodic.*

Proof. Again we can assume the direction is vertical. If the foliation is not uniquely ergodic then by Theorem 3.8 $g_t\phi$ eventually leaves every compact set of moduli space. Again this says that $g_t e$ leaves every compact set of G/Γ and that Γ has a parabolic element fixing infinity. By Lemma 5.6 the leaves in the vertical direction are all closed, which contradicts the minimality. \square

5.3 Asymptotics

Let ϕ be a quadratic differential. Recall that we defined $N_2(\phi, T)$ as the number of homotopy classes of closed geodesics of length less than T or equivalently, the number of maximal cylinders whose circumference has length at most T . By Theorem 4.9 there are upper and lower quadratic bounds on $N_2(\phi, T)$ for any ϕ . Veech (op. cit.) has shown the following.

Theorem 5.12 *Suppose $\Gamma(\phi)$ is a lattice. Then there exists a constant $c = c(\phi)$ such that $N_2(T) \sim cT^2$.*

The proof is using the Eisenstein series of $\Gamma(\phi)$. The constant c has been computed in the case of isosceles triangles T_n with angles $(\pi/n, \pi/n, (n-2)\pi/n)$. In that case

$$c = \frac{n(n^2 - 1)}{48\pi(n - 2)|T_n|},$$

where $|T_n|$ is the area of the triangle.

Gutkin and Judge [31] have given another proof of the same theorem using the mixing of the horocyclic flow. The constant c is given in their work in geometrical terms.

5.4 Covers

Constructing a flat surface M from a rational billiard polygon P it may happen that some of the cone angles are equal to 2π or, equivalently, the points are not zeroes of the quadratic differential. For example, if P is a square then M is a flat torus with no singular points. From the billiard view point this means that one can continuously define the extension of the billiard trajectory through the vertex of P ; this is possible if and only if the angle is of the form π/n . Adding such a point to Σ , the set of points that are required to be preserved by the affine diffeomorphisms, may change the Veech group. It is particularly important to keep track of removable singular points in the study of covers of flat structures.

Definition 5.13 *Let $(X_i, \phi_i), i = 1, 2$ be flat structures and $\Sigma_i \subset X_i, i = 1, 2$ sets that contain the zeroes of X_i . Then (X_1, ϕ_1, Σ_1) covers (X_2, ϕ_2, Σ_2) if there exists a continuous map $f : X_1 \rightarrow X_2$, called a cover, that sends Σ_1 to Σ_2 and which is given, in local coordinates on the complement of these sets, by parallel translation. The multiplicity of a cover is the number of preimages of a non-singular point in X_2 (independent of the point).*

In particular, a cover defines a holomorphic cover of the Riemann surface X_2 by the Riemann surface X_1 , branched over Σ_2 .

Definition 5.14 *Two covers $f_1 : X_1 \rightarrow Y$ and $f_2 : X_2 \rightarrow Y$ are called isomorphic if there exists an isomorphism of the flat structures $g : X_1 \rightarrow X_2$ such that $f_1 = f_2 g$. Similarly one defines an isomorphism of covers $f_1 : Y \rightarrow X_1$ and $f_2 : Y \rightarrow X_2$.*

The following finiteness property holds (see [76] or [30, 31] for a proof).

Proposition 5.15 *A given flat structure admits a finite number of isomorphic classes of covers with a given degree. Likewise, the number of isomorphic classes of covers realizable by a given flat structure is finite.*

The next construction provides examples of covers of flat structures associated with billiards in polygons. Let P_1 and P_2 be rational polygons, and assume that P_2 tiles P_1 by reflections. This means that P_1 is partitioned into a number of isometric copies of P_2 , each two either disjoint, or having a common vertex or a common side, and if two of these polygons have a side in common then they are symmetric with respect to this side.

For example, a right triangle with angle π/n tiles by reflections a regular n -gon – see figure 4 for $n = 8$. Another example is given by the above mentioned Gutkin's almost integrable billiard polygons: if such a polygon is drawn on the integer square lattice in the plane, then it is tiled by the unit square.

Let (X_i, ϕ_i) be the flat structures associated with the polygon P_i , $i = 1, 2$, and assume that P_2 tiles P_1 . Let Σ_1 be the set of zeroes of ϕ_1 (cone points with angles greater than 2π).

Lemma 5.16 *There is a set Σ_2 containing the zeroes of ϕ_2 such that (X_1, ϕ_1, Σ_1) covers (X_2, ϕ_2, Σ_2) .*

Proof. Denote by G_1 and G_2 the groups generated by the linear parts of the reflections in the sides of the polygons P_1 and P_2 . Then $G_1 \subset G_2$. Each copy of P_2 , involved in the tiling of P_1 , is identified with P_2 by an isometry; this isometry is a composition of reflections in the sides of P_2 . These identifications, combined, define a projection p from P_1 to P_2 . Given a point $x \in P_1$, let $\alpha(x) \in G_2$ be the linear part of the isometry that takes P_2 to the tile (a copy of P_2) that contains x .

Consider the map $g : P_1 \times G_1 \rightarrow P_2 \times G_2$ that sends (x, β) to $(p(x), \beta\alpha(x))$. Notice that the surfaces X_i are quotient spaces of $P_i \times G_i$, $i = 1, 2$, and the map g determines a map $f : X_1 \rightarrow X_2$. This is the desired cover. The multiplicity equals n/m where n is the number of tiles in the tiling and $m = [G_2 : G_1]$.

Some singular points of the first structure may project to removable singular points of the second one. This happens only if P_2 has a vertex angle π/n . One then adds these points to Σ_2 . \square

Note that one does not need to add removable singular points if none of the angles of P_2 is of the form π/k .

As an example to Lemma 5.16 consider again a regular octagon P_1 tiled by right triangles P_2 with angle $\pi/8$. The flat surface M_2 for the triangle is the regular octagon whose opposite sides are identified, while the flat surface M_1 for the regular octagon is a union of 8 octagons with some gluings of the sides. Thus M_1 covers M_2 with multiplicity 8.

The next result relates the Veech groups of flat structures one of which covers the other; it was obtained independently by Gutkin and Judge [30], [31] and by Vorobets [76]. Recall that two subgroups of G are commensurate if they share a common subgroup that has a finite index in each.

Theorem 5.17 *If (X_1, ϕ_1, Σ_1) covers (X_2, ϕ_2, Σ_2) then $\Gamma(\phi_1)$ and $\Gamma(\phi_2)$ are commensurate. In particular, one of the groups $\Gamma(\phi_1)$ and $\Gamma(\phi_2)$ is a lattice if and only if the other is.*

Proof. Consider the set of triples $S = \{(X, \phi, f)\}$ where (X, ϕ) is a flat structure and $f : X_1 \rightarrow X$ is a cover. Let \bar{S} be the set of equivalence classes of such triples considered up to isomorphism of covers.

Let $g \in Aff^+(\phi_1)$ and let $A = a(g)$ be its derivative. Then fg is a cover of the flat structure $(X, A^{-1}(\phi))$ by (X_1, ϕ_1) . This defines a right action of $Aff^+(\phi_1)$ on S , and this action descends to an action on \bar{S} . Consider the subgroup $Aff_0^+(\phi_1) \subset Aff^+(\phi_1)$ that consists of the affine transformations acting trivially on \bar{S} . By Proposition 5.15 \bar{S} is finite, therefore this subgroup has finite index. Let $\Gamma_0 \subset \Gamma(\phi_1)$ consist of the derivatives of the elements of $Aff_0^+(\phi_1)$; then Γ_0 has finite index too.

Let $g \in Aff_0^+(\phi_1)$ and $A = a(g)$. Then, for every $(X, \phi, f) \in S$ the flat structures (X, ϕ) and $(X, A^{-1}\phi)$ are isomorphic. Therefore $A \in \Gamma(\phi)$. Since (X_1, ϕ_1) covers (X_2, ϕ_2) , it follows that $\Gamma_0 \subset \Gamma(\phi_2)$. Thus $\Gamma_0 \subset \Gamma(\phi_1) \cap \Gamma(\phi_2)$ and has finite index in $\Gamma(\phi_1)$. To show that Γ_0 has finite index in $\Gamma(\phi_2)$ as well one applies a similar argument to the set of covers of (X_2, ϕ_2) . \square

Corollary 5.18 ([35]). *If a rational polygon P_2 has no angles of the form π/n and P_2 tiles P_1 by reflections then the respective Veech groups are commensurable.*

Hubert and Schmidt [35] have shown that the Veech group associated to the isosceles triangle with angles $2\pi/n, (n-2)\pi/2n, (n-2)\pi/2n$ is not a lattice even though the right triangle with angle π/n tiles it with one reflection. This is due to the appearance of removable singular points that we discussed above.

6 Interval exchange transformations

6.1 Topological structure of orbits

In this section we discuss the topological structure of orbits of an interval exchange transformation – see Definition 1.11; we refer to [39] for a detailed exposition; see also

[43] and [13]. Without loss of generality, we may assume that end-points of all the intervals involved are discontinuity points of the interval exchange transformation T involved; otherwise adjacent segments could be joined into one and T would be an exchange of a smaller number of intervals.

The following definitions provide analogs of saddle connection and metric cylinder in the present setting.

Definition 6.1 *An orbit segment $(x, Tx, \dots, T^{k-1}x)$ is called a connecting segment of the interval exchange transformation T if the points x and $T^{k-1}x$ are end-points of some intervals I_i and I_j but none of the intermediate points $Tx, \dots, T^{k-2}x$ is an end-point of any of the intervals I_1, \dots, I_n .*

Definition 6.2 *An interval $I = [a, b) \subset [0, 1)$ is called rigid if every positive iterate of T is continuous on I .*

If x is a periodic point of an interval exchange transformation T then there exists maximal rigid interval containing x and consisting of periodic points with the same period. The end-points of this interval belongs to connecting segments.

Definition 6.3 *An interval exchange transformation is called generic if it has no rigid intervals.*

One has the following criterion for genericity. Denote by λ the Lebesgue measure on $[0, 1)$.

Lemma 6.4 *The following properties are equivalent:*

- (i) T is generic;
- (ii) T does not have periodic orbits;
- (iii) $\lim_{k \rightarrow \infty} \sup_{i_0, i_1, \dots, i_k} \lambda(I_{i_0} \cap TI_{i_1} \cap \dots \cap T^k I_{i_k}) = 0$.

Proof. We already mentioned that if T has a periodic point then it has a rigid interval. Thus (i) implies (ii). If T has a rigid interval then (iii) clearly does not hold. Thus (iii) implies (i).

It remains to deduce (iii) from (ii). Let E be the union of the orbits of the left end-points of the intervals I_1, \dots, I_n . Then (iii) is equivalent to E being dense. Assume that (iii) does not hold. Then the set $[0, 1) - \bar{E}$ is a non-empty open set. Each component of this set is an open interval, and T exchanges these components. Since there are only finitely many components of a given length, each component is periodic, and (ii) does not hold. \square

Definition 6.5 *A point x is called generic if it is a continuity point for all (positive and negative) iterates of T .*

One can prove (see [39]) that a generic point is either periodic or the closure of its orbit is a finite union of intervals. The latter union is called a *transitive component* of the interval exchange transformation. We refer to [39] for the following structural result.

Theorem 6.6 *Let T be an exchange of n intervals. Then the interval $[0, 1)$ splits into a finite union of connecting segments and at most $2n - 2$ disjoint open invariant sets, each of which is either a transitive or a periodic component, and each of these components is a finite union of open intervals.*

6.2 Number of invariant measures. Lack of mixing

The first result of this section concerns the number of invariant measures for an interval exchange transformation. Assume that T is a generic exchange of n intervals. The idea of the proof of the next result belongs to V. Oseledets. We follow an elegant exposition in [39]; see also [13] and [63].

Theorem 6.7 *There exist at most n mutually singular T -invariant Borel probability measures.*

Proof. Denote by ξ the partition of $[0, 1)$ into the intervals I_1, \dots, I_n . Due to condition (iii) in Lemma 6.4, ξ is a one-sided generator for T . Therefore an invariant measure μ is determined by its values on the elements of the partitions

$$\xi_k = \xi \vee T\xi \vee \dots \vee T^{k-1}\xi, \quad k = 1, 2, \dots$$

Claim: μ is determined by its values on I_1, \dots, I_n .

To prove this claim, consider two finite partitions η and ν of the interval $[0, 1)$ into subintervals, and let μ be a nonatomic Borel probability measure on $[0, 1)$. Then the values of μ on the elements of the partition $\eta \vee \nu$ are uniquely determined. Indeed, let $[a, b)$ be an element of $\eta \vee \nu$. Then a is the end-point of an interval from either η or ν . Therefore $\mu([0, a))$ is determined. The same applies to $\mu([0, b))$, thus $\mu([a, b))$ is determined as well.

Now the italicized claim above is proved inductively: setting $\eta = \xi_k, \nu = T^k\xi$, one has: $\eta \vee \nu = \xi_{k+1}$, and the claim follows.

Assign a point of an $(n - 1)$ -dimensional simplex

$$\Delta = \{(x_1, \dots, x_n) : \sum x_j = 1\}$$

to a T -invariant measure μ as follows: $x_j = \mu(I_j)$. We obtain a map F from the space of invariant probability measures to Δ . This map is affine, continuous in the weak * topology and, as has been shown, injective. It remains to notice that if μ_1, \dots, μ_k are mutually singular measures then $F(\mu_1), \dots, F(\mu_k)$ are linearly independent. Indeed, assume that a relation holds:

$$a_1 F(\mu_1) + \dots + a_k F(\mu_k) = F(a_1 \mu_1 + \dots + a_k \mu_k) = 0. \quad (6.1)$$

For every $i = 1, \dots, k$ there is a set $U_i \subset [0, 1)$ such that $\mu_i(U) > 0$ but $\mu_j(U) = 0$ for all $j \neq i$. Thus $(a_1 \mu_1 + \dots + a_k \mu_k)(U_i) = a_i \mu_i(U_i)$, and (6.1) implies that $a_i = 0$. \square

The estimate can be improved if we add the assumption of irreducibility of the permutation. Recall, an exchange of n intervals (I_1, \dots, I_n) is determined by the pair (\bar{i}, σ) where the vector $\bar{i} = (i_1, \dots, i_n)$ consists of the lengths $i_j = |I_j|$ and $\sigma \in S_n$ is the permutation corresponding to T . We say that σ is *irreducible*, if for no $k < n$ one has $\sigma\{1, \dots, k\} = \{1, \dots, k\}$. W. Veech proved the following theorem [74].

Theorem 6.8 *Assume that T is an interval exchange with irreducible permutation. Then the number of mutually singular T -invariant Borel probability measures does not exceed $n/2$.*

A similar bound was found by Katok [36] in the context of flows on surfaces.

Next we discuss mixing properties of interval exchange transformations.

Definition 6.9 *A measure-preserving transformation $T : (I, \mu) \rightarrow (I, \mu)$ is called mixing if for every measurable sets A and B one has:*

$$\lim_{n \rightarrow \infty} \mu(T^{-n}(A) \cap B) = \mu(A)\mu(B).$$

The following result is due to A. Katok [37], see also [13].

Theorem 6.10 *No interval exchange transformation is mixing with respect to any invariant Borel probability measure μ .*

The proof makes use of the next lemma which also has other applications in the study of interval exchange transformations. Let T be an exchange of n intervals. Consider an interval $J = [a, b) \subset [0, 1)$, and denote by $T_J : J \rightarrow J$ the first return map.

Lemma 6.11 *The map T_J is an exchange of at most $n + 2$ intervals.*

Proof. By the Poincaré recurrence theorem, almost every point of J returns to J . Denote by Σ the union of the discontinuity set of T and the end-points of J ; this set contains at most $n + 1$ points. Consider the set

$$\Omega = \{x \in J : T_J(x) = T^k(x) \quad \text{and} \quad T^l(x) \notin \Sigma \quad \text{for} \quad l = 0, \dots, k - 1; \quad k = 1, 2, \dots\}.$$

Then Ω is an open dense subset of J .

Consider a maximal interval $J_\alpha \subset \Omega$ and let y be its left end-point. The restriction of T_J on J_α is T^k , for some k . Then an iterate $T^l(y)$, $l = 0, \dots, k$, belongs to the set Σ ; otherwise J_α can be extended beyond y . Each of the points of Σ can appear in this way as an iterate of at most one left end-point of a maximal interval J_α . Thus Ω is the complement of at most $n + 1$ points, and T_J is an exchange of at most $n + 2$ intervals. \square

Proof of Theorem 6.10. It is enough to consider an ergodic measure μ . Such a measure is either concentrated on a finite set or non-atomic. In the former case T is not mixing. Assume that μ is the Lebesgue measure; we will show that the general case reduces to this one. Let T be an exchange of n intervals.

Fix $J \subset [0, 1)$. According to Lemma 6.11, T_J is an exchange of intervals J_1, \dots, J_s where $J = J_1 \cup \dots \cup J_s$:

$$T_J|_{J_i} = T^{k_i}, \quad i = 1, \dots, s; \quad s \leq n + 2. \quad (6.2)$$

Set: $J_i^m = T^m(J_i)$. Then, ignoring finite sets,

$$[0, 1) = \cup_{i=1}^s \cup_{m=0}^{k_i-1} J_i^m;$$

indeed, the set on the right hand side is T -invariant. Denote this partition of the interval $[0, 1)$ by ξ_J . Choosing J sufficiently small, one makes ξ_J arbitrarily fine.

Now one repeats the argument for the maps T_{J_i} , $i = 1, \dots, s$. One has:

$$J_i = \cup_{j=1}^{s_i} J_{ij}, \quad s_i \leq n + 2, \quad \text{and} \quad T_{J_i}|_{J_{ij}} = T^{k_{ij}}.$$

Set: $J_{ij}^m = T^m(J_{ij})$. Then

$$[0, 1) = \cup_{i=1}^s \cup_{j=1}^{s_i} \cup_{m=0}^{k_{ij}-1} J_{ij}^m.$$

Notice that $T^{k_{ij}}(J_{ij}^m) \subset J_i^m$, and therefore $J_{ij}^m \subset T^{-k_{ij}}(J_i^m)$. Thus

$$J_i^m = \cup_{j=1}^{s_i} J_{ij}^m \subset \cup_{j=1}^{s_i} T^{-k_{ij}}(J_i^m).$$

Let B be a set, measurable with respect to the partition ξ_J . Then

$$B \subset \cup_{i=1}^s \cup_{j=1}^{s_i} T^{-k_{ij}}(B).$$

Since T is measure-preserving and $s, s_i \leq n+2$, one has, for some k_{ij} ,

$$\mu(B \cap T^{k_{ij}}(B)) = \mu(T^{-k_{ij}}(B) \cap B) \geq \frac{1}{(n+2)^2} \mu(B). \quad (6.3)$$

Fix a set A such that

$$\mu(A) < \frac{1}{5(n+2)^2}. \quad (6.4)$$

For each positive integer N one can choose a subinterval $J \subset [0, 1)$ so small that

- (i) there exists a set B , measurable with respect to the partition ξ_J , and such that $\mu(A \Delta B) < \mu(A)^2$;
- (ii) the numbers k_i in (6.2) are all greater than N .

To satisfy (ii), one chooses a point x such that the points $T^i(x)$, $i = 0, \dots, N-1$, are points of continuity of T . Then these points are all distinct (otherwise T would have a periodic orbit and a periodic component of positive measure), and J can be taken as a sufficiently small neighborhood of x .

It follows from (i) that $\mu(B) \geq \mu(A) - \mu(A)^2$. Applying (6.3) to the set B and taking (i) and (6.4) into account, one obtains, for some $k_{ij} \geq k_i > N$,

$$\begin{aligned} \mu(A \cap T^{k_{ij}}(A)) &\geq \mu(B \cap T^{k_{ij}}(B)) - 2\mu(A \Delta B) \geq \frac{1}{(n+2)^2} \mu(B) - 2\mu(A)^2 \geq \\ &\frac{1}{(n+2)^2} \mu(A) - \frac{1}{(n+2)^2} \mu(A)^2 - 2\mu(A)^2 \geq \mu(A)^2 \left(5 - \frac{1}{(n+2)^2} - 2 \right) \geq 2\mu(A)^2. \end{aligned}$$

Since N was arbitrary, T is not mixing.

To complete the proof it remains to justify the reduction of an invariant non-atomic probability measure μ to the Lebesgue measure λ . Define a map $\phi : [0, 1] \rightarrow [0, 1]$ by the formula: $\phi(x) = \mu([0, x])$. Then ϕ is monotone, continuous, surjective and

$\phi_*\mu = \lambda$. Although ϕ is not necessarily injective, it is an isomorphism between the measure spaces $([0, 1], \mu)$ and $([0, 1], \lambda)$.

Define a map $S : [0, 1] \rightarrow [0, 1]$ by the formula: $S(x) = \phi(T(y))$ where $y \in \phi^{-1}(x)$; in other words, $S = \phi \circ T \circ \phi^{-1}$. To see that S is well defined, note that $\phi^{-1}(x)$ is either a point or an interval. In the latter case, $\phi T \phi^{-1}(x)$ is either a point or an interval. Actually, the latter holds, as follows from the next equalities:

$$\lambda(\phi T \phi^{-1}(x)) = \mu(T \phi^{-1}(x)) = \mu(\phi^{-1}(x)) = \lambda(\{x\}) = 0.$$

Thus S is well defined. Since S preserves the Lebesgue measure and orientation, and is continuous and one-to-one except for at most n points, S is an exchange of at most n intervals. \square

Using the reduction of the billiard flow on an invariant surface to an interval exchange transformation described in Section 1.7, or more generally from the flow on a flat surface, one obtains the next corollary – see [36] for details.

Corollary 6.12 *For every rational polygon (flat surface) the billiard flow on an invariant surface is not mixing.*

6.3 Ergodicity of interval exchange transformations

In this section and the next, assume that the permutation associated to the interval exchange is irreducible. This is an obvious necessary condition for an interval exchange to be uniquely ergodic. Since $i_1 + \dots + i_n = 1$, the space of irreducible interval exchange transformations is the product of the $(n-1)$ -dimensional simplex Δ^{n-1} and the set of irreducible permutations S_n^0 . Keane asked whether a typical (with respect to Lebesgue measure λ on Δ^{n-1}) interval exchange is uniquely ergodic.

Theorem 6.13 *Let σ be irreducible. For λ a.e. $\bar{i} \in \Delta^{n-1}$ the corresponding interval exchange T is uniquely ergodic.*

This theorem was first proved independently in [51] and [73]. Later entirely different proofs were given in [46], [60], and [9]. It also follows from Theorem 3.5 as noted in [47]. Our purpose here is to give a sufficient condition for unique ergodicity due to M. Boshernitzan [9]. Instead of taking the first return map to a subinterval, one iterates the interval exchange. T^k is an interval exchange of essentially kn intervals. Denote by $m(k)$ the length of the shortest interval of T^k .

One defines a subset $A \subset N = \{1, 2, \dots\}$ to be *essential* if for every $l \geq 2$ there exists $\alpha > 1$ such that the system of equations

- (1) $n_{i+1} > 2n_i, 1 \leq i \leq l-1$
- (2) $n_l \leq \alpha n_1$
- (3) $n_i \in A$

has an infinite number of solutions. The Boshernitzan criterion reads as follows.

Theorem 6.14 *If for some $\epsilon > 0$ the set $\{k : m(k) \geq \frac{\epsilon}{k}\}$ is essential then T is uniquely ergodic.*

6.4 Asymptotic flag of an interval exchange transformation

In this section we describe, without proofs, recent results obtained by A. Zorich, M. Kontsevich and G. Forni [79], [80], [81], [82], [48], [18] (these results can be also stated in the setting of measured foliations on surfaces).

Consider an interval exchange transformation T , pick a generic subinterval $J \subset [0, 1)$ and a generic point $x \in [0, 1)$. Generically, T is ergodic, therefore one has the following equality:

$$\#\{i : 0 \leq i \leq N-1, T^i(x) \in J\} = |J|N + o(N).$$

Zorich observed in computer experiments that the above error term, generically, grows as a power of N : error term $\sim O(N^\alpha)$. In other words,

$$\frac{\log \#\{i : 0 \leq i \leq N-1, T^i(x) \in J\} - |J|N}{\log N} = \alpha. \quad (6.5)$$

The exponent $\alpha < 1$ depends only on the permutation associated with the interval exchange transformation T . This observation lead to the theory discussed in this section.

Renormalization procedure for interval exchange transformations. According to Lemma 6.11, the first return map T_J to a subinterval $J \subset [0, 1)$ is again an interval exchange transformation, although generally on more than n intervals. If J can be chosen so that T_J is an exchange on n intervals, then if renormalize J to the unit length, then the procedure that assigns T_J to T determines a transformation of $\Delta^{n-1} \times S_n^0$. An example of such a procedure is provided by the *Rauzy induction* [59]. A modification of Rauzy induction was constructed in [80]; and the corresponding

transformation of $\Delta^{n-1} \times S_n^0$ is called the *generalized Gauss map*. The set S_n^0 decomposes into invariant subsets under the map G ; these subsets are called *Rauzy classes*. Following arguments by Veech in [73], it is proved in [80] that for every Rauzy class R the generalized Gauss map is ergodic with respect to an absolutely continuous invariant probability measure μ on $\Delta^{n-1} \times R$.

The generalized Gauss map plays a role, similar to that of the Teichmüller geodesic flow for quadratic differentials.

Lyapunov exponents and the asymptotic flag. Given a point $x \in J_q$, denote by $B_{pq}(x)$ the number of points in the trajectory segment $x, T(x), \dots, T^{k-1}(x)$ that belong to I_p ; here $T_J(x) = T^k(x) \in J$ is the first return point. For fixed (p, q) the number $B_{pq}(x)$ does not depend on x , and we suppress x from the notation. We thus get an $n \times n$ matrix B .

Zorich proves in [81] that the matrix-valued function B determines a measurable cocycle on $\Delta^{n-1} \times R$:

$$\int \log^+ \|B^{-1}\| d\mu < \infty.$$

The cocycle B^{-1} has the following spectrum of Lyapunov exponents:

$$\theta_1 > \theta_2 \geq \theta_3 \geq \dots \geq \theta_g \geq 0 = \dots = 0 \geq -\theta_g \geq \dots \geq -\theta_2 > -\theta_1,$$

where g depends on the Rauzy class R and the multiplicity of zero is $n - 2g$. Conjecturally, the above spectrum of Lyapunov exponents is simple; the simplicity of θ_2 is proved in [18].

Given $(\bar{i}, \sigma) \in \Delta^{n-1} \times R$, set

$$B^{(k)}(\bar{i}, \sigma) = B(\bar{i}, \sigma)B(G(\bar{i}, \sigma)) \dots B(G^{k-1}(\bar{i}, \sigma)).$$

The cocycle B^{-1} determines a flag of subspaces in \mathbf{R}^n , depending on a point in $\Delta^{n-1} \times R$,

$$H_1(\bar{i}, \sigma) \subset H_2(\bar{i}, \sigma) \subset \dots \subset H_g(\bar{i}, \sigma) \subset H(\bar{i}, \sigma) \subset \mathbf{R}^n.$$

The flag is defined for μ almost every (\bar{i}, σ) by the following conditions:

$$\lim_{k \rightarrow \infty} \frac{\log \|B^{(k)}(\bar{i}, \sigma)^{-1}v\|}{k} = -\theta_j \quad \text{for all } v \in H_j, v \notin H_{j-1},$$

$$\lim_{k \rightarrow \infty} \frac{\log \|B^{(k)}(\bar{i}, \sigma)^{-1}v\|}{k} > 0 \quad \text{for all } v \notin H.$$

Main result. Given an exchange T of n intervals, consider the following vector-valued function counting returns to the intervals (I_1, \dots, I_n) :

$$S(x, N) = (S_1(x, N), \dots, S_n(x, N)), \quad S_i(x, N) = \sum_{k=0}^{N-1} \chi_i(T^k(x))$$

where $x \in [0, 1)$ and χ_i is the indicator function of the interval I_i .

The main result by Zorich is the following theorem.

Theorem 6.15 *For a fixed Rauzy class R and μ almost every interval exchange transformation $(\bar{i}, \sigma) \in \Delta^{n-1} \times R$, the counting vector-function $S(x, N)$ enjoys the following properties.*

For every $x \in [0, 1)$ one has:

$$\lim_{N \rightarrow \infty} \frac{S(x, N)}{N} = \bar{i},$$

and the one-dimensional subspace H_1 is spanned by the vector \bar{i} .

For every covector $f \in \text{Ann}(H_j)$, $f \notin \text{Ann}(H_{j+1})$, $j = 1, \dots, g-1$, and every $x \in [0, 1)$ one has:

$$\limsup_{N \rightarrow \infty} \frac{\log |(f, S(x, N))|}{\log N} = \frac{\theta_{j+1}}{\theta_j}.$$

For every covector $f \in \text{Ann}(H)$, $\|f\| = 1$, and every $x \in [0, 1)$ one has:

$$|(f, S(x, N))| \leq \text{Const},$$

and the constant does not depend on either f, x or N .

A particular case $j = 1$ provides an explanation of the experimental observation (6.5). Fix $k \in \{1, \dots, n\}$ and consider the covector

$$f_k = (0 \dots 010 \dots 0) - i_k(1 \dots 1)$$

where the first vector has 1 at the k -th position. Since the space H_1 is spanned by the vector \bar{i} , the covector f_k lies in $\text{Ann } H_1$. One has:

$$((0 \dots 010 \dots 0), S(x, N)) = S_k(x, N), \quad ((1 \dots 1), S(x, N)) = N,$$

therefore Theorem 6.15 implies that

$$\limsup_{N \rightarrow \infty} \frac{\log |S_k(x, N) - i_k N|}{\log N} = \frac{\theta_2}{\theta_1}.$$

This is (6.5) with $J = I_k$ and $\alpha = \theta_2/\theta_1$.

7 Miscellaneous results

7.1 Stable periodic trajectories

The material in this subsection is taken from the paper [23].

Definition 7.1 *A periodic billiard trajectory in a polygon Q is called stable if an arbitrary small perturbation of the polygon Q leads to a perturbation of this trajectory but not to its destruction.*

For example, a 2-periodic trajectory in a square, perpendicular to a side, is not stable: a deformation of the square to a quadrilateral without parallel sides destroys this trajectory. The Fagnano 3-periodic trajectory in an acute triangle, connecting the foot points of the altitudes, is stable, and the same holds for the 6-periodic trajectory which is the double of the Fagnano orbit – see figure 1 again.

What follows is a criterion for the stability of periodic trajectories. Label the sides of the polygon Q by $1, 2, \dots, k$. Without loss of generality, assume that the periodic trajectory has an even number of links $2n$ (otherwise, double it). Then the trajectory is encoded by the sequence i_1, \dots, i_{2n} of the labels of the consecutively visited sides of Q .

Lemma 7.2 *The trajectory is stable if and only if the labels in the sequence i_1, \dots, i_{2n} can be partitioned into pairs of identical labels such that the label in each pair appears once at an odd position and once at an even one.*

Proof. Denote by α_i the angle, made by i -th side of Q with a fixed direction. Unfold the periodic trajectory to a straight line l . Then the $(2n)$ -th copy of Q along this line, denoted by Q_{2n} , is parallel to Q and, moreover, the parallel translation t that takes Q to Q_{2n} preserves l . The composition of two reflections is a rotation through the angle, twice that between the axes of reflection. It follows that Q_{2n} is obtained from Q by a rotation through the angle

$$2(\alpha_{i_1} - \alpha_{i_2} + \alpha_{i_3} - \dots + \alpha_{i_{2n-1}} + \alpha_{i_{2n}}),$$

and this angle is a multiple of 2π .

Thus for the trajectory to be stable it is necessary that the variation of this angle satisfies the equation

$$\delta\alpha_{i_1} - \delta\alpha_{i_2} + \delta\alpha_{i_3} - \dots + \delta\alpha_{i_{2n-1}} + \delta\alpha_{i_{2n}} = 0$$

for every perturbation of the polygon. An arbitrary variation of Q gives rise to an arbitrary variation of the directions of its sides. Hence the above relation holds only if its terms cancel pairwise, as claimed.

Conversely, assume that the condition of the lemma on the coding of a periodic trajectory is satisfied. Using the same notation as before, one has: $t(l) = l$. Let Q' be a polygon, sufficiently close to Q , whose sides are labeled the same way as in Q . Reflect Q' consecutively in its sides according to the sequence i_1, \dots, i_{2n} . Then Q'_{2n} is parallel to Q' . Let t' be the respective parallel translation; this translation is close to t . Choose a point $x \in l \cap Q \cap Q'$ and let l' be the line through x , invariant under t' . If Q' is sufficiently close to Q then l' is the unfolding of a billiard trajectory in Q' labeled i_1, \dots, i_{2n} , and since $t'(l') = l'$ this trajectory is $(2n)$ -periodic. \square

The lemma implies the next property of irrational billiards.

Corollary 7.3 *If the angles of a billiard k -gon are maximally independent over the rationals, that is, the dimension of the linear space over \mathbf{Q} , generated by the angles, equals $k - 1$, then every periodic billiard trajectory in this polygon is stable.*

Proof. Consider a periodic billiard trajectory, labeled i_1, \dots, i_{2n} . Then

$$\alpha_{i_1} - \alpha_{i_2} + \alpha_{i_3} - \dots + \alpha_{i_{2n-1}} + \alpha_{i_{2n}} = \pi m, \quad m \in \mathbf{Z}.$$

Rewrite this as

$$(\alpha_{i_1} - \alpha_1) - (\alpha_{i_2} - \alpha_1) + (\alpha_{i_3} - \alpha_1) - \dots + (\alpha_{i_{2n-1}} - \alpha_1) + (\alpha_{i_{2n}} - \alpha_1) = \pi m.$$

If the labels do not satisfy the condition of the above lemma, that is, do not cancel pairwise, one obtains a non-trivial relation over \mathbf{Q} on the angles $\alpha_2 - \alpha_1, \alpha_3 - \alpha_1, \dots, \alpha_k - \alpha_1$. It remains to notice that the linear space over \mathbf{Q} , generated by these angles, coincides with that, generated by the angles of the billiard polygon. \square

7.2 Encoding billiard trajectories. Polygonal billiards have zero entropy

The encoding of billiard trajectories by the consecutively visited sides of the billiard polygon provides a link between billiard and symbolic dynamics. Let us consider this encoding in more detail. We follow the paper [22].

Assume that the billiard k -gon Q is simply connected (most of the results hold without this assumption), and its sides are labeled $1, \dots, k$. The phase space of the

billiard transformation T consists of the inward unit vectors whose foot points are on the boundary ∂Q and whose forward orbits never hit a vertex of Q . Given a phase point x , denote by $w(x)$ the infinite sequence of labels $1, \dots, k$ that encodes the forward T -orbit of x . Given an infinite word w , denote by $X(w)$ the set of phase points x with $w(x) = w$. Call a subset $S \subset X(w)$ a strip if all vectors in S are parallel and their foot points constitute an interval on a side of Q . An open strip is defined analogously. The billiard transformation takes a strip to a strip, preserving its width.

Start with two simple observations. First, if $w(x) = w(y)$ then the vectors x and y are parallel. If not, the unfolded trajectories of x and y linearly diverge, so a vertex of a copy of Q will fall into the angle between them. First time this happens the reflections of the two trajectories occur at different sides of the polygon, thus $w(x) \neq w(y)$.

Fig. 12

Secondly, if x and y are parallel vectors in $X(w)$ then every parallel vector, whose foot point lies between those of x and y also belongs to $X(w)$. This is true because Q is simply connected. It follows that $X(w)$ is the maximal strip corresponding to the word w . The trajectories of its boundary points come arbitrarily close to vertices of Q .

The next theorem from [22] shows that the encoding is one-to-one on the set of non-periodic billiard trajectories.

Theorem 7.4 *If w is a $(2n)$ -periodic word then each vector from $X(w)$ has a $2n$ -periodic trajectory and $X(w)$ is an open strip. If w is aperiodic then the set $X(w)$ consists of at most one point.*

Proof. Start with the first claim. Consider $x \in X(w)$ and unfold its trajectory to the line l . One claims that the $(2n)$ -th copy of Q along this trajectory denoted, as before, by Q_{2n} , is parallel to Q . If not, x and $T^{2n}(x)$ have their foot points on the same side of Q but they are not parallel. According to the remark, preceding the theorem, $w(x) \neq w(T^{2n}(x))$ which contradicts the periodicity of w .

Let t be the parallel translation that takes Q to Q_{2n} . One wants to show that $t(l) = l$. Suppose not. Then, for a sufficiently great integer k , the polygon $t^k(Q)$ does not intersect l . This is a contradiction because this polygon is Q_{2nk} . Finally, the strip $X(w)$ is open because a periodic trajectory stays a bounded distance from the vertices of the polygon.

Now consider the second claim. Assume, to the contrary, that the set $X(w)$ is a maximal strip of non-zero width; call this strip S . Let $x \in S$ be the vector whose base point lies exactly in the middle of S .

The proof is especially simple if Q is a rational polygon, and we start with this case. Since Q is rational, the trajectory of x has a finite set of directions. It follows that this trajectory will visit some side of Q infinitely often with some fixed direction. Call these points $x_i = T^{n_i}(x)$ and let $S_i = T^{n_i}(S)$. All these strips have the same width, therefore they must intersect. Suppose S_i and S_j intersect. Then the maximal width strips \bar{S}_i and \bar{S}_j , containing S_i and S_j , respectively, are also parallel and intersect. These maximal strips cannot coincide, otherwise w would be periodic. Therefore a part of the boundary of \bar{S}_i lies inside \bar{S}_j (or vice versa). Since the trajectories of the boundary points of the maximal strip \bar{S}_i come arbitrarily close to vertices of Q , there are vertices inside \bar{S}_i , a contradiction.

Consider a general polygon Q . One cannot claim anymore that the trajectory of x visits some side of Q infinitely often with some fixed direction. This claim will be replaced by a weaker recurrence result.

Let Y be the forward limit set of x under the mapping T . Since the width of S is positive, the trajectories of the points from Y stay bounded distance from the vertices of Q . Therefore the restriction of T to Y is continuous. Note that Y is closed and bounded, hence compact.

One uses the strengthened version of the Poincaré recurrence theorem, due to H. Furstenberg [20]. This theorem implies that there exists a uniformly recurrent point $y \in Y$. This means the following: for every neighborhood U of y there is a constant C such that the return times $n_i > 0$, defined by $T^{n_i}(y) \in U$, satisfy the inequalities $n_{i+1} - n_i < C$.

Fix a sufficiently small $\epsilon > 0$ and consider the unfolded maximal strip S' of y together with its ϵ neighborhood S'_ϵ . Since S' is maximal, some vertices of copies of Q , unfolded along this strip, will fall into S'_ϵ . *Claim:* they fall with uniformly bounded gaps into each of the two components of $S'_\epsilon - S'$.

This claim follows from the uniform recurrence of y . More specifically, let z be the leftmost point of S' . Then there is $m > 0$ such that the foot point of $T^m(z)$ lies within $\epsilon/2$ of a vertex v of Q ; let m be the first such number. Let U be a neighborhood of y such that for every $u \in U$ the m -th iterate $T^m(u)$ is $\epsilon/2$ close to $T^m(y)$; such a neighborhood exists since T is continuous at y . Let n_i be the return times of y to U . Then $T^{m+n_i}(y)$ is $\epsilon/2$ close to $T^m(y)$, and therefore $T^{m+n_i}(z)$ is $\epsilon/2$ close to $T^m(z)$. It follows that the foot points of $T^{m+n_i}(z)$ lie within ϵ from the vertex v . The claim is proved.

To finish the proof of the theorem recall that y is a forward limit of x : there is a sequence $n_i \rightarrow \infty$ such that $x_i = T^{n_i}(x) \rightarrow y$. The argument, given above in the case of a rational polygon Q , shows that no two vectors x_i and x_j are parallel. Consider the intersection of the unfolded maximal strip for x_i with $S'_\epsilon - S'$. Let α_i be the angle between these strips of width ϵ and, say, δ . Then $\alpha_i \rightarrow 0$ as $i \rightarrow \infty$. The intersection is a parallelogram which, by elementary geometry, contains a rectangle of width ϵ and length $l_i = (\delta - \epsilon \cos \alpha_i) / \sin \alpha_i$. Since the gaps between the vertices in $S'_\epsilon - S'$ are uniformly bounded and $l_i \rightarrow \infty$ as $i \rightarrow \infty$, a vertex will eventually appear inside the unfolded maximal strip for x_i . This is a contradiction. \square

As a consequence, the closure of every non-periodic billiard orbit contains a vertex of the billiard polygon.

An important corollary of the theorem concerns the entropy of polygonal billiards.

Corollary 7.5 *The metric entropy of the billiard mapping with respect to any invariant measure is zero.*

Proof. Morally, the argument goes as follows: the "past" is uniquely determined by the "future", therefore the entropy vanishes.

More specifically, consider the set Σ of words in $1, \dots, k$ that encode the billiard trajectories for time from $-\infty$ to ∞ , and let S be the shift transformation of Σ . Then the encoding map conjugates the billiard transformation and S . By the theorem, the shift S has a one-sided generator, the partition into k parts according to the value of the zero symbol in a word. This means that for this partition η the σ -algebra of measurable sets is generated by $\bigvee_{i=0}^{\infty} S^{-i}(\eta)$. It is a standard result in ergodic theory that the metric entropy of a transformation with a one-sided generator vanishes (see, e.g., [13]). \square

The variational principle implies that the topological entropy vanishes as well.

The zero entropy result was proved in [8] and [63] for the canonical billiard measure, and in [37] for topological entropy (see also [22] and [28, 29]). The latter works concerns a broader class of transformations, the so-called, polygon exchanges.

The zero entropy result implies that a number of quantities, associated with a polygonal billiard, grow slower than exponentially: the number of different words of length n in Σ , the number of strips of n -periodic trajectories or periodic trajectories of length not greater than L , etc. Conjecturally, all these quantities have polynomial growth.

7.3 Complexity of billiard trajectories in rational polygons

Consider an aperiodic billiard trajectory in a rational polygon in a direction θ . According to Theorem 7.4 its code w is an aperiodic sequence. Generically this sequence enjoys a weaker periodicity property, called quasiperiodicity, described in the next result.

Proposition 7.6 *For all but countably many directions θ each finite segment of w appears in w infinitely many times.*

Proof. Let i_1, \dots, i_n be a segment of w , and let x be a phase point on the trajectory under consideration whose first n successive reflections occur in the sides labeled i_1, \dots, i_n . There is a neighborhood U of x such that for every phase point $y \in U$ the first n reflections occur in the same sides of the polygon. By Theorem 1.8, for all but countably many directions θ , the trajectory is dense on the invariant surface. It follows that x returns to U infinitely many times, and each time the segment i_1, \dots, i_n reappears in w . \square

A convenient way to measure the complexity of aperiodic trajectories and their codes is provided by the next definition.

Definition 7.7 *The complexity function $p(n)$ of an infinite sequence w is the number of distinct n -element segments of w .*

We start with the complexity of billiard trajectories in a square; the results are due to Hedlund and Morse [33]. We slightly modify the encoding using only two symbols, say, 0 and 1, to indicate that a trajectory reflects in a horizontal or a vertical side, respectively.

Theorem 7.8 *For every aperiodic trajectory one has: $p(n) = n + 1$.*

The sequences with complexity $p(n) = n + 1$ are called Sturmian.

Proof. Start with the following observation. Every aperiodic trajectory is dense in the respective invariant torus. Therefore $p(n)$ can be computed as the number of different initial segments of length n in the codes $w(x)$ where x ranges over the phase vectors having a fixed direction.

Unfold the billiard trajectory to a line l . Partition the square grid in the plane into "ladders", going in the south-east – north-west direction, as shown in figure 13. The

n -th symbol in the code of the trajectory is 0 or 1, according as l meets a horizontal or a vertical segment of n -th ladder.

Fig. 13

Let (e_1, e_2) be the orthonormal frame. Consider the linear projection of the plane onto the diagonal $x + y = 0$ whose kernel is parallel to the line l . Factorize the diagonal by the translation through $e_1 - e_2$ and identify the quotient space with the unit circle \mathbf{S}^1 . The vertices of the first ladder are the lattice points (a, b) with $a + b = 1$ or $a + b = 2$. Since l has an irrational slope the projections of the vertices of the first ladder partition the circle into two irrational arcs. Let T be the rotation of \mathbf{S}^1 through the length of an arc, that is, through the projection of e_1 .

The number of different initial n -segments corresponding to the lines, parallel to l , equals the number of segments into which the projections of the vertices of the first n ladders partition \mathbf{S}^1 . Each ladder is obtained from the first one by the translation through e_1 . It follows that $p(n)$ equals the cardinality of the orbit $T^i(0)$, $i = 0, \dots, n$. Since T is an irrational rotation of the unit circle all points of this orbit are distinct, and $p(n) = n + 1$. \square

There are generalizations of the above theorem to multi-dimensional cubes – see [2] and [4].

Next we consider the complexity of billiard trajectories in convex rational polygons. The following result is due to P. Hubert ([34]).

Let Q be a convex rational k -gon with the angles $\pi m_i/n_i$, $i = 1, \dots, k$ (m_i and n_i are coprime), and let N be the least common multiple of n_i 's. Let θ be a direction such that there are no generalized diagonals in this direction. Consider a billiard trajectory in direction θ which avoids the vertices. Let w be its coding in the alphabet $\{1, \dots, k\}$, and denote by $p(n)$ the complexity of w .

Theorem 7.9 *For all sufficiently great n one has: $p(n) = n(k - 2)N + 2N$.*

In particular, if Q is a square, one obtains: $p(n) = 4(n + 1)$. This does not contradict the previous theorem because the encoding there was different and, actually, 4 times less precise: the pair of parallel sides were labeled identically.

Proof. Consider the invariant surface M of the billiard flow. Recall that M is the result of pairwise identifying the sides of $2N$ copies of Q . Thus M has the structure of a complex; denote by v, e and f the number of vertices, edges and faces. Note that $e = Nk$ and $f = 2N$. The edges can be labeled by pairs (i, j) , $i = 1, \dots, k$, $j = 1, \dots, N$

so that the (i, j) -th edge corresponds to the i -th side of Q . Let $E \subset M$ be the union of the edges.

The billiard trajectory under consideration can be encoded in an obvious way in the alphabet $\{(i, j)\}$; denote its code by W . The projection $(i, j) \rightarrow i$ sends W to w (we will always denote by W words in the alphabet $\{(i, j)\}$ and by w their projections). In other words, the new encoding takes into account not only the side of Q in which a reflection occurs but also the angle of reflection that can take finitely many values.

Let $P(n)$ be the complexity of W .

Claim: $P(n) = n(k - 2)N + 2N$ for all $n \geq 1$.

Arguing as in the proof of the preceding theorem, $P(n)$ equals the number of distinct initial segments of length n in the codes $W(x)$ where x ranges over points of E .

Next, given a word W_n of length n , consider the set $X(W_n)$ consisting of the points $x \in E$ such that the initial segment of $W(x)$ is W_n . Then, arguing as in the previous section, $X(W_n)$ is an interval on an edge from E .

Denote by T the first return map of the directional flow F_θ to E (that is, the billiard map). Consider the set of initial segments of length $(n + 1)$ in $W(x)$, $x \in X(W_n)$. The word W_n will have different successors if and only if $T^n(X(W_n))$ contains a vertex V of M . We say that W_n splits at V .

To find $P(n + 1) - P(n)$ one needs to learn how many words of length n split at each vertex. Let V be a vertex with cone angle $2\pi c$. Then there are c incoming saddle connections. Trace each one back n steps, and let W^1, \dots, W^c be the respective n -length words. Then $T^n(X(W^\alpha))$ contains V for each $\alpha = 1, \dots, c$.

The words W^α are all different since their $(n - 1)$ -st letters are distinct. The latter hold because in the polygon Q there is only one trajectory that starts at a given side in a given direction and goes straight to a given vertex.

Thus all the n -length words that split at V are W^1, \dots, W^c . It follows that

$$P(n + 1) - P(n) = \sum c(V),$$

sum over all vertices. To evaluate this sum note that the Euler characteristic of M equals the sum of indices of singular points of the flow F_θ ; the index at vertex V equals $1 - c(V)$. Hence

$$\chi(M) = \sum (1 - c(V)) = v - \sum c(V).$$

On the other hand, $\chi(M) = v - e + f$, and it follows that $\sum c(V) = (k - 2)N$. Therefore $P(n + 1) - P(n) = (k - 2)N$, and since $P(1) = kN$, the claim follows.

To finish the proof of the theorem it remains to show that $P(n) = p(n)$ for all sufficiently great n . Suppose not; then there exist arbitrary long words $W \neq W'$ that are the initial segments of the codes of phase points $x \in E$ such that $w = w'$.

The billiard trajectories, corresponding to W and W' , start at the same side of Q but make different angles with it: otherwise, since $w = w'$, they would reflect at the same sides of Q and meet them at the same angles which would imply that $W = W'$. Note that the angle between the two trajectories is bounded below by a constant α depending on Q and θ .

Unfold the trajectories to straight lines. Since these lines linearly diverge, there is a constant m , depending on Q and α , such that after at most m reflections of Q along either of the unfolded trajectories a vertex of a copy of Q will fall into the angle between the lines. Therefore if the length of W and W' is greater than m the codes w and w' are distinct, a contradiction. \square

If Q is a general k -gon let N be the least common denominator of its π -rational angles and s be the number of its distinct π -irrational angles. Then one has the next upper bound for the complexity $p_\theta(n)$ of the billiard trajectories in Q in a fixed initial direction θ – see [32].

Theorem 7.10 *For all n one has:*

$$p_\theta(n) \leq kNn\left(1 + \frac{n}{2}\right)^s.$$

Using techniques, similar to the ones in the proof of Theorem 7.9, S. Troubetzkoy obtained in [67] a complexity lower bound for arbitrary polygonal billiards. Unlike the previous results, his theorem concerns billiard trajectories in all directions, so that $p(n)$ now means the number of distinct words of length n in the coding in the alphabet $\{1, \dots, k\}$ of all billiard orbits in a k -gon.

Theorem 7.11 *For every polygon there exists a constant c such that $p(n) \geq cn^2$ for all $n \geq 0$.*

A similar estimate is proved in [67] for d -dimensional polyhedra with the exponent 2 replaced by d .

7.4 Periodic trajectories in some irrational billiards

The following elementary argument shows that every rational polygon has a periodic billiard trajectory of a special kind; it was found independently by A. Stepin ([23]) and by M. Boshernitzan ([10]).

Shoot the billiard ball in the direction, perpendicular to a side of the polygon Q . By Poincaré's recurrence theorem, for almost every initial position, the ball will return to the original side at an angle, arbitrarily close to $\pi/2$. Since the set of possible directions of the ball is discrete in a rational polygon, this angle will be equal to $\pi/2$. After the ball bounces off of the side it will backtrack the same trajectory, so the trajectory is periodic.

Surprisingly, a variation of this argument applies to some irrational polygons. This was observed by B. Cipra, R. Hanson and A. Kolan ([14]), and then generalized by E. Gutkin and S. Troubetzkoy ([32]). The next result is due to the former authors.

Theorem 7.12 *Given a right triangle with π -irrational acute angles, almost every (in the sense of measure) billiard trajectory, that starts at a side of the right angle in the direction, perpendicular to this side, returns to the same side in the same direction.*

Proof. The idea is to apply the construction which gave the invariant surfaces of the billiard flow in rational polygons. A significant difference with the rational case is that the invariant surface will not be compact.

Start with reflecting the triangle in the sides of the right angle to obtain a rhombus R . The study of the billiard inside the triangle reduces to that inside the rhombus. Let α be the acute angle of R .

Fig. 14

Consider the beam of horizontal trajectories, starting at the upper half of the vertical diagonal of the rhombus, and construct the invariant surface of the phase space, corresponding to this beam. This surface consists of rhombi, obtained from R by the action of the group $A(R)$, whose sides are pasted pairwise in an appropriate manner. Since α is π -irrational, this surface is not compact. The invariant surface is partially foliated by the parallel trajectories from the horizontal beam, and this foliation has an invariant transversal measure, "the width of a beam".

Each rhombus involved is obtained from the original one by a rotation through the angle $n\alpha$ where $n \in \mathbf{Z}$. Such a rhombus will be referred to as the n -th rhombus and denoted by R_n ; in particular, $R = R_0$. A trajectory from the beam under consideration may leave the n -th rhombus through a side which has one of the two possible directions. Call such a side positive if the trajectory enters an $(n + 1)$ -st rhombus and negative if it enters an $(n - 1)$ -st one.

One wants to show that almost every trajectory returns to R_0 (where they hit the vertical diagonal in the perpendicular direction). We will prove that for every $\epsilon > 0$ the relative measure of the set of trajectories that do not return to R_0 is less than ϵ .

Let ϵ be given. Since α is π -irrational, there exists $n > 0$ such that the vertical projection of the positive side of R_n is less than ϵ . This implies that the relative measure of the set of trajectories that make it to $(n + 1)$ -st rhombi is less than ϵ .

The rest of trajectories are bound to stay in R_0, R_1, \dots, R_n ; call this set of trajectories S . The union of the rhombi from 0 through n is compact, therefore the Poincaré recurrence theorem applies to S . It follows that almost every trajectory in S is recurrent, that is, returns to R_0 . \square

In fact, the above proof gives more: for every direction in a rhombus almost every billiard trajectory in this direction returns arbitrarily close to itself and in the same direction. A relevant definition is proposed in [32].

Definition 7.13 *Given a billiard polygon Q , a direction is called recurrent if Lebesgue almost every phase point with this direction returns to the same direction. A polygon Q is called strongly recurrent if every direction is recurrent.*

In particular, every rational polygon is strongly recurrent. An argument, similar to the preceding proof, gives the next result from [32], applicable, in particular, to all parallelograms.

Theorem 7.14 *Let Q be a polygon whose sides have one of the two fixed directions. Then Q is strongly recurrent. In particular, the billiard orbits, perpendicular to a side of Q , are periodic with probability one.*

7.5 A non-periodic trajectory that is not dense in the configuration space

An example of such a trajectory was constructed by G. Galperin in [21].

Fig. 15

Consider a centrally symmetric hexagon (so that its opposite sides are parallel and congruent) such that the bisector of the angle $BAC = \alpha$ is perpendicular to the side AF and the the bisector of the angle $ABF = \beta$ is perpendicular to the side BC .

Since the bisector of α is perpendicular to BC , the reflection of AB in BC is the line BF . Thus the reflection of any vertical segment leading to BC is parallel to BF .

Similarly, the reflection of vertical segments in CD are parallel to AC . After another reflection the segment becomes vertical again. Choose a line, perpendicular to the side AB ; let X, Y, Z be the projections on this line of the points A, E, F , respectively. The set of vertical lines is identified with a horizontal segment XY , and the second iteration of the billiard transformation induces the exchange of the intervals XZ and ZY . An exchange of two intervals is equivalent to a circle rotation. For generic angles α and β each orbit of this circle rotation is dense. Therefore each vertical billiard trajectory is dense in the hexagon.

Extend the non-vertical sides of the hexagon to obtain a parallelogram. Then the two triangles, added to the hexagon, are never visited by any vertical trajectory, that starts inside the original hexagon.

Similar but more technical considerations prove the following result by Galperin.

Theorem 7.15 *There exists an acute angle α_0 such that for almost every (in the sense of measure) $\alpha \in (0, \alpha_0)$ the right triangle with the acute angle α contains a non-periodic and not everywhere dense billiard trajectory.*

References

- [1] P. Arnoux. Ergodicité Générique des Billards Polygonaux. Sémin. Bourbaki, No 696, (1987-88).
- [2] P. Arnoux, C. Mauduit, I. Shiokawa, J.-I. Tamura. Complexity of sequences defined by billiards in the cube. Bull. SMF, 122, (1994), 1-12.
- [3] E. Aurell, C. Itzykson. Rational billiards and algebraic curves. J. Geom. and Phys., 5, (1988), 191-208.
- [4] Yu. Baryshnikov. Complexity of trajectories in rectangular billiards. Comm. Math. Phys., 174, (1995), 43-56.
- [5] A. Beardon. The geometry of discrete groups. Springer-Verlag, 1983.
- [6] L. Bers. Finite dimensional Teichmüller spaces and generalizations, Bull. AMS, 5, (1981), 131-172.
- [7] S. Bleiler, A. Casson. Automorphisms of surfaces after Nielsen and Thurston, Cambridge Univ. Press, 1987.

- [8] C. Boldrighini, M. Keane, F. Marchetti. Billiards in polygons. *Ann. of Prob.*, 6, (1978), 532-540.
- [9] M. Boshernitzan, A condition for a minimal interval exchange transformation to be uniquely ergodic. *Duke J. Math.* 52, (1985), 723-752.
- [10] M. Boshernitzan. Billiards and rational periodic directions in polygons. *Amer. Math. Monthly*, 99, (1992), 522-529.
- [11] M. Boshernitzan, G. Galperin, T. Kruger, S. Troubetzkoy. Periodic billiard orbits are dense in rational polygons. *Trans. AMS*, 350, (1998), 3523-3535.
- [12] Y. Cheung. Ph.D. dissertation. University of Illinois at Chicago.
- [13] I. Cornfeld, S. Fomin, Ya. Sinai. *Ergodic theory*. Springer-Verlag, 1982.
- [14] B. Cipra, R. Hanson, A. Kolan. Periodic trajectories in right triangle billiards. *Phys. Rev.*, E 52, (1995), 2066-2071.
- [15] C. Earle, F. Gardiner. *Teichmüller disks and Veech's \mathcal{F} -structures*. *Contemp. Math.*, 201, AMS, Providence, 1997, 165-189.
- [16] A. Eskin, H. Masur. Pointwise asymptotic formulas on flat surfaces, *Ergod. Th. Dyn. Syst.*, 21, (2001), 443-478.
- [17] A. Fathi, F. Laudenbach, V. Poenaru. *Travaux de Thurston sur les surfaces*. Asterisque, 66-67, 1979.
- [18] G. Forni. Deviation of ergodic averages for area preserving flows on surfaces of higher genus. Preprint.
- [19] R. Fox, R. Kershner. Geodesics on a rational polyhedron. *Duke Math. J.*, 2, (1936), 147-150.
- [20] H. Furstenberg. *Recurrence in Ergodic Theory and Combinatorial Number Theory*. Princeton Univ. Press, 1981.
- [21] G. Galperin. Nonperiodic and not everywhere dense billiard trajectories in convex polygons and polyhedrons. *Comm. Math. Phys.*, 91, (1983), 187-211.
- [22] G. Galperin, T. Kruger, S. Troubetzkoy. Local instability of orbits in polygonal and polyhedral billiards. *Comm. Math. Phys.*, 169, (1995), 463-473.

- [23] G. Galperin, A. Stepin, Ya. Vorobets. Periodic billiard trajectories in polygons: generating mechanisms. *Russ. Math. Surv.*, 47, No 3, (1992), 5-80.
- [24] S. Glashow, L. Mittag, Three rods on a ring and the triangular billiard. *J. Stat. Phys.* 87 (1997), 937–941.
- [25] E. Gutkin. Billiards in polygons. *Physica D*, 19, (1986), 311-333.
- [26] E. Gutkin. Billiards in polygons: survey of recent results. *J. Stat. Phys.*, 83, (1996), 7-26.
- [27] E. Gutkin. Billiards on almost integrable polyhedral surfaces. *Ergod. Th. Dyn. Syst.*, 4, (1984), 569-584.
- [28] E. Gutkin, N. Haydn. Topological entropy of generalized polygon exchanges. *Bull. AMS*, 32, (1995), 50-57.
- [29] E. Gutkin, N. Haydn. Topological entropy of polygon exchange transformations and polygonal billiards. *Erg. Th. Dyn. Syst.*, 17, (1997), 849-867.
- [30] E. Gutkin, C. Judge. The geometry and arithmetic of translation surfaces with applications to polygonal billiards. *Math. Res. Lett.*, 3, (1996), 391-403.
- [31] E. Gutkin, C. Judge. Affine mappings of translation surfaces: geometry and arithmetic. *Duke Math. J.* 103, (2000) 191-213
- [32] E. Gutkin, S. Troubetzkoy. Directional flows and strong recurrence for polygonal billiards. in *Pitman Res. Notes Math. Ser.*, 362, 1996, 21-45.
- [33] G. Hedlund, M. Morse. Symbolic dynamics 2. Sturmian trajectories. *Amer. J. of Math.*, 62, (1940), 1-42.
- [34] P. Hubert. Complexité de suites définies par des billards rationnels. *Bull. Soc. Math. France*, 123, (1995), 257-270.
- [35] P. Hubert, T. Schmidt, Veech groups and polygonal coverings. *J. Geom. Phys.* 35 (2000), no. 1, 75–91
- [36] A. Katok, Invariant measures for flows on oriented surfaces. *Sov. Math. Dokl.* 14, (1973), 1104-1108

- [37] A. Katok. Interval exchange transformations and some special flows are not mixing. *Isr. Math. J.*, 35, (1980), 301-310.
- [38] A. Katok. The growth rate for the number of singular and periodic orbits for a polygonal billiard. *Comm. Math. Phys.*, 111, (1987), 151-160.
- [39] A. Katok, B. Hasselblatt. *Introduction to the modern theory of dynamical systems*. Cambridge Univ. Press, 1995.
- [40] A. Katok, B. Hasselblatt. Principal structures. *This Handbook*.
- [41] A. Katok, A. Zemlyakov. Topological transitivity of billiards in polygons. *Math. Notes*, 18, (1975), 760-764.
- [42] M. Keane. Coding problems in ergodic theory. *Proc. Int. School of Math. Phys.*, Univ. Camerino, 1974.
- [43] M. Keane. Interval exchange transformation. *Math. Zeitschrift*, 141, (1975), 25-31.
- [44] M. Keane, Nonergodic interval exchange transformations. *Israel J. Mathematics*, 26, (1977), 188-196.
- [45] R. Kenyon, J. Smillie. Billiards in rational-angled triangles. *Comm. Math. Helv.*, 75, (2000), 65-108.
- [46] S. Kerckhoff. Simplicial systems for interval exchange transformations and measured foliations, *Ergod. Th. Dyn. Syst.* 5, (1985), 257-271.
- [47] S. Kerckhoff, H. Masur, J. Smillie. Ergodicity of Billiard Flows and Quadratic Differentials. *Ann. of Math.*, 124, (1986), 293-311.
- [48] M. Kontsevich. Lyapunov exponents and Hodge theory. *Adv. Ser. Math. Phys.*, 24, 318–332. World Sci. Publ., 1997.
- [49] M. Kontsevich, A. Zorich. Connected components of the moduli spaces of holomorphic differentials with prescribed singularities. Preprint.
- [50] I. Kra. On the Nielsen-Thurston-Bers type of some self-maps of Riemann surfaces. *Acta. Math.*, 146, (1981), 231-270.

- [51] H.Masur. Interval exchange transformations and measured foliations. *Ann. of Math.*, 115, (1982), 169-200.
- [52] H. Masur, Closed trajectories for quadratic differentials with an application to billiards. *Duke Math. J.*, 53, (1986), 307-314.
- [53] H. Masur. Lower bounds for the number of saddle connections and closed trajectories of a quadratic differential. In *Holomorphic Functions and Moduli*, Vol.1, D. Drasin ed., Springer-Verlag, 1988, 215-228.
- [54] H.Masur. The growth rate of trajectories of a quadratic differential. *Erg. Th. Dyn. Syst.*, 10, (1990), 151-176.
- [55] H. Masur. Hausdorff dimension of the set of nonergodic foliations of a quadratic differential. *Duke Math. J.*, 66, (1992), 387-442.
- [56] H. Masur, J.Smillie. Hausdorff dimension of sets of nonergodic foliations, *Ann. of Math.*, 134, (1991), 455-543.
- [57] H. Masur, J. Smillie. Quadratic differential with prescribed singularities and pseudo-Anosov diffeomorphisms. *Comm. Math. Helv.*, 68, (1993), 289-307.
- [58] J. Oxtoby, S. Ulam. Measure preserving homeomorphisms and metrical transitivity, *Ann. of Math.*, 42, (1941), 874-920.
- [59] G. Rauzy. Echanges d'intervalles et transformations induites. *Acta Arith.*, 34, (1979), 315-328.
- [60] M. Rees. An alternative approach to the ergodic theory of measured foliations on surfaces, *Ergod. Th. Dyn. Syst.*, 1, (1981), 461-488.
- [61] R. Richens, M. Berry. Pseudointegrable systems in classical and quantum mechanics. *Physica D*, 2, (1981), 495-512.
- [62] E. Sataev. On the number of invariant measures of flows on orientable surfaces. *Izv. Acad. Sci. USSR*, 9, (1975), 860-878.
- [63] Ya. Sinai. *Introduction to ergodic theory*. Princeton Univ. Press, 1976.
- [64] K. Strebel. *Quadratic differentials*. Springer-Verlag, 1984.

- [65] S. Tabachnikov. Billiards. SMF Panoramas et Synthèses, 1995, No 1.
- [66] W. Thurston. On geometry and dynamics of diffeomorphisms of surface. Bull Amer. Math. Soc., 18, (1988), 417-431.
- [67] S. Troubetzkoy. Complexity lower bounds for polygonal billiards. Chaos, 8, (1998), 242-244.
- [68] W. Veech. Strict ergodicity in zero dimensional dynamical systems and the Kronecker-Weyl theorem mod 2. Trans. AMS, 140, (1969), 1-34.
- [69] W. Veech. Teichmuller geodesic flow. Ann. of Math., 124, (1986), 441-530.
- [70] W. Veech. Moduli spaces of quadratic differentials. J. D'Analyse Math., 55, (1990), 117-171.
- [71] W. Veech. Teichmuller curves in moduli space. Eisenstein series and an application to triangular billiards. Invent. Math., 97, (1990), 117-171.
- [72] W. Veech. The billiard in a regular polygon, Geom.Funct.Anal., 2, (1992), 341-379.
- [73] W. Veech. Gauss measures for transformations on the space of interval exchange maps. Ann. Math., 115, (1982), 201-242.
- [74] W. Veech. Interval exchange transformations. J. d'Analyse Math., 33, (1978), 222-272.
- [75] Ya. Vorobets. Ergodicity of billiards in polygons. Mat. Sb., 188, (1997), No 3, 65-112.
- [76] Ya. Vorobets. Plane structures and billiards in rational polygons: the Veech alternative. Russ. Math. Surv., 51, (1996), No 5, 3-42.
- [77] P. Walters. An Introduction to Ergodic Theory. Springer-Verlag, 1982.
- [78] C. Ward. Calculation of Fuchsian groups associated to billiards in a rational triangle. Ergod. Th. and Dyn. Syst., 18, (1998), 1019-1042.
- [79] A. Zorich. Asymptotic flag of an orientable measured foliation on a surface. in Geometric Study of Foliations, World Sci., 1994, 479-498.

- [80] A. Zorich. Finite Gauss measure on the space of interval exchange transformations. Lyapunov exponents. *Ann. Inst. Fourier*, 46, (1996), 325-370.
- [81] A. Zorich. Deviation for interval exchange transformations. *Ergod. Th. and Dyn. Syst.*, 17, (1997), 1477-1499.
- [82] A. Zorich. How do the leaves of a closed 1-form wind around a surface. *AMS Transl., ser. 2*, v. 197, 135-178, AMS, Providence, 1999.