

THE HOROCYCLE FLOW AT PRIME TIMES

PETER SARNAK AND ADRIÁN UBIS

ABSTRACT. We prove that the orbit of a non-periodic point at prime values of the horocycle flow in the modular surface is dense in a set of positive measure. For some special orbits we also prove that they are dense in the whole space—assuming the Ramanujan/Selberg conjectures for GL_2/\mathbb{Q} . In the process, we derive an effective version of Dani’s Theorem for the (discrete) horocycle flow.

1. INTRODUCTION

If (X, T) is a dynamical system, for any $x \in X$ one can ask about the distribution of points $P_x = \{T^p x : p \text{ prime}\}$ in the orbit $\theta_x = \{T^n x : n \geq 1\}$. For example if X is finite then this is equivalent to Dirichlet’s Theorem on primes in an arithmetic progression. If (X, T) is ergodic, Bourgain [5] shows that for almost all x , $T^p x$ with p prime, satisfies the Birkhoff ergodic Theorem and hence is equidistributed. If (X, T) is ‘chaotic’, for example if it has a positive entropy then there may be many x ’s for which $T^p x$ is poorly distributed in $\overline{\theta_x}$. For example if $T : [0, 1] \rightarrow [0, 1]$ is the doubling map $x \mapsto 2x$ then one can construct an explicit (in terms of its binary expansion) ξ such that $\overline{\theta_\xi} = [0, 1]$ but $T^p \xi \rightarrow 0$ as $p \rightarrow \infty$.

The setting in which one can hope for a regular behaviour on restricting to primes is that of unipotent orbits in a homogeneous space. Let G be a connected Lie group, Γ a lattice in G and $u \in G$ an Ad_G unipotent element, then Ratner’s Theorem [31] says that if $X = \Gamma \backslash G$ and $T : X \rightarrow X$ is given by

$$(1.1) \quad T(\Gamma g) = \Gamma gu,$$

then $\overline{\theta_x}$, with $x = \Gamma g$, is homogeneous and the orbit xu^n , $n = 1, 2, \dots$ is equidistributed in $\overline{\theta_x}$ w.r.t. an algebraic measure $d\mu_x$. In the case that μ_x is the normalized volume measure $d\mu_G$ on X it is conjectured in [12] that $\overline{P_x} = X$ and in fact that xu^p , $p = 2, 3, 5, 7, 11, \dots$ is equidistributed w.r.t. $d\mu_G$. Care should be taken in formulating this conjecture in the intermediate cases where $\overline{\theta_x}$ is not connected as there may be local

congruence obstructions, but with the ‘obvious’ modifications this conjecture seems quite plausible. In intermediate cases where $\overline{\theta}_x$ is one of

- (i) finite
- (ii) a connected circle or more generally a torus
- (iii) a connected nilmanifold $\Gamma \backslash N$,

\overline{P}_x and the behaviour of xu^p , $p = 2, 3, 5, \dots$ is understood. Case (i) requires no further comment while for (ii) it follows from Vinogradov’s work that the points are equidistributed w.r.t dt , the volume measure on the torus. The same is true for (iii) as was shown recently by Green and Tao [13, 14]; in order to prove this, apart from using Vinogradov’s methods they had to control sums of the type $\sum_n e(\alpha n[\beta n])$, which are similar to Weyl sums but behave in a more complex way.

Our purpose in this paper is to examine this problem in the basic case of $X = SL(2, \mathbb{Z}) \backslash SL(2, \mathbb{R})$. According to Hedlund [16], $\overline{\theta}_x$ is either finite, a closed horocycle of length l , $0 < l < \infty$, or is all X . The first two cases correspond to (i) and (ii). In the last case we say that x is *generic*. By a theorem of Dani [8] the orbit xu^n , $n = 1, 2, 3, \dots$ is equidistributed in its closure w.r.t. one of the corresponding three types of algebraic measures. For $N \geq 1$ and $x \in X$ define the probability measure $\pi_{x,N}$ on X by

$$(1.2) \quad \pi_{x,N} := \frac{1}{\pi(N)} \sum_{p < N} \delta_{xu^p}$$

where for $\xi \in X$, δ_ξ is the delta mass at ξ and $\pi(N)$ is the number of primes less than N . We are interested in the weak limits ν_x of the $\pi_{x,N}$ as $N \rightarrow \infty$ (in the sense of integrating against continuous functions on the one-point compactification of X). If x is generic then the conjecture is equivalent to saying that any such ν_x is $d\mu_G$. One can also allow x , the initial point of the orbit, to vary with N in this analysis and in the measures in (1.2). Of special interest is the case $x = \Gamma g$

$$g = H_N := \begin{bmatrix} N^{-\frac{1}{2}} & 0 \\ 0 & N^{\frac{1}{2}} \end{bmatrix} \quad \text{and} \quad u = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix},$$

when $H_N u^j$, $0 \leq j \leq N - 1$ is a periodic orbit for T of period N . These points are spread evenly on the unique closed horocycle in X whose length is N . They also comprise a large piece of the Hecke points in X corresponding to the Hecke correspondence of degree N

$$C_N = \left\{ \frac{1}{\sqrt{N}} \begin{bmatrix} a & b \\ c & d \end{bmatrix} : ad = N, a, d > 0, b \pmod{d} \right\}.$$

We can now state our main results. The first asserts that ν_x does not charge small sets with too much mass, that is ν_x is uniformly absolutely continuous with respect to $d\mu_G$.

Theorem 1.1 (Non-concentration at primes). *Let x be generic and ν_x a weak limit of $\pi_{x,N}$, then*

$$d\nu_x \leq 10 d\mu_G.$$

Remark: As we said before, probably what truly happens is that $d\nu_x = d\mu_G$. If on the other hand we allow n to vary not over primes but over almost primes, then the quantitative equidistribution that we develop to prove Theorem 1.1 can be used together with a lower bound sieve (see [11], Chapter 12) to prove the density of the orbit. More precisely let x be generic, then the points $T^n x$, as n varies over numbers with at most 10 prime factors, are dense in X .

As a consequence of Theorem 1.1 we deduce that $\overline{P_x}$ has to be big.

Corollary 1.2 (Large closure for primes). *Let x be generic, then*

$$\text{Vol}(\overline{P_x}) \geq \frac{1}{10},$$

and if $U \subset X$ is an open set with $\text{Vol}(U) > 1 - 1/10$ then $xu^p \in U$, for a positive density of primes p .

In the case of ‘Hecke orbits’ P_{H_N} , we can prove more;

Theorem 1.3 (Prime Hecke orbits are dense). *Let ν be a weak limit of the measures $\pi_{H_N,N}$. Assuming the Ramanujan/Selberg Conjectures concerning the automorphic spectrum of GL_2/\mathbb{Q} (see for example the Appendix in [35]) we have*

$$\frac{1}{5} d\mu_G \leq d\nu \leq \frac{9}{5} d\mu_G.$$

Theorem 1.3 has an application to a variant of Linnik’s problem on projections of integral points on the level 1 surface for quadratic forms in 4-variables. Let $N \geq 1$ and denote by M_N the set of 2×2 matrices whose determinant equals N . Denote by π the projection $A \mapsto \frac{1}{\sqrt{N}}A$ of $M_N(\mathbb{R})$ onto $M_1(\mathbb{R})$. Using their ergodic methods Linnik and Skubenko [27] show that the projection of the integer points $M_N(\mathbb{Z})$ into $M_1(\mathbb{R})$ become dense as $N \rightarrow \infty$. In quantitative form they show that for U a (nice) compact subset of $M_1(\mathbb{R})$

$$(1.3) \quad |\{A \in M_N(\mathbb{Z}) : \pi(A) \in U\}| \sim \sigma_1(N)\mu(U)$$

as $N \rightarrow \infty$, where $\mu(u)$ is the ‘Hardy-Littlewood’ normalized Haar measure for $SL_2(\mathbb{R})$ on $M_1(\mathbb{R})$ and $\sigma_1(N) = \sum_{d|N} d$. (1.3) can also be

proved using Kloosterman's techniques in the circle method [25] as well as using Hecke Correspondences in GL_2 as explained in [33]. Using the last connection we establish the following Corollary whose formulation is cleanest when N is prime and which we assume is the case in the corollary. For

$$A = \begin{bmatrix} a & b \\ c & d \end{bmatrix} \in M_N(\mathbb{Z})$$

we let $b_1(A)$ be the slope of the kernel of A (which is a line) modulo N . That is if $N \mid a$ set $b_1(A) = \infty$ and otherwise $b_1(A)$ is the unique integer $0 \leq b_1 < N$ satisfying $b_1 \equiv \bar{a}b \pmod{N}$ where $\bar{a}a \equiv 1 \pmod{N}$.

Corollary 1.4. *Assume the Ramanujan/Selberg Conjectures for GL_2 . Let U be a (nice) compact subset of $M_1(\mathbb{R})$ and $\epsilon > 0$. Then for N prime sufficiently large*

$$\frac{1/5 - \epsilon}{\log N} \leq \frac{|\{A \in M_N(\mathbb{Z}) : \pi(A) \in U, b_1(A) \text{ is prime}\}|}{|\{A \in M_N(\mathbb{Z}) : \pi(A) \in U\}|} \leq \frac{9/5 + \epsilon}{\log N}.$$

In particular the projections of the points $A \in M_N(\mathbb{Z})$ with $b_1(A)$ prime, become dense in $M_1(\mathbb{R})$.

We end the introduction with an outline of the contents of the sections and of the techniques that we use. The proofs of Theorems 1.1 and 1.3 use of sieve methods. These reduce sums over primes $\sum_{p \leq N} f(xu^p)$, to the study of linear sums over progressions $\sum_{n \leq N/d} f(xu^{nd})$, and bilinear sums $\sum_{n \leq \min(N/d_1, N/d_2)} f_1(xu^{nd_1})f_2(xu^{nd_2})$. A critical point in the analysis is to allow d to be as large as possible, this is measured by the level of distribution α ; $d \leq N^\alpha$ (respectively $\max(d_1, d_2) \leq N^\alpha$). The first type of sums are connected with equidistribution in (X, T) and the second type with joinings of (X, T^{d_1}) with (X, T^{d_2}) .

The effective rate of equidistribution of long pieces of unipotent orbits has been studied in the case of general compact quotients $\Gamma \backslash SL(2, \mathbb{R})$. For continuous such orbits this is due to Burger [7] while for discrete ones to Venkatesh [39]. Both make use of the spectral gap in the decomposition of $SL(2, \mathbb{R})$ acting by translations on $L^2(\Gamma \backslash SL(2, \mathbb{R}))$. One can quantify Venkatesh's method to obtain a positive level of distribution for the linear sums and then follow the analysis in the proof of Theorem 1.1 to obtain the analogues of Theorem 1.1 and Corollary 1.2 for such Γ 's (the constant 10 is replaced by a number depending on the spectral gap).

For $\Gamma \backslash SL(2, \mathbb{R})$ noncompact but of finite volume, due to the existence of periodic orbits of the horocycle flow one cannot formulate a simple uniform rate of equidistribution in Dani's Theorem. In a preprint

[38] A. Strömbergsson gives an effective version of Dani’s Theorem for continuous orbits in terms of the excursion rate of geodesics; he uses Burger’s approach.

We only learned of [38] after completing our formulation and treatment of an effective Dani theorem for the continuous flows, see Theorems 4.6 and 4.7 in section 4. One of the main results of this paper is an effective Dani Theorem for discrete unipotent orbits of (X, T^s) , see Theorems 4.12 and 5.2. These give a quantitative equidistribution w.r.t. algebraic measures of long pieces of such orbits and allowing s to be large. This discrete case is quite a bit more complex both in its formulation and its proof. It requires a series of basic Lemmas (see section 2) which use the action of $SL_2(\mathbb{Z})$ to give quantitative approximations of pieces of horocycle orbits by periodic horocycles, much like the approximation of reals by rationals in the theory of diophantine approximation. Critical to this analysis are various parameters associated with a given piece of horocycle orbit. The resulting approximations allow us to approach the level of distribution sums by taking f, f_1, f_2 to be automorphic forms and expanding them in Fourier series in the cusp. The burden of the analysis is in this way thrown onto the Fourier coefficients of these forms. This leads us to Theorem 4.12 which gives a suitable level of distribution in the linear sums and with which we can apply an upper bound sieve (Brun, Selberg) and deduce Theorem 1.1 and Corollary 1.2.

Theorem 1.3 involves a lower bound sieve and in particular a level of distribution for bilinear sums. This is naturally connected with effective equidistribution of 1-parameter unipotent orbits in $(SL(2, \mathbb{Z}) \backslash SL(2, \mathbb{R})) \times (SL(2, \mathbb{Z}) \backslash SL(2, \mathbb{R}))$ which is a well known open problem since Ratner’s paper [30]. For the special Hecke points that are taken in Theorem 1.3 our Fourier expansion approach converts the bilinear sums into sums of products of shifted coefficients of these automorphic forms (“Shifted convolution”). In Section 3 we review the spectral approach to this well studied problem; in particular we use the recent treatments in [3, 4] which are both convenient for our application and also allow for a critical improvement over [35] in the level aspect. Proposition 3.1 gives a slight improvement over [4] and also [29], in this q aspect, and it is optimal under the Ramanujan/Selberg Conjectures. Concerning the linear sums for these Hecke orbits we establish a level of $1/2$ (see the discussion at the end of section 7). This is the optimal level that can be proved by automorphic form/spectral methods. To analyze the sum over primes we use the sieve developed by Duke-Friedlander-Iwaniec [10]. For an asymptotics for the sum over primes (i. e. a “prime number theorem”) they require a level of distribution of $1/3$ for the bilinear

sums, given the level of $1/2$ that we have for the linear sums. Using the best bounds towards Ramanujan/Selberg Conjecture for GL_2/\mathbb{Q} we establish a level of $\alpha = 3/19$ for these bilinear sums. This falls short of the $1/3$ mark as well as the $1/5$ mark which is needed to get a lower bound in the sum over primes. However assuming the Ramanujan/Selberg Conjecture this does give a strong enough level of distribution to deduce Theorem 1.3.

2. HOROCYCLE APPROXIMATION

The group $G = SL(2, \mathbb{R})$ can be parametrized through its Lie algebra. The Lie algebra \mathfrak{g} are the 2×2 real matrices of zero trace. In this way, the so-called Iwasawa parametrization $g = h(x)a(y)k(\theta)$ with $x, y, \theta \in \mathbb{R}$ and

$$h(x) = \begin{bmatrix} 1 & x \\ 0 & 1 \end{bmatrix}, \quad a(y) = \begin{bmatrix} y^{\frac{1}{2}} & 0 \\ 0 & y^{-\frac{1}{2}} \end{bmatrix}, \quad k(\theta) = \begin{bmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{bmatrix}.$$

corresponds to the Lie algebra basis

$$(2.1) \quad R = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix}, \quad H = \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix}, \quad V = \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix}$$

in the sense that $h(x) = \exp(xR)$, $a(e^{2u}) = \exp(uH)$ and $k(\theta) = \exp(\theta V)$. Moreover it is unique when restricting θ to $[-\pi, \pi)$.

We can explicitly define a left G -invariant metric on G as

$$d_G(g, h) = \inf \left\{ \sum_{i=0}^{n-1} \psi(x_i, x_{i+1}) : x_0, \dots, x_n \in G; x_0 = g; x_n = h \right\}$$

with $\psi(x, y) = \min(\|x^{-1}y - I\|, \|y^{-1}x - I\|)$ and $\|\cdot\|$ any norm. Let us fix the norm

$$\left\| \begin{bmatrix} a & b \\ c & d \end{bmatrix} \right\| = \sqrt{2a^2 + (b+c)^2 + 4c^2 + 2d^2}.$$

for concreteness. Then, we can describe the metric in terms of the Iwasawa parametrization as

$$d_G s^2 = \frac{dx^2 + dy^2}{y^2} + d\theta^2$$

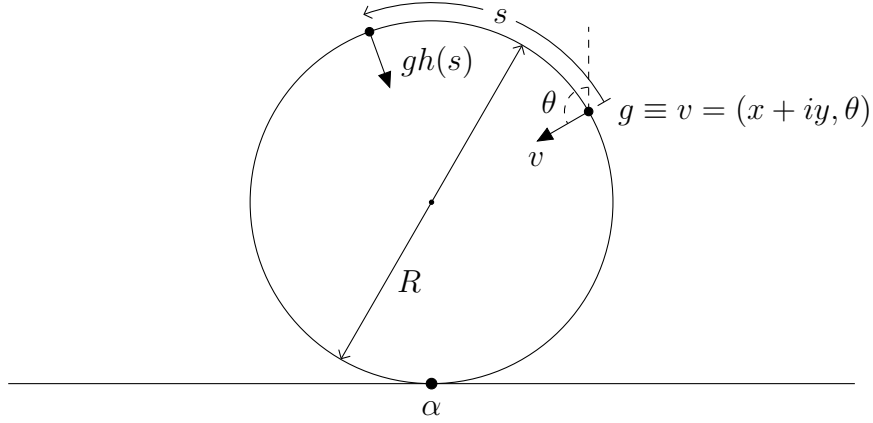
We can also write any Haar measure in G as a multiple of

$$d\mu_G = \frac{dx}{y} \frac{dy}{y} \frac{d\theta}{\pi}.$$

Since G is unimodular, this measure is both left and right G -invariant.

By sending $(x+iy, \theta)$ to $h(x)a(y)k(\theta/2)$ we see that $\{\pm I\} \backslash G$ endowed with this metric is isometric to the unit tangent bundle $T_1\mathbb{H}$ of the

FIGURE 1. The horocycle flow



Poincaré upper half-plane \mathbb{H} with the metric $ds^2 = y^{-2}(dx^2 + dy^2)$. We shall use both notations to refer to an element of $\{\pm I\} \backslash G$, and we shall even use $x + iy$ as a shorthand for $(x + iy, 0)$. In this way, we can express multiplication in $\{\pm I\} \backslash G$ as

$$(2.2) \quad \begin{bmatrix} a & b \\ c & d \end{bmatrix} (z, \theta) = \left(\frac{az + b}{cz + d}, \theta - 2 \arg(cz + d) \right).$$

Now, we define the discrete *horocycle flow* at distance s as the transformation $g \mapsto gh(s)$. The name comes from the fact that a horocycle is a circle in \mathbb{H} tangent to $\partial\mathbb{H}$, and the horocycle flow sends a point with tangent vector pointing towards the center of the horocycle S to the unique point at distance s “to the right” whose tangent vector also points to the center of S . In terms of our parametrization, we can write

$$(2.3) \quad gh(t + \cot \theta) = h\left(\alpha - \frac{Rt}{t^2 + 1}\right) a\left(\frac{R}{t^2 + 1}\right) k(-\operatorname{arccot} t)$$

for $g = h(x)a(y)k(\theta)$, where

$$R = y(\sin \theta)^{-2}$$

is the diameter of the horocycle and

$$\alpha = x - yW$$

its point of tangency with $\partial\mathbb{H}$, where $W = \cot \theta$ (see Figure 1).

Now, we consider the homogeneous space $X = \Gamma \backslash G$, with the metric induced from the one in G , namely $d_X(\Gamma g, \Gamma h) = \min_{\gamma \in \Gamma} d_G(\gamma g, h)$. Then, we have that X is isometric to $T_1(\Gamma \backslash \mathbb{H})$. So, considering the

identification $g = (z, \theta)$, we can set

$$D_X = \{(z, \theta) : |z| \geq 1, -1/2 \leq \Re z \leq 1/2, -\pi \leq \theta \leq \pi\}$$

as a fundamental domain for X .

We also have that the horocycle flow in G descends to X , and when doing so its behaviour becomes more complex. As we noted in the introduction, Dani proved that for any $\xi \in X$ the orbit generated by the horocycle flow, $\{\xi h(s)^n : n = 0, 1, 2, \dots\}$, is dense in either

- (i) a discrete periodic subset of a closed horocycle
- (ii) a closed horocycle
- (iii) the whole space,

and we can explicitly state which possibility happens in terms of $\xi = \Gamma g$: (i) for α and sR^{-1} rational numbers, (ii) for α rational and sR^{-1} irrational and (iii) for α irrational. Moreover, in each case the orbit becomes equidistributed in its closure w.r.t. the algebraic probability measure supported there.

If one considers the continuous version of the horocycle flow, $\{\xi h(t) : t \in \mathbb{R}_{\geq 0}\}$, only the possibilities (ii) and (iii) can occur, depending just on the rationality of α .

Our purpose in sections 4 and 5 is to analyze which of the possibilities is the “nearest” for a finite orbit

$$\{\xi h(s)^n : 0 \leq sn \leq T\};$$

or $\{\xi h(t) : t \in [0, T]\}$ in the continuous case, in terms of the parameters α , sR^{-1} and T . In preparation for that, we define some quantities associated to the piece of horocycle $\{\xi h(t) : t \in [0, T]\}$, which will be useful for applying spectral theory, and moreover allow us to decide between (i), (ii) and (iii) for ξ . Let $Y_T(g)$ the “Euclidean distance” from the piece of horocycle $P_{g,T} = \{gh(t) : t \in [0, T]\}$ to the border $\partial\mathbb{H}$, namely

$$Y_T(g) = \inf\{y : h(x)a(y)k(\theta) \in P_{g,T}\},$$

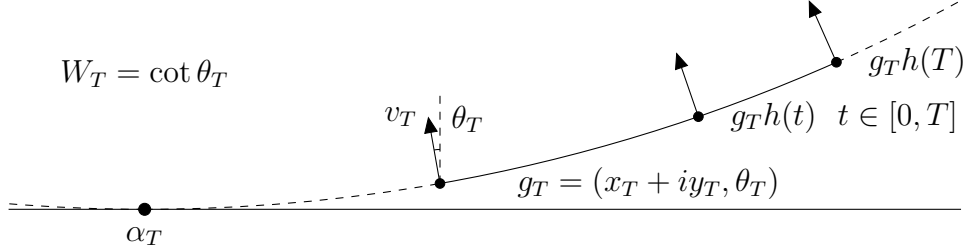
which coincides with the y associated to one of the extremes of the piece. Due to (2.3) we have

$$(2.4) \quad Y_T(h(x)a(y)k(\theta)) \asymp \min(y, \frac{R}{T^2}),$$

where the symbol \asymp is defined as follows.

Definition 2.1 (Notation for bounds). We shall use the notation $f = O(g)$ or $f \ll g$ meaning $|f| \leq C|g|$ for some constant $C > 0$; we shall also write $f \asymp g$ as a substitute for $f \ll g \ll f$. Finally, we shall use the notation $f < g^{O(1)}$ and $f < g^{-O(1)}$, with $g > 0$, meaning $|f| < g^C$ and $|f| < (g^{-1})^C$ respectively, for some constant $C > 0$. The

FIGURE 2. Piece of horocycle in highest position



implicit constant C will not depend on any other variable unless in a statement that contains an implication of the kind “if $f_1 = O(g_1)$ then $f_2 = O(g_2)$ ”; in that case the constant implicit in $O(g_2)$ depends on the one in $O(g_1)$. On the other hand, whenever we use the notation $g^{o(1)}$ or $g^{-1/|o(1)|}$, we mean that the function in $o(1)$ depends just on g and goes to zero as $g \rightarrow \infty$.

Now, the key concept is the following.

Definition 2.2 (Fundamental period). Let $g \in G$ and $T \geq 1$. We define the *fundamental period* of $\xi = \Gamma g$ at distance T as y_T^{-1} , where

$$y_T = y_T(\Gamma g) = \sup\{Y_T(\gamma g) : \gamma \in \Gamma\}.$$

The point of this definition is that we want to approximate our piece of horocycle by a closed one. Now, a closed horocycle has the shape $\Gamma a(y)H$ with H the closed subgroup $H = \{h(t) : t \in \mathbb{R}\}$; the period of this closed horocycle is y^{-1} —in the sense that $\Gamma a(y)h(\cdot)$ is a periodic function of period y^{-1} . We shall see that the closed horocycle of period y_T^{-1} is near to our original piece of horocycle $\{\Gamma gh(t) : t \in [0, T]\}$.

From the discontinuity of the Γ action one can deduce that the supremum in the definition of y_T is actually reached. Moreover, for $g \in U$ an open dense subset of G one can assure that this happens for a unique point $g_T = h(x)a(*)k(*)$ with $-1/2 < x \leq 1/2$, of the shape either γg or $\gamma gh(T)$ for some $\gamma \in \Gamma$. This defines the key parameters θ_T , α_T , y_T and W_T associated to g_T (see Figure 2).

Let us define the following equivalence relation in the space of pieces of horocycles of length T : $P \sim P'$ if there exists $\gamma \in \Gamma$ such that $P' = \gamma P$ as sets. We can identify a piece P with its point g which is nearest to $\partial\mathbb{H}$ —in case that both extremes are at the same distance from the border, we choose the left one. In this way, we can see the space of pieces of horocycles of length T as a subset of $PSL_2(\mathbb{R})$. Then,

we have just showed that

$$D_{X,T} = \overline{\{g_T : g \in U \cap D_X\}};$$

is a fundamental domain for this equivalence relation. We give an explicit description of this fundamental domain beginning by realizing y_T arithmetically; for that purpose we make the following definitions.

Definition 2.3 (Torus distance). Let $\alpha \in \mathbb{R}$. We define the *integral part* of α , and write $[\alpha]$, as the nearest integer to α . We also define its *fractional part* as $\{\alpha\} = \alpha - [\alpha]$. Finally, we define $\|\alpha\|$ as the absolute value of $\{\alpha\}$.

Definition 2.4 (Rational approximation). Let $\alpha \in \mathbb{R}$ and $U > 0$. We define

$$\kappa_U(\alpha) = \min\{m \in \mathbb{N} : \|m\alpha\| \leq U^{-1}\}.$$

We have the following property regarding the previous definition.

Lemma 2.5. *Let $\alpha \in \mathbb{R}$, $U > 0$ and $q \in \mathbb{N}$. If $\|q\alpha\| \leq 1/U$ then either $\|q\alpha\|$ is a multiple of $\|\kappa_U(\alpha)\alpha\|$ or $q \geq U/2$.*

Proof. We can assume $U > 2$. Writing $k = \kappa_U(\alpha)$, we have

$$\left|\alpha - \frac{a}{q}\right| \leq \frac{1}{Uq} \quad \left|\alpha - \frac{b}{k}\right| \leq \frac{1}{Uk}$$

for some integers a, b , $(b, k) = 1$. Thus, either $a/q = b/k$ or

$$\frac{1}{qk} \leq \left|\frac{a}{q} - \frac{b}{k}\right| \leq \frac{1}{Uq} + \frac{1}{Uk} \leq \frac{2}{Uk}$$

which implies $q \geq U/2$. In the former case, for some $\lambda \in \mathbb{N}$ we have $1/2 > |q\alpha - a| = \lambda|k\alpha - b|$, so that $\|q\alpha\| = |q\alpha - a|$. \square

From (2.4) we know that $y_T(\Gamma g)$ is always $\gg (2T)^{-2}$. It is natural that we can improve on that by translating g by different elements $\gamma \in \Gamma$. The following result reflects as far as we can go.

Lemma 2.6 (Period realization). *Let $T \geq 5$. For any g with $y = y(g) \gg 1$ we have*

$$y_T(\Gamma g) \asymp \min(y, RT^{-2}) + T^{-1} \min\left(\frac{U}{\kappa_U(\alpha)}, \frac{U^{-1}}{\|\kappa_U(\alpha)\alpha\|}\right)^2.$$

with $U = (T/R)^{1/2}$. Therefore, for any $g \in G$ we have

$$y_T(\Gamma g) \gg T^{-1}.$$

Proof. We can rewrite the Γ action (2.2) as

$$y_\gamma = \frac{1}{R} \frac{1}{c^2 + 2c\epsilon \cos \theta + \epsilon^2}, \quad R_\gamma = \frac{R}{(c\alpha + d)^2},$$

where γ is the matrix there—with R_γ and y_γ the corresponding parameters associated to γg —and $\epsilon = \pm(yR_\gamma)^{-1/2}$, the sign given by the one of $(c\alpha + d) \sin \theta$. Let us first treat the case $R \geq T$. We want to show that $y_T \asymp \min(y, RT^{-2})$. If $y \leq RT^{-2}$, since $y \gg 1$ it is clear that $y_T(\Gamma g) \asymp y$. If $y > RT^{-2}$, let us suppose that $y_T > CRT^{-2}$ for a large constant C . If g_T equals either γg or $\gamma gh(T)$, we deduce from (2.4) that $R_\gamma \gg CR$, so that $c \neq 0$ and $d = -[c\alpha]$ in the definition of R_γ . Therefore

$$|y_\gamma - \frac{1}{c^2 R}| \ll \frac{1}{c^2 R} (R_\gamma y)^{-1/2},$$

and since $R_\gamma y \gg CRy \gg CTy \gg CT \gg C$ we get that $y_\gamma \ll c^{-2}R^{-1} \ll RT^{-2}$ which is in contradiction with our assumption.

Now let us treat the case $T \geq R$. Choosing $c = \kappa_U(\alpha)$ and $d = -[\kappa_U(\alpha)\alpha]$ in the previous formulas, we have $R_\gamma = R\|\kappa_U(\alpha)\alpha\|^{-2} \geq T$ and then $y_\gamma \asymp \kappa_U(\alpha)^{-2}R^{-1}$. Considering (2.4) this clearly implies

$$Y_T(\gamma g) \asymp \min\left(\frac{1}{\kappa_U(\alpha)^2 R}, \frac{R}{T^2 \|\kappa_U(\alpha)\alpha\|^2}\right).$$

which equals the second term in the sum of the Lemma's statement. Since $y_T(\Gamma g) \geq Y_T(\gamma g)$, it only remains to prove that $y_T(\Gamma g) \leq CY_T(\gamma g)$ for some constant $C > 1$. Let us suppose this is not the case; then there exists $\gamma_* \in \Gamma$ such that $Y_T(\gamma_* g) > CY_T(\gamma g)$ and since $Y_T(\gamma g) \gg T^{-1}$ by (2.4) it follows that $\|c_* \alpha\| < U^{-1}$ which by definition of $\kappa_U(\alpha)$ implies $c_* \geq \kappa_U(\alpha)$. Also we could repeat the previous reasoning to show that

$$Y_T(\gamma_* g) \asymp \min\left(\frac{1}{c_*^2 R}, \frac{R}{T^2 \|c_* \alpha\|^2}\right).$$

Finally, by applying Lemma 2.5 with $q = c_*$ we have $Y_T(\gamma_* g) = O(Y_T(\gamma g))$ which is a contradiction. \square

Can we get something better than the bound $O(T)$ for the fundamental period? Not in general, think for instance of the case $g = a(T^{-1})$: then the piece of horocycle is just the closed horocycle of length T and $y_T = T^{-1}$ in this case. This is not a coincidence, as the following result shows.

Lemma 2.7 (Domain description). *For any $T \geq 2$ and $c > 0$, we define the set*

$$B_c(T) = \{h(x)a(y)k(\theta) : -1/2 \leq x \leq 1/2, y^{-1} < Tc^{-1}, |\theta| < T^{-1}c^{-1}\}.$$

We have that

$$B_{c_1}(T) \subset D_{X,T} \subset B_{c_2}(T)$$

for some positive constants c_1, c_2 and any $T \geq 2$.

Proof. The inclusion $D_{X,T} \subset B_{c_2}(T)$ comes just from Lemma 2.6 and the fact that by (2.3) the lowest point in any piece of horocycle of length T has $\theta = O(T^{-1})$.

Let $g = h(x)a(y)k(\theta) \in B_{c_1}(T)$ for some large c_1 . Suppose that $g \notin D_{X,T}$, then there exists $\gamma \in \Gamma$ with $c = c_\gamma \neq 0$ such that $\gamma g \in D_{X,T}$. So, since g and $g' = h(x')a(y')k(\theta') := gh(T)$ are in $P_{g,T}$, by Lemma 2.6 we have

$$\min(a(\gamma g), a(\gamma g')) = y_T(g) \geq \frac{\epsilon}{T}$$

for some $\epsilon > 0$. On the other hand, one can check that $|x' - x| \gg yT$ and $y' \asymp y$. Therefore, by (2.2) we have

$$\min(a(\gamma g), a(\gamma g')) \ll \frac{y}{\max(|cx + d|, |cx' + d|)^2} \ll \frac{y}{|c|^2|x - x'|^2} \ll \frac{1}{yT^2}$$

which is $O(\frac{1}{c_1 T})$, giving a contradiction for c_1 large enough. \square

Remark: This lemma implies that in any case, a large part of the piece of horocycle lies at height $O(y_T)$, thus showing that it is near to the closed horocycle of period y_T^{-1} . Moreover, it says that the fundamental domain is essentially

$$|x| \leq 1/2, \quad y^{-1} \ll T \ll |W|$$

with y^{-1} describing the period of the associated closed horocycle and $W = \cot \theta$ measuring the distance to it; $\theta = 0$ being the extreme case in which the piece is actually a closed horocycle.

From (2.3) we can write

$$gh(t + \cot \theta) = h(\alpha - Rt^{-1})a(Rt^{-2}) + O(t^{-1})$$

meaning that both points are at distance $O(|t|^{-1})$ in G . Since g_T can be either γg or $\gamma gh(T)$, we can always parametrize any piece of horocycle of length T as

$$(2.5) \quad h \left(\alpha_T + \frac{y_T W_T}{1 \pm t/W_T} \right) a \left(\frac{y_T}{(1 \pm t/W_T)^2} \right) + O \left(\frac{W_T^{-1}}{1 \pm t/W_T} \right) \quad t \in [0, T].$$

As we said before, we shall understand the piece of horocycle in terms of the parameters α_T, y_T and W_T (and s in the discrete case). But we are mainly interested in a fixed Γg and letting s change, and in that situation we will be able to express the results just in terms of α and sR^{-1} . To do that, we need to relate both kind of conditions.

First we write the necessary result for the continuous orbit. In order to read the following result, it is convenient to keep in mind that for any $g \in D_X$ with $y(g) < \delta^{-1}$ either g or $\begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix} g$ satisfies that its coefficients (R, α, y^{-1}) are bounded by $O(\delta^{-3})$.

Lemma 2.8 (Fundamental period and continuous Dani). *Let $0 < \delta < 1/2$ and $g_0 \in G$ with coefficients (R, α, y^{-1}) bounded by δ^{-1} . Then, $y_T^{-1} < \delta^{-O(1)}$ if and only if there exists a positive integer $q < \delta^{-O(1)}$ such that $\|q\alpha\| < \delta^{-O(1)}T^{-1}$.*

Remark. For our application to orbits at prime values we will always have $\delta^{-1} = (\log T)^A$ for some constant $A > 0$ in this lemma as well as in later statements. Moreover, throughout the whole paper we can assume that $(\delta^{-1})^c < T$ for some large constant $c > 1$, because otherwise the results are trivial.

Proof. Since we always have $y_T^{-1} \ll T$, we can assume $\delta^{-O(1)} \ll T$ in the proof. One can check that the constants implicit in the statement of Lemma 2.6 have a polynomial dependency on the constant in $y \gg 1$. Thus, applying it to our case, since $y > \delta$, we will get $\delta^{-O(1)}$ as implicit constants. Therefore, since $\delta < R, y < \delta^{-1}$, we have that Lemma 2.8 is a direct consequence of Lemma 2.6 and Lemma 2.5. \square

Now we write the analogous result for a discrete orbit. Probably it is better to skip it on first reading, at least until one arrives at Theorem 4.12. Before, let us fix our notation for inverses modulo a number.

Definition 2.9 (Modular inverses). Let q be an integer different from zero. For any $a \in \mathbb{Z}$ coprime to q , we define \bar{a} as the integer between 1 and q such that $\bar{a}a \equiv 1 \pmod{q}$.

Lemma 2.10 (Fundamental period and discrete Dani). *Let $N \geq 1$, $0 < \delta < 1/2$. Let $g_0 \in G$ with coefficients $(R_0, \alpha_0, y_0^{-1})$ bounded by δ^{-1} . Then the following statements are equivalent, unless $s < \delta^{-O(1)}N^{-1}$:*

- (i) *There exists $q \in \mathbb{N}$ and $\gamma \in \Gamma$ with the coefficients of γ and q bounded by $(\delta^{-1}\tau(q_2))^{O(1)}$ such that $g = \gamma g_0$ satisfies*

$$\|q \frac{s}{R}\| < (\delta^{-1}\tau(q_2))^{O(1)}N^{-1}, \quad \left\| \left[q \frac{s}{R} \right] \alpha \right\| < (\delta^{-1}\tau(q_2))^{O(1)}(sN^2)^{-1}.$$

where \tilde{q}_2 and q_2 are the denominators in the expressions as reduced fractions of $\frac{[[q \frac{s}{R}]\alpha]}{[q \frac{s}{R}]}$ and $\frac{[q \frac{s}{R}]}{qq_2^2}$ respectively.

- (ii) *There exists $y < 1$ and an integer q' with $y^{-\frac{1}{2}}/(\tau(q'_2)\delta^{-1})^{O(1)} < q' < y^{-\frac{1}{2}}(\tau(q'_2)\delta^{-1})^{O(1)}$ such that $\Gamma g_0 = \Gamma(x + iy, \theta)$ with*

$$\|q'x\| + N\|q'sy\| + (sN)^2q'|\theta|y < y^{\frac{1}{2}}(\delta^{-1}\tau(q'_2))^{O(1)},$$

where q'_2 is the denominator in the expression as a reduced fraction of $\frac{[q'sy]}{q'}$.

Remarks. The proof actually gives $q'_2 = q_2$. One can check that conditions in (i) assure that the coefficients of γ and q are always bounded by $(1+s)^{o(1)}\delta^{-O(1)}$. One can see in the proof that from (i) we actually get a y in (ii) satisfying $y \gg (1+s)^{-2-o(1)}$.

Proof. Let us begin by demonstrating that (ii) implies (i). Let us write explicitly the Iwasawa decomposition

$$(2.6) \quad g = h(x)a(y)k(\theta) = \begin{bmatrix} -xy^{-\frac{1}{2}}\sin\theta + y^{\frac{1}{2}}\cos\theta & xy^{-\frac{1}{2}}\cos\theta + y^{\frac{1}{2}}\sin\theta \\ -y^{-\frac{1}{2}}\sin\theta & y^{-\frac{1}{2}}\cos\theta \end{bmatrix},$$

or equivalently, with $W = \cot\theta$,

$$(2.7) \quad g = \begin{bmatrix} -R^{-\frac{1}{2}}\alpha & R^{-\frac{1}{2}}\alpha W + R^{\frac{1}{2}} \\ -R^{-\frac{1}{2}} & R^{-\frac{1}{2}}W \end{bmatrix}.$$

From (ii) we have that

$$x = \frac{a'_1}{q'_1} + yO(M) \quad q'_1 \mid q', (a'_1, q'_1) = 1, M = (\tau(q'_2)\delta^{-1})^{O(1)}.$$

On the other hand, suppose that $y^{-1} > 4M^2s^2$. Then, $s\sqrt{y} < 1/2M$, so $q'sy = (q'\sqrt{y})(s\sqrt{y}) < M(1/2M) < 1/2$ which implies $\|q'sy\| = q'sy$, thus $q'_2 = 1$. But then, (ii) gives $M = \delta^{-O(1)}$ and $s < \delta^{-O(1)}/N$.

All this means that we can assume $y^{-1} < 4M^2s^2$. Then considering the matrix

$$\gamma_{a'_1/q'_1} = \begin{bmatrix} -\overline{a'_1} & d'_1 \\ q'_1 & -a'_1 \end{bmatrix}$$

in Γ , we have

$$g^* = \gamma_{a'_1/q'_1}g = \begin{bmatrix} -\overline{a'_1}y^{\frac{1}{2}}(1 + (yq'_1\overline{a'_1})^{-1}\theta O(M)) & * \\ q'_1y^{\frac{1}{2}}(1 + \theta O(M)) & q'_1y^{\frac{1}{2}}O(M) \end{bmatrix}.$$

Since $|\theta| < M(sN)^{-2}$, this gives

$$R^{-1} = q_1^2y\left(1 + \frac{O(M)}{s^2N^2}\right), \quad \alpha = -\frac{\overline{a'_1}}{q'_1} + \frac{1}{q_1^2y} \frac{O(M)}{s^2N^2},$$

with $\alpha = \alpha(g^*)$, $R = R(g^*)$ —the corresponding parameters associated to g^* . It also gives the inequalities $y(g^*)^{-1} \ll Mq_1^2y \ll M$, so if $q_1^2y < M^{-1}$ then g^* is in the fundamental domain, but then since $y(g_0) < \delta^{-1}$ this implies $q_1^2y \gg M^{-1}$. Now,

$$\frac{q'}{q'_1} \frac{s}{R} = q'_1(q'sy) + \frac{O(M)}{sN^2}$$

so that $q = q'/q'_1 < M$ satisfies

$$\|q \frac{s}{R}\| < \frac{M}{N}, \quad [q \frac{s}{R}] = q'_1 [q' sy]$$

and also, using the expression for α , we have

$$[[q \frac{s}{R}]\alpha] = -\overline{a'_1} [q' sy], \quad \|[q \frac{s}{R}]\alpha\| < \frac{M}{sN^2}.$$

On the other hand

$$\frac{[[q \frac{s}{R}]\alpha]}{[q \frac{s}{R}]} = \frac{-\overline{a'_1}}{q'_1}, \quad \frac{[q \frac{s}{R}]}{qq_1^2} = \frac{[q' sy]}{q'}$$

so $\tilde{q}_2 = q'_1$ and $q_2 = q'_2$. Finally, we have that the coefficients of the matrix g^* are bounded by M , so since $\Gamma g^* = \Gamma g_0$ we have $g^* = \gamma g_0$, γ with coefficients bounded by M and (i) follows.

Now let us prove that (i) implies (ii). Let $s = s(g)$, $R = R(g)$ and $W = W(g)$. We have

$$(2.8) \quad q \frac{s}{R} = [q \frac{s}{R}] + \frac{M}{N}, \quad \alpha = \frac{a_1}{q_1} + \frac{M}{s^2 N^2}$$

with $M = (\tau(q_2)\delta^{-1})^{O(1)}$ for some coprime integers a_1, q_1 , with $q_1 \mid [q \frac{s}{R}]$. We also have $|W|R^{-\frac{1}{2}} = y(g)^{-\frac{1}{2}} < M$. Let us consider $g_* = \gamma_{a_1/q_1} g = (x + iy, 0)k(\theta)$. Via (2.7)

$$(2.9) \quad g_* = \begin{bmatrix} * & -R^{\frac{1}{2}} \overline{a_1} (1 + \frac{M}{a_1 q_1}) \\ \frac{q_1 M}{s^2 N^2} & R^{\frac{1}{2}} q_1 (1 + \frac{M}{s^2 N^2}) \end{bmatrix},$$

so considering the components in the lower row we have (by (2.6))

$$|\tan \theta| < \frac{M}{s^2 N^2}, \quad y = \frac{1}{R q_1^2} (1 + \frac{M}{s^2 N^2})$$

so taking $q' = qq_1$ we have

$$q' sy = \frac{1}{q_1} q \frac{s}{R} (1 + \frac{M}{s^2 N^2}) = \frac{[q \frac{s}{R}]}{q_1} (1 + \frac{M}{s^2 N^2}) = \frac{[q \frac{s}{R}]}{q_1} + \frac{M}{q_1 N}$$

so $\|q' sy\| < My^{\frac{1}{2}} N^{-1}$. We also have $M^{-1} y^{-\frac{1}{2}} < q' < My^{-\frac{1}{2}}$ and $(sN)^2 q' y |\theta| < y^{\frac{1}{2}} M$. Now, by the second column in (2.9) and by (2.6) we have

$$x + y \tan \theta = \frac{-R^{-\frac{1}{2}} \overline{a_1} (1 + \frac{M}{a_1 q_1})}{R^{\frac{1}{2}} q_1 (1 + \frac{M}{s^2 N^2})} = -\frac{\overline{a_1}}{q_1} + \frac{M}{q_1^2}$$

so $q'x = -q\overline{a_1} + y^{\frac{1}{2}} M$ and then $\|q'x\| < My^{\frac{1}{2}}$. Finally

$$\frac{a_1}{q_1} = \frac{[[q \frac{s}{R}]\alpha]}{[q \frac{s}{R}]}, \quad \frac{[q' sy]}{q'} = \frac{[q \frac{s}{R}]}{qq_1^2}$$

so $q_1 = \tilde{q}_2$ and $q_2' = q_2$. □

3. AUTOMORPHIC FORMS

A key to our analysis of the averages of functions along pieces of long periodic horocycles is the use of automorphic forms. We will need the sharpest known estimates for periods of the type

$$(3.1) \quad \int_0^1 f\left(\begin{bmatrix} 1 & x \\ 0 & 1 \end{bmatrix} g\right) e(-hx) dx$$

where $h \in \mathbb{Z}$, $g \in G$ and f is a mildly varying function on $L_0^2(\Gamma_0(q)\backslash G)$ and $q \geq 1$ is an integer. Here the subzero indicates that $\int_{\Gamma_0(q)\backslash G} f(g) dg = 0$ and $\Gamma_0(q)$ denotes the standard Hecke congruence subgroup of $SL_2(\mathbb{Z})$.

If $h = 0$ and $g = a(y)$ then (3.1) measures the equidistribution of the closed horocycle of period $1/y$ in $\Gamma_0(q)\backslash G$. This can be studied using Eisenstein series and the precise rate of equidistribution is tied up with the Riemann Hypothesis (see [32] for the case $q = 1$, the rate is $O_\epsilon(y^{\frac{1}{4}+\epsilon})$ if and only if RH is true). For $h \neq 0$ the size of (3.1) is controlled by the full spectral theory of $L_0^2(\Gamma_0(q)\backslash G)$ and in particular the Ramanujan/Selberg conjectures for GL_2/\mathbb{Q} (see the appendix to [35]). In this case (3.1) is closely related to the much studied shifted convolution problem. In order to bench-mark the upper bound that we are aiming for, consider the case that $f \in L_0^2(\Gamma_0(q)\backslash G)$ is K invariant, that is $f(gk(\theta)) = f(g)$. Thus $f = f(z)$ with $z = x + iy \in \mathbb{H}$ and (3.1) is the period

$$(3.2) \quad \int_0^1 f(x + iy) e(-hx) dx.$$

We are assuming that f is smooth and in $L_0^2(\Gamma_0(q)\backslash \mathbb{H})$. We use the Sobolev norms

$$(3.3) \quad \|f\|_{W^{2d}}^2 = \int_{\Gamma_0(q)\backslash \mathbb{H}} |f(z)|^2 dA(z) + \int_{\Gamma_0(q)\backslash \mathbb{H}} |\Delta^d f(z)|^2 dA(z)$$

where $d \geq 0$ is an integer, $dA = \frac{dx dy}{y^2}$ is the area form and Δ the Laplacian for the hyperbolic metric.

Expanding f in the Laplacian spectrum of $\Gamma_0(q)\backslash \mathbb{H}$ (see [19]) with ϕ_j an orthonormal basis of cusp forms and $E_j(z, s)$ the corresponding Eisenstein series, $j = 1, 2, \dots, \nu(q)$, yields

$$(3.4) \quad f(z) = \sum_{j \neq 0} \langle f, \phi_j \rangle \phi_j(z) + \sum_{j=1}^{\nu} \frac{1}{2\pi} \int_{-\infty}^{\infty} \langle f, E_j(\cdot, \frac{1}{2} + it) \rangle E_j(z, \frac{1}{2} + it) dt.$$

Note that $\phi_0(z) = 1/\sqrt{\text{Vol}(\Gamma_0(q)\backslash\mathbb{H})}$ does not appear since we are assuming that

$$(3.5) \quad \int_{\Gamma_0(q)\backslash\mathbb{H}} f(z) dA(z) = 0.$$

Let λ_j denote the Laplacian eigenvalue of ϕ_j and write $\lambda_j = \frac{1}{4} + t_j^2$. ϕ_j may be expanded in a Fourier series (see [19])

$$(3.6) \quad \phi_j(z) = \sum_{n \neq 0} \rho_j(n) y^{\frac{1}{2}} K_{it_j}(2\pi|n|y) e(nx).$$

Using the Atkin-Lehner level raising operators one can choose the orthonormal basis ϕ_j to consist of new forms and old forms and then normalizing the coefficients in (3.4) amounts to bounding the residues of $L(s, \phi_j \times \phi_j)$ at $s = 1$. These are known to be bounded above and below by $(\lambda_j q)^{\pm\epsilon}$ respectively (see [17, 20]). In this way one has the bound (see [22] for details when q is square free which essentially is the case of interest to us, and [3, 4] for the general q)

$$(3.7) \quad \rho_j(h) \ll_{\epsilon} (\lambda_j q h)^{\epsilon} q^{-\frac{1}{2}} \cosh\left(\frac{\pi t_j}{2}\right) h^{\theta}$$

for any $\epsilon > 0$, and where θ is an acceptable exponent for the Ramanujan/Selberg Conjecture. $\theta = 7/64$ is known to be acceptable [24] while $\theta = 0$ is what is conjectured to be true.

It follows from (3.7), (3.4) and (3.6) that

$$(3.8) \quad \int_0^1 f(x+iy) e(-hx) dx \ll_{\epsilon} \frac{y^{\frac{1}{2}}(qh)^{\epsilon} h^{\theta}}{\sqrt{q}} \sum_{j \neq 0} |\langle f, \phi_j \rangle| |K_{it_j}(2\pi|h|y) \lambda_j^{\epsilon}| \cosh\left(\frac{\pi t_j}{2}\right) + \text{cts}$$

where the term ‘‘cts’’ is a similar contribution from the continuous spectrum and for which $\theta = 0$ is known since the coefficients of Eisenstein series are unitary divisor sums.

The Bessel function K satisfies the inequalities, say for $v \ll 1$ (see [1])

$$(3.9) \quad \begin{aligned} K_{\nu}(v) &\ll_{\epsilon} v^{-\nu-\epsilon} & 0 \leq \nu \leq 1/2, \\ K_{it}(v) &\ll_{\epsilon} v^{-\epsilon} & 0 \leq t \leq 1, \\ e^{\frac{\pi}{2}t} K_{it}(v) &\ll_{\epsilon} v^{-\epsilon} & 1 \leq t < \infty. \end{aligned}$$

Hence from (3.7) we have that for $0 < |hy| \ll 1$

(3.10)

$$\int_0^1 f(x+it)e(-hx) dx \ll_\epsilon \frac{y^{\frac{1}{2}-\theta-\epsilon}(qh)^\epsilon}{\sqrt{q}} \sum_{j \neq 0} |\langle f, \phi_j \rangle| \lambda_j^\epsilon$$

$$(3.11) \quad \ll_\epsilon \frac{y^{\frac{1}{2}-\theta-\epsilon}(qh)^\epsilon}{\sqrt{q}} \left(\sum_{j \neq 0} |\langle f, \phi_j \rangle|^2 \lambda_j^2 \right)^{\frac{1}{2}} \left(\sum_{j \neq 0} \lambda_j^{-2+2\epsilon} \right)^{\frac{1}{2}}.$$

Weyl's law for $\Gamma_0(q) \backslash \mathbb{H}$ gives the uniform bound

$$(3.12) \quad \sum_{\lambda_j \leq \lambda} 1 \ll \text{Vol}(\Gamma_0(q) \backslash \mathbb{H}) \lambda$$

for $\lambda \geq 1$ and $q \geq 1$. Hence

$$(3.13) \quad \sum_{j \neq 0} \lambda_j^{-2+2\epsilon} \ll \text{Vol}(\Gamma_0(q) \backslash \mathbb{H}) = q \prod_{p|q} \left(1 + \frac{1}{p}\right) \text{Vol}(\Gamma_0(1) \backslash \mathbb{H}) \ll q^{1+\epsilon}.$$

We conclude that for $|h|y \ll 1$ and $\epsilon > 0$

$$(3.14) \quad \int_0^1 f(x+iy)e(-hx) dx \ll_\epsilon y^{\frac{1}{2}-\theta-\epsilon} q^\epsilon \|f\|_{W^2}.$$

For $\theta = 0$ (3.14) is sharp, that is it cannot be improved. To see this take for example $h = 1$ in

$$(3.15) \quad \int_0^1 f(x+iy)e(-hx) dx = y^{\frac{1}{2}} \sum_{j \neq 0} \langle f, \phi_j \rangle \rho_j(h) K_{it_j}(2\pi|h|y) + \text{cts.}$$

Choosing

$$f(z) = \sum_{0 \leq t_j \leq 1} \overline{\rho_j(1) K_{it_j}(2\pi y)} \phi_j(z)$$

yields

$$\|f\|_{W^2}^2 \asymp \sum_{0 \leq t_j \leq 1} |\rho_j(1)|^2 |K_{it_j}(2\pi y)|^2$$

while

$$(3.16) \quad \int_0^1 f(x+iy)e(-hx) dx = y^{\frac{1}{2}} \sum_{0 \leq t_j \leq 1} |\rho_j(1)|^2 |K_{it_j}(2\pi y)|^2 \asymp y^{\frac{1}{2}} \|f\|_{W^2}^2.$$

Recall that Iwaniec [20] shows that

$$(3.17) \quad |\rho_j(1)| \gg_\epsilon (\lambda_j q)^{-\epsilon} \cosh\left(\frac{\pi t_j}{2}\right) / \sqrt{q},$$

hence changing y a little if need be to make sure that $|K_{it_j}(2\pi y)| \gg 1$ for most $t_j \leq 1$, we see that for this f

$$(3.18) \quad \|f\|_{W^2} \gg_\epsilon (q\lambda_j)^{-\epsilon}.$$

It follows from (3.18) and (3.16) that (3.14) is sharp when $\theta = 0$.

On the other hand if $\theta > 0$ then (3.14) can be improved for q in the range $y^{-\frac{1}{2}+2\theta} \leq q \leq y^{-\frac{1}{2}}$. To do so we estimate the j -sum in (3.16) using the Kuznetsov formula for $\Gamma_0(q)\backslash\mathbb{H}$, rather than invoking the sharpest bound (3.7) for the individual coefficients. One applies Kuznetsov with suitably chosen positive (on the spectral side) test functions and then using only Weil's upper bound for the Kloosterman sums that appear on the geometric side of the formula, we get (see [20]); for $X \geq 1$, $h \neq 0$

$$(3.19) \quad \sum_{0 < \lambda_j < \frac{1}{4}} |\rho_j(h)|^2 X^{4|t_j|} \ll 1 + \frac{|h|^{\frac{1}{2}} X}{q}$$

and for $\lambda \geq 1$

$$(3.20) \quad \sum_{\frac{1}{4} \leq \lambda_j \leq \lambda} |\rho_j(h)|^2 \cosh(\pi t_j) \leq \lambda \left(1 + \frac{|h|^{\frac{1}{2}}}{q}\right).$$

Now for $0 < |hy| \ll 1$, using (3.9) and (3.15) we have

$$(3.21) \quad \int_0^1 f(x+iy)e(hx) dx \ll y^{\frac{1}{2}} \left(\sum_{0 < \lambda_j < \frac{1}{4}} |\rho_j(h)|^2 |hy|^{-2|t_j|} \right)^{\frac{1}{2}} \left(\sum_{\lambda_j < \frac{1}{4}} |\langle f, \phi_j \rangle|^2 \right)^{\frac{1}{2}} \\ + y^{\frac{1}{2}} \left(\sum_{\lambda_j \geq \frac{1}{4}} |\rho_j(h)|^2 |\lambda_j|^{-2} \right)^{\frac{1}{2}} \left(\sum_{\lambda_j \geq \frac{1}{4}} |\langle f, \phi_j \rangle|^2 \lambda_j^2 \right)^{\frac{1}{2}}.$$

Taking $X = |hy|^{-\frac{1}{2}}$ in (3.19) and applying it to the first term on the right hand side of (3.21) and applying (3.20) for the second term we arrive at

$$(3.22) \quad \int_0^1 f(x+iy)e(-hx) dx \ll y^{\frac{1}{2}} \|f\|_{W^2} \left(1 + \frac{y^{-\frac{1}{4}}}{\sqrt{q}} + \frac{|h|^{\frac{1}{4}}}{\sqrt{q}}\right) \ll y^{\frac{1}{2}} \|f\|_{W^2} \left(1 + \frac{y^{-\frac{1}{4}}}{\sqrt{q}}\right)$$

since $|h| \ll y^{-1}$.

We combine (3.14) with (3.22) to arrive at our strongest unconditional estimate; for $0 < |hy| \ll 1$ and $\epsilon > 0$,

$$(3.23) \quad \int_0^1 f(x+iy)e(-hx) dx \ll_\epsilon (y^{-1}q)^\epsilon y^{\frac{1}{2}} \|f\|_{W^2} \min(y^{-\theta}, 1 + \frac{y^{-\frac{1}{4}}}{\sqrt{q}}).$$

For the application to the level of equidistribution in type I and II sums connected to sieving as we do in section 7, (3.23) can be improved slightly when summing over h in certain ranges; we leave that discussion for section 7.

In generalizing (3.14) and (3.23) to functions on $\Gamma_0(q)\backslash G$ it is natural to use the spectral decomposition of $L^2_0(\Gamma_0(q)\backslash G)$ into irreducibles under the action of G by right translation on this space. This is the path chosen in [3] and [4] and we will follow these treatments closely modifying it as needed for our purposes.

The Lie algebra \mathfrak{g} of G consists of the two by two matrices of trace zero. An element X of \mathfrak{g} gives rise to a left invariant differential operator on $C^\infty(G)$;

$$(3.24) \quad D_X f(g) = \frac{d}{dt} f(g \exp(tX))_{t=0}.$$

The operators D_H, D_R, D_L with $H, R, L = R - V$ in (2.1) generate the algebra of left invariant differential operators on $C^\infty(G)$. For $k \geq 0$ define the Sobolev k -norms on functions on $\Gamma_0(q)\backslash G$ by

$$(3.25) \quad \|f\|_{W^k} := \sum_{\text{ord}(D) \leq k} \|Df\|_{L^2(\Gamma_0(q)\backslash G)},$$

where D ranges over all monomials in D_H, D_R and D_L of degree at most k . For a unitary representation π of G , \mathfrak{g} acts on the associated Hilbert space V_π and one can define Sobolev norms of smooth vectors in the same way. The center of the algebra of differential operators is generated by the Casimir operator ω which in our basis is given by

$$(3.26) \quad \omega = -\frac{1}{4}(D_H D_H + 2D_R D_H + 2D_H D_R).$$

In Iwasawa coordinates it is given by

$$(3.27) \quad \omega = -y^2 \left(\frac{\partial^2}{\partial x^2} + \frac{\partial^2}{\partial y^2} \right) + y \frac{\partial^2}{\partial x \partial \theta}.$$

We decompose $L^2_0(\Gamma_0(q)\backslash G)$ under right translation by $g \in G$ into irreducible subrepresentations π of G . This extends the decomposition in (3.4) to

$$(3.28) \quad L^2_0(\Gamma_0(q)\backslash G) = \int_{\widehat{G}} V_\pi d\mu(\pi),$$

where π ranges over \widehat{G} the unitary dual of G and μ is a measure on \widehat{G} (it depends on q of course) which corresponds to this spectral decomposition. It consists of a cuspidal part on which the spectrum is

discrete (cuspidal $\Leftrightarrow \int_0^1 f\left(\begin{bmatrix} 1 & x \\ 0 & 1 \end{bmatrix} g\right) dx = 0, \forall g$) and a continuous part corresponding to an integral over unitary Eisenstein series. According to (3.28) for $f \in L_0^2(\Gamma_0(q)\backslash G)$

$$(3.29) \quad f = \int_{\widehat{G}} f_\pi d\mu(\pi)$$

and

$$(3.30) \quad \|f\|_{L_0^2(\Gamma_0(q)\backslash G)}^2 = \int_{\widehat{G}} \|f_\pi\|_{V_\pi}^2 d\mu(\pi).$$

Since π is irreducible and ω commutes with the G -action it follows that ω acts on smooth vectors in V_π by a scalar which we denote by λ_π . Weyl's law for the principal and complementary series representations of \widehat{G} together with the dimensions of the discrete series representation in $L_0^2(\Gamma_0(q)\backslash G)$ imply that for $T \geq 1$,

$$(3.31) \quad \mu\{\pi : |\lambda_\pi| \leq T\} \ll \text{Vol}(\Gamma_0(q)\backslash G)(1+T).$$

Note that for a, b non-negative integers and $f \in V_\pi$ smooth we have from $\omega f = \lambda_\pi f$ that

$$(3.32) \quad \|f\|_{W^a} \leq (1 + \lambda_\pi)^{-b} \|f\|_{W^{a+b}}.$$

Using the Hecke operators, Atkin-Lehner theory and choosing suitable bases to embed the space of new forms of a given level $t \mid q$ to level q (see [22, 4]) we can further decompose $L_0^2(\Gamma_0(q)\backslash G)$ into an orthogonal direct sum/integral of irreducibles on which the Fourier coefficients satisfy the analogue of (3.7). In more detail if π is an irreducible constituent in the above decomposition and π is of level q (by which we mean that V_π has a vector which is a classical modular new form of level q) then the Whittaker functional

$$(3.33) \quad W_f(y) := \int_0^1 f(h(t)a(y))e(-t) dt$$

is non-zero as a function of f a smooth vector in V_π . Moreover as is shown in [4]

$$(3.34) \quad \frac{\langle f, f \rangle}{\text{Vol}(\Gamma_0(q)\backslash G)} = c_\pi \langle W_f, W_f \rangle$$

where the second inner product is the standard inner product on $L^2(\mathbb{R}^*, \frac{dy}{y})$ and in the Kirillov model for π and c_π satisfies ¹

$$(3.35) \quad (\lambda_\pi q)^{-\epsilon} \ll_\epsilon c_\pi \ll_\epsilon (\lambda_\pi q)^\epsilon.$$

¹note our normalization of the Sobolev norms on $C^\infty(\Gamma_0(q)\backslash G)$ and that of [4] differ by a factor of $\text{Vol}(\Gamma_0(q)\backslash G)$

Moreover the m -th coefficient ($m \neq 0$) for our period integral satisfies the relation; for $f \in V_\pi$

$$(3.36) \quad \int_0^1 f(h(t)a(y))e(-mt) dt = \frac{\lambda_\pi(|m|)}{\sqrt{|m|}} W_f(|m|y).$$

Here $\lambda_\pi(m)$ is the eigenvalue of the m -th Hecke operator on V_π . Recall that we are assuming that these satisfy

$$(3.37) \quad \lambda_\pi(m) \ll \tau(|m|)|m|^\theta$$

and similarly that the Laplace eigenvalue $\lambda_\pi = \frac{1}{4} + t_\pi^2$ (in the case that π is spherical), if $it_\pi > 0$ then

$$(3.38) \quad it_\pi \leq \theta.$$

The invariant differential operators on V_π induce an action on W_f , namely

$$D_X W_f = W_{D_X f}$$

and using $D_H = 2y \frac{d}{dy}$, $D_R = 2\pi iy$ and $D_L = \frac{1}{2\pi i} \left(\frac{-\lambda_\pi}{y} + y \frac{d^2}{dy^2} \right)$ we get using the various normalizations and (3.34), (3.35), (3.36), (3.37) and (3.38) (see [4]) that for $b \geq 0$ fixed, $m \neq 0$ and $f \in V_\pi$

$$(3.39) \quad \int_0^1 f(h(t)a(y))e(-mt) dt \ll_{\epsilon, b} \left| \frac{\lambda_\pi q}{my} \right|^\epsilon \frac{1}{\sqrt{q}} \frac{\tau(|m|)y^{\frac{1}{2}-\theta}}{1+|my|^b} \|f\|_{W^{4+b}}.$$

As mentioned above, in [4] it is shown how this analysis may be extended to any π not just the ones of level q . The same bounds (3.39) may be established for the Eisenstein spectrum and in that case one can take $\theta = 0$. Hence if $f \in L_0^2(\Gamma_0(q)\backslash G)$ is smooth we may apply (3.29), (3.30) and (3.31) to arrive at

$$(3.40) \quad \begin{aligned} \int_0^1 f(h(t)a(y))e(-mt) dt &= \int_{\widehat{G}} \int_0^1 f_\pi(h(t)a(y))e(-mt) dt d\mu(\pi) \\ &\ll_{\epsilon, b} \frac{q^{\epsilon-\frac{1}{2}}}{|my|^\epsilon} \frac{\tau(|m|)y^{\frac{1}{2}-\theta}}{1+|my|^b} \int_{\widehat{G}} \|f\|_{W^{4+b}} \lambda_\pi^\epsilon d\mu(\pi) \\ &\ll_{\epsilon, b} \left| \frac{q}{my} \right|^\epsilon \frac{\tau(|m|)y^{\frac{1}{2}-\theta}}{1+|my|^b} \|f\|_{W^{6+b}}. \end{aligned}$$

For $m = 0$ only the Eisenstein series enter and a similar analysis of the constant term yields that for smooth f 's in $L_0^2(\Gamma_0(q)\backslash G)$

$$(3.41) \quad \int_0^1 f(h(t)a(y)) dt \ll q^\epsilon y^{\frac{1}{2}-\epsilon} \|f\|_{W^6}.$$

(3.40) and (3.41) are the generalizations of (3.14) to $\Gamma_0(q)\backslash G$ and as we have noted these bounds are essentially contained in [4]. The slight

improvement (3.23) can also be incorporated into this $\Gamma_0(q)\backslash G$ analysis and this gives our main estimate for the period integral:

Proposition 3.1. *For $\epsilon > 0$, $b \geq 0$ fixed and θ admissible for the Ramanujan/Selberg conjecture for GL_2/\mathbb{Q} , if $f \in C^\infty(\Gamma_0(q)\backslash G)$ and $m \in \mathbb{Z}$, then for $m \neq 0$*

$$\int_0^1 f(h(t)a(y))e(-mt) dt \ll_{\epsilon,b} \left(\frac{q}{|m|y}\right)^\epsilon \frac{\tau(|m|)y^{\frac{1}{2}}\|f\|_{W^{6+b}}}{1+|my|^b} \min(y^{-\theta}, \frac{1}{y^\epsilon} + \frac{y^{-\frac{1}{4}}}{y^\epsilon q^{\frac{1}{2}}})$$

while for $m = 0$, assuming $\int_X f d\mu_G = 0$,

$$\int_0^1 f(h(t)a(y)) dt \ll_\epsilon (qy^{-1})^\epsilon y^{\frac{1}{2}}\|f\|_{W^6}.$$

Remark. For $\Gamma = \Gamma_0(1)$ the inequality (3.38) is satisfied for $\theta = 0$, and then in the bound for $m \neq 0$ in Proposition 3.1 we can substitute $y^{-\theta}$ by $|m|^\theta$ (see [3]). For π continuous we can substitute $y^{-\theta}$ by $\min(1, |t_\pi|)$ (see [3, 6]).

We can improve slightly Proposition 3.1 in average as follows. First, by Parseval we get

$$(3.42) \quad \sum_m |\tilde{f}(m, y)|^2 = \int_0^1 |f(h(x)a(y))|^2 dx \leq \|f\|_{L^\infty}^2$$

where $\tilde{f}(m, y) = \int_0^1 f(h(t)a(y))e(-mt) dt$. Moreover, integrating by parts b times and using the identity $h(\epsilon)a(y) = a(y)h(\epsilon/y)$ we get the bound $|\tilde{f}(m, y)| \ll_b (|m|y)^{-b} |\widetilde{D_R^b f}(m, y)|$, so by Cauchy's inequality and Parseval we have

$$(3.43) \quad \sum_{|m|>M} |\tilde{f}(m, y)|^2 \ll_b (My)^{-b} \|D_R^b f\|_{L^\infty}^2.$$

To end this section we record some bounds for sums over Hecke eigenvalues $\lambda_\pi(n)$, that will be needed later. We restrict to $\Gamma = \Gamma_0(1)$ the full modular group as this is what will be used. Since such a π is either a fixed holomorphic form or a fixed Maass form it is known that (for the holomorphic case it is classical while for the Maass case see [15]): for $T \geq 1$ and $x \in \mathbb{R}$

$$(3.44) \quad \sum_{m \leq T} \lambda_\pi(m)e(mx) \ll_{\pi,\epsilon} T^{\frac{1}{2}+\epsilon}.$$

We need control on the dependence in (3.44) of the implied constant. This can be done in terms of the eigenvalue λ_π of π . It is also convenient

for us here and elsewhere to work with the smooth normalized sums;

$$(3.45) \quad \sum_m \frac{\lambda(m)}{\sqrt{m}} e(mx) \psi(mu)$$

for $u > 0$ and small.

We can use the set up above with Whittaker functions to bound the sums in (3.45). For $f \in V_\pi$ one can control the L^∞ -norm of f by its Sobolev norms (see [3]):

$$(3.46) \quad \|f\|_{L^\infty} \ll \tilde{\lambda}_\pi^{3+o(1)} \|f\|_{W^3}$$

for π discrete, with $\tilde{\lambda}_\pi = \max(1, |\lambda_\pi|)$, and

$$(3.47) \quad \sup_{g=h(x)a(y)k(\theta) \in G} \frac{|f(g)|}{y^{1/2} + y^{-1/2}} \ll \tilde{\lambda}_\pi^{3+o(1)} \|f\|_{W^3}$$

for π continuous.

Let π be discrete. Suppose that W_f has support in $(1, 2)$ for some $f \in V_\pi$. From (3.34) and the action of the differential operators on W_f we deduce

$$(3.48) \quad \|f\|_{W^b} \ll_b \tilde{\lambda}_\pi^{b+o(1)} \|W_f\|_{W^{2b}},$$

with the norms for W_f being the usual Sobolev norms for real functions. On the other hand, for any $\psi \in C_0^\infty((1, 2))$, there exists $f \in V_\pi$ with $W_f = \psi$. This comes from the fact that if $f \in V_\pi$ then $f_g(x) = f(xg) \in V_\pi$ for any $g \in G$, and $W_{f_{a(y)}}(u) = W_f(yu)$, $W_{f_{h(x)}}(u) = e(x)W_f(u)$. Then, due to (3.46) and (3.48), for any $0 < u < 1$ we have

$$(3.49) \quad \sum_m \frac{\lambda_\pi(m)}{\sqrt{m}} e(mx) \psi(mu) = f(h(x)a(u)) \ll \tilde{\lambda}_\pi^{6+o(1)} \|\psi\|_{W^6}.$$

In the same way we have ²

$$(3.50) \quad \left(\sum_m \left| \frac{\lambda_\pi(m)}{\sqrt{m}} \right|^2 |\psi(mu)|^2 \right)^{\frac{1}{2}} = \left(\int_0^1 |f(h(t)a(u))|^2 dt \right)^{\frac{1}{2}} \ll \tilde{\lambda}_\pi^{3+o(1)} \|\psi\|_{W^6}.$$

For π continuous, the suitable norm for the Hilbert space V_π satisfies $\min(1, |t_\pi|) \|f\|_{V_\pi} \asymp \|W_f\|_{L^2}$ and again from the action of the differential operators on W_f we have

$$(3.51) \quad \min(1, |t_\pi|) \|f\|_{W^b} \ll_b \tilde{\lambda}_\pi^{b+o(1)} \|W_f\|_{W^{2b}}.$$

²Note that in terms of the dependence in λ_π this is much weaker than Iwaniec's [20], $\sum_{m \leq T} |\lambda_\pi(m)|^2 \ll_\epsilon T \tilde{\lambda}_\pi^\epsilon$ and its extensions (7.8) and (7.10)

Proceeding as in the discrete case, from (3.47) and (3.51) and sending $h(x)a(u)$ to the fundamental domain gives

$$(3.52) \quad \min(1, |t_\pi|) \sum_m \frac{\lambda_\pi(m)}{\sqrt{m}} e(mx) \psi(mu) \ll \min\left(\frac{u^{-1/2}}{q}, \frac{u^{1/2}}{\|qx\|}\right) \tilde{\lambda}_\pi^{6+o(1)} \|\psi\|_{W^6},$$

where q is the natural number smaller than $u^{-1/2}$ for which the quantity in the formula is the largest. This last formula can be improved (by using Poisson Summation in \mathbb{R}^2 , or by Voronoi formula) whenever $1/\|qx\| > 1/uq$, in the sense of adding the decay factor

$$\left(\frac{1/\|qx\|}{1/uq}\right)^{-b}$$

for any fixed $b > 0$ to the right of the inequality—but paying with a factor $O_b(\tilde{\lambda}_\pi^{O(b)} \|\psi\|_{W^{O(b)}})$.

4. EFFECTIVE EQUIDISTRIBUTION: CONTINUOUS ALGEBRAIC MEASURES

In this section we prove effective versions of Dani Theorem. We begin with some useful notation. First we clarify the notion of continuous algebraic measure.

Definition 4.1 (Continuous algebraic measure). We shall say that an algebraic measure on X is *continuous* if it arises as limit when T goes to infinity of the probability measures carried by the pieces of horocycle $\Gamma gh(st)$, $t \in [0, T]$, for some fixed $g \in G$, $s \in \mathbb{R}$.

Since we are going to talk about probability measures, it is a good idea to normalize sums.

Definition 4.2 (Expectation). Let $J \subset \mathbb{R}$ be a finite set. Let $f : \mathbb{R} \rightarrow \mathbb{C}$. We define the *expectation* of f in J , and we write it as

$$\mathbb{E}_{x \in J} f(x),$$

as $\sum_{x \in J} f(x) / \sum_{x \in J} 1$. If J is a bounded subinterval of \mathbb{R} , we define the expectation in the same way but using integrals.

To deal with bounded functions that vary slowly, especially near the cusp, we use the following Lipschitz norms.

Definition 4.3 (Lipschitz norm). Let $f : X \rightarrow \mathbb{C}$. We define the *Lipschitz norm* of f as

$$\|f\|_{\text{Lip}} = \|f\|_{L^\infty} + \sup_{x \neq x' \in X} \frac{|f(x) - f(x')|}{\widehat{d}_X(x, x')}.$$

with \widehat{d}_X the metric on X defined as

$$\widehat{d}_X(x, x') = \min(d_X(x, x'), e^{-d_X(x, \Gamma I)} + e^{-d_X(x', \Gamma I)}),$$

where I is the identity matrix in G .

Notice that the use of the metric \widehat{d}_X instead of d_X in the previous definition implies that f can be seen as a function in the one-point compactification of X .

Finally, we define the quantitative concept of equidistribution that we shall use to describe our main results

Definition 4.4 (Effective equidistribution). Let $0 < \delta < 1/2$. Let $g : J \rightarrow X$, with J either a subinterval of \mathbb{Z} or of \mathbb{R} . We say that $g(x), x \in J$ is δ -equidistributed w.r.t. a Borel probability measure μ on X if

$$|\mathbb{E}_{x \in J} f(g(x)) - \int f d\mu| \leq \delta \|f\|_{\text{Lip}}$$

for any $f : X \rightarrow \mathbb{C}$.

Remark. This concept of quantitative equidistribution is much weaker than the one used by Green and Tao in [13].

The Lipschitz norm controls the Sobolev norms in the following sense: let $\delta > 0$; for any $f \in C(X)$ with $\|f\|_{\text{Lip}} = 1$ there exists another function $f_\delta \in C(X)$ such that

$$(4.1) \quad \|f - f_\delta\|_{L^\infty} < \delta$$

and $\|Df_\delta\|_{L^\infty} \ll_{\text{ord}D} \delta^{-\text{ord}D}$ for any left-invariant differential operator. This will allow us from now on to substitute a function with bounded Lipschitz norm by one with “bounded” Sobolev norms³. We can build f_δ as follows: for two functions $f, w \in L^2(SL_2(\mathbb{R}))$ define its convolution as $f * w(g) = \int_{SL_2(\mathbb{R})} f(gt)w(t^{-1}) d\mu(t)$; if f is Γ -invariant then so is $f * w$. Moreover we can also write $f * w(g) = \int_{SL_2(\mathbb{R})} f(t)w(t^{-1}g) d\mu(t)$ by the left invariance of $d\mu$. Pick $\psi \in C_0^\infty((0, 1))$ a non-negative function and define $\psi_\delta(g) = c_\delta \psi(\delta^{-1}d_{SL_2(\mathbb{R})}(g, I))$ with c_δ the constant that gives $\int_{SL_2(\mathbb{R})} \psi_\delta(g^{-1}) d\mu(g) = 1$. Finally, defining $f_\delta = f * \psi_\delta$ one can check that it satisfies the desired properties.

Before beginning with our quantitative results, let us recall the following lemma on cancellation in oscillatory integrals.

³Note that in Sections 4, 5 and 6 the entire discussion takes place on X ; that is $\Gamma = SL_2(\mathbb{Z})$. It is only in Section 7 that the level q of congruence subgroups is relevant.

Lemma 4.5 (Integration by parts). *Let $A > 1$. Let $\eta \in C_0^\infty(1, 2)$, $F \in C^\infty(\mathbb{R})$, η with bounded derivatives and $|F'| \asymp A$, $F^{(j)} \ll_j A$ for any $j \in \mathbb{N}$. Then*

$$\int \eta(t)e(F(t)) dt \ll A^{-\frac{1}{|\alpha(1)|}}.$$

Proof. Write $F(t) = Af(t)$ and integrate by parts several times. Since the derivatives of f are bounded, the result follows. \square

We are ready to give our first quantitative version of Dani's Theorem for a continuous orbit.

Theorem 4.6. *Let $0 < \delta < 1/2$, $T > 1$. For any $\xi \in X$ there exists a positive integer $j < \delta^{-O(1)}$ such $\xi h(t)$, t in any subinterval of $[0, T]$ of length T/j , is δ -equidistributed w.r.t. a continuous algebraic measure in X . Moreover, it is the volume measure unless $y_T^{-1} < \delta^{-O(1)}$.*

Remark. Notice that the use of the metric \widehat{d}_X in the definition of δ -equidistribution implies that Theorem 4.6 says nothing about the part of the horocycle that is near the cusp; the same can be said about any other result in sections 4 and 5.

Proof. We can assume that δ^{-1} is not larger than a small power of T , because otherwise we could choose $j = T/\delta$ and then for any $t_0 \in \mathbb{R}$, $\xi h(t)$, $t \in [t_0, t_0 + \delta]$ would be δ -equidistributed w.r.t. the algebraic measure carried by the point $\xi h(t_0)$.

If $y_T^{-1} > \delta^{-O(1)}$ we will prove the result with $j = 1$; then, in that case we can assume that $|1 \pm t/W_T| \gg \delta$ since this only deletes at most a proportion δ of the orbit. On the other hand, if for all t in a certain interval J we have $y_T/|1 \pm t/W_T|^2 \gg \delta^{-1}$ then by (2.5) we deduce that $\widehat{d}_X(\xi h(t), \xi h(t_0)) \ll \delta$ for any fixed $t_0 \in J$, and then $\xi h(t)$, $t \in J$ will be δ -equidistributed w.r.t. the algebraic measure carried by $\xi h(t_0)$; this implies that if $y_T^{-1} \leq \delta^{-O(1)}$ we can assume that $1/|1 \pm t/W_T| \ll \delta^{-O(1)}$. So in any case we can assume we are in the range $|1 \pm t/W_T|^{-1} < \delta^{-O(1)}$.

Now, if the Theorem is false we can assume that the Lipschitz norm of f is 1 and that it has $\delta^{-O(1)}$ -bounded Sobolev norms. Because of (2.5), there exists $y_T \leq y^* < y_T \delta^{-O(1)}$ and $\eta \in C_0^\infty(\mathbb{R})$ with $\|\eta\|_{W^j} < \delta^{-O(j)}$ and support an interval I outside of $1/|1 \pm t/W_T| < \delta^{-O(1)}$ such that

$$|\mathbb{E}_{t \in I} \eta\left(\frac{t}{T}\right) f_0\left(h\left(\alpha_T + \frac{y_T W_T}{1 \pm t/W_T}\right) a(y^*)\right)| \gg \delta,$$

where f_0 equals $f - \int_0^1 f(h(x)a(y^*)) dx$ in the case $y_T^{-1} \ll \delta^{-O(1)}$ and $f - \int_X f d\mu_G$ otherwise. Notice that this implies, by Proposition 3.1 in

the second case, that $\int_0^1 f_0(h(t)a(y^*)) dt \ll \delta^2$. But then, expanding in Fourier series $f_0(h(x)a(y)) = \sum_m \tilde{f}_0(m, y)e(mx)$ we get

$$\sum_{m \neq 0} |\tilde{f}_0(m, y^*)| \left| \int \eta(t) e\left(\frac{my_T W_T}{1 \pm tTW_T^{-1}}\right) dt \right| \gg \delta.$$

The derivative of the phase in the exponential is of size $|m|y_T T \gg |m|$, and then integrating by parts (Lemma 4.5) and using Proposition 3.1 we have for any large $b > 0$ that

$$\delta^{-O(1)} \sum_{m=1}^{\infty} y_T^{1/2} \tau(m) m^\theta (my_T T)^{-b} \gg_b \delta.$$

This implies that $\delta^{-O(1)} T^{-1/2} \gg \delta^{-O(1)} y_T^{1/2} (y_T T)^{-b} \gg_b \delta$, which is a contradiction. \square

Remark: This result is related to Theorem 1 in [38] for the particular case $\Gamma = SL(2, \mathbb{Z})$. It could also be proven by using quantitative mixing, as done in [39] for $\Gamma \backslash G$ compact but taking care of the position of the orbit.

When we consider Γg fixed, we can make the result more explicit.

Theorem 4.7 (Effective Dani, continuous orbit). *Let $\xi \in X$ fixed. The continuous algebraic measure appearing in Theorem 4.6 is the volume measure unless there exists an integer $q < \delta^{-O(1)}$ such that*

$$\|q\alpha\| < \delta^{-O(1)} T^{-1}.$$

Proof. This follows from Theorem 4.6, Lemma 2.8 and the remark before it. \square

Remark. Here and hereafter whenever we speak of “algebraic measure” it will mean one of the three algebraic measures appearing in Dani’s Theorem (described in section 2).

As a corollary of this result we obtain Dani’s Theorem for the continuous flow: if α is irrational, for any $0 < \delta < 1/2$ we can find $T = T_\delta$ large enough so that no q exists as in the statement, and then $\Gamma gh(t)$, $t \in [0, T]$ is δ -equidistributed w.r.t. the volume measure.

Next we begin our study of the discrete orbit $\xi h(s)^n$, $n = 0, 1, \dots, N$. Our aim in the rest of this section is to understand, in terms of ξ when any large piece of this orbit is δ -equidistributed w.r.t. a continuous algebraic measure. We will be able to do it for not very large s ; essentially in the range $s < (sN)^{1/5}$.

In the proofs we shall use the following version of the Stationary Phase principle (a more precise one can be found in [18]):

Lemma 4.8 (Stationary Phase). *Let $\eta \in C_0^\infty((1, 2))$. Let $A > 1$, $F \in C^\infty$ with $F'' \asymp A$ and $F^{(j)} \ll_j A$ for $j \geq 3$ in the interval $(1/2, 4)$. Then, we have*

$$\int \eta(t)e(F(t)) dt = \eta^*(t_F)e(F(t_F))A^{-1/2} + O(A^{-\frac{1}{|\alpha(1)|}})$$

where t_F is the only point where F' vanishes, and $\eta^* \in C_0^\infty((1/2, 4))$ depends just on A and F'' , and η^* and its derivatives are bounded (in terms of η only).

Proof. First, let us write

$$\int \eta(t)e(F(t)) dt = e(F(t_F)) \int \eta(t_F + s)e(J(t_F, s)) ds$$

where

$$J(t, s) = F(t + s) - F(t) - F'(t)s = \int_0^s \int_0^u F''(t + v) dv du.$$

Let us choose $C_1 = A^{1/8}$. For some $\psi \in C_0^\infty((-2, -1) \cup (1, 2))$ we can write

$$1(s) = \sum_{j=-\infty}^{\infty} \psi\left(\frac{s}{2^j A^{-1/2}}\right) = \psi_0\left(\frac{s}{A^{-1/2}}\right) + \psi_1\left(\frac{s}{C_1 A^{-1/2}}\right) + \sum_{1 < 2^j < C_1} \psi\left(\frac{s}{2^j A^{-1/2}}\right),$$

with ψ_0 and ψ_1 coming from the sum over $2^j \leq 1$ and $2^j \geq C_1$ respectively. Decomposing the integral accordingly, one can easily show that the part corresponding to ψ_1 is bounded by $C_1^{-1/|\alpha(1)|} < A^{-1/|\alpha(1)|}$, and the rest can be expressed as

$$A^{-1/2}e(F(t_F))\eta^*(t_F)$$

with

$$\eta^*(t) = \eta_0^*(t) + \sum_{1 < 2^j < C_1} \eta_{2^j}(t),$$

where

$$\eta_C(t) = C \int \psi_C(t, s)e(C^2 J_C(t, s)) ds.$$

with $\psi_C(t, s) = \psi(s)\eta(t + CA^{-1/2}s)$ and $J_C(t, s) = C^{-2}J(t, CA^{-1/2}s)$, and η_0^* is defined as η_{2^0} but changing ψ by ψ_0 . The first thing to note is that the support of η^* is contained in $(1/2, 4)$ and the second is that η^* just depends on F'' , A and η . It is easy to show that η^* and its derivatives are bounded. For the rest, we have that the support of $\psi_C(t, \cdot)$ is contained in $(1, 2)$; moreover $\frac{\partial}{\partial s} J_C \asymp 1$, and the partial derivatives of both ψ_C and J_C in t and s are uniformly bounded. Then,

by integrating by parts several times, one shows that $(\frac{\partial}{\partial t})^j \eta_C(t) \ll_j C^{-1}$ and then $(\frac{\partial}{\partial t})^j \eta^*(t) \ll_j 1$ \square

In order to prove a quantitative version of Dani's Theorem for the discrete orbit $\Gamma gh(s)^n$, $n = 0, 1, \dots, N$, we split the analysis in three cases: when the piece of orbit is near to a closed horocycle (θ_T and y_T^{-1} small, $T = sN$), when it is far from any closed horocycle (θ_T and y_T^{-1} large), and the intermediate case.

We begin by the "near" case. Here, one can parametrize the orbit essentially as $\Gamma h(p(n))a(y_T)$ with $p(n)$ a quadratic polynomial. The following allows us to handle the distribution for such a sequence.

Proposition 4.9 (Near to a closed horocycle). *Let $0 < \delta < 1/2$. The sequence $s_n = \alpha + \beta n + \omega n^2 + iy$, $n \leq N$ is δ -equidistributed with respect to the algebraic measure on the closed horocycle of period y^{-1} unless either $\min(y^2/N|\omega|, Ny) < \delta^{-O(1)}$ or there exists a natural number q smaller than $(\delta/\tau(q_2))^{-O(1)}(1 + y^{-1/2})$ such that*

$$(4.2) \quad \|q\alpha\| + N\|q\beta\| + N^2q|\omega| < y^{\frac{1}{2}}(\delta/\tau(q_2))^{-O(1)},$$

where q_2 is the denominator in the expression as reduced fraction of $[q\beta]/q$.

Proof. The result is trivial for $y > \delta^{-1}$, since then for any $x_1, x_2 \in \mathbb{R}$ we have $|f(x_2 + iy) - f(x_1 + iy)| \leq \delta \|f\|_{\text{Lip}}$. Otherwise, let us suppose that $s_n, n \leq N$ does not satisfy the condition in the statement and that both $y^2/N|\omega|$ and Ny are larger than δ^{-c} for $c > 1$ a large constant. Then

$$|\mathbb{E}_{n \leq N} f(\alpha + \beta n + \omega n^2 + iy) \eta(n/N)| \gg \delta,$$

for some f with $\int_0^1 f(x + iy) dx = 0$, $\|f\|_{W^j} \ll_j \delta^{-j}$ and $\eta \in C_0^\infty((0, 1))$ with $\|\eta^{(j)}\|_{L^\infty} \ll_j \delta^{-j}$. Take $M = y^{-1} \delta^{-\sqrt{c}}$; there exist $q_2 \leq M$ and a_2 coprime to q_2 such that

$$\beta = \frac{a_2}{q_2} + \epsilon \quad |\epsilon| \leq \frac{1}{q_2 M}.$$

By splitting into arithmetic progressions modulo q_2 we have

$$(4.3) \quad \delta^{-O(1)} |\mathbb{E}_{b \leq q_2} \mathbb{E}_{j \leq N/q_2} f(\alpha + \frac{a_2 b}{q_2} + \epsilon q_2 j + \omega q_2^2 j^2 + iy) \eta(\frac{j}{N/q_2})| > 1.$$

Since $\|F^{(j)}\|_{L^\infty} \ll_j \delta^j$, with $F(t) = f(\alpha + a_2 b/q_2 + \epsilon q_2 t + \omega q_2^2 t^2 + iy) \eta(\frac{t}{N/q_2})$, we can see (for instance applying Poisson Summation and integrating by parts) that it is possible to substitute the inner sum by

the integral

$$(4.4) \quad \delta^{-O(1)} |\mathbb{E}_{b \leq q_2} \int f(\alpha + \frac{a_2 b}{q_2} + \epsilon N t + \omega N^2 t^2 + iy) \eta(t) dt| > 1.$$

At this point, let us remark that from the beginning we could assume the support of η to be at distance δ from the point $-\epsilon(2\omega N)^{-1}$. By the spectral decomposition (3.29) and the Fourier series expansion (taking into account that the zero coefficient vanishes) we have

$$\delta^{-O(1)} \left| \sum_{m \neq 0} \frac{\lambda_\pi(m)}{\sqrt{|m|}} e(m\alpha) W_f(|m|y) \mathbb{E}_{b \leq q_2} e(\frac{ma_2 b}{q_2}) I(m) \right| > 1$$

for some π with $\lambda_\pi \ll \delta^{-O(1)}$ and $f \in V_\pi$ with $\|f\|_{W_j} \ll_j \delta^{-O(j)}$ and $I(m) = \int \eta(t) e(mN(\epsilon t + \omega N t^2)) dt$. Then

$$\delta^{-O(1)} \left| \sum_{k \neq 0} \frac{\lambda_\pi(q_2 k)}{\sqrt{|q_2 k|}} e(kq_2 \alpha) W_f(|kq_2|y) I(kq_2) \right| > 1.$$

By the multiplicativity of the Hecke eigenvalues

$$(4.5) \quad \lambda_\pi(ab) = \sum_{d|(a,b)} \mu(d) \lambda_\pi(a/d) \lambda_\pi(b/d)$$

we have

$$\delta^{-O(1)} \frac{|\lambda_\pi(q_2/d)|}{\sqrt{dq_2}} \left| \sum_{m \neq 0} \frac{\lambda_\pi(m)}{\sqrt{|m|}} e(mdq_2 \alpha) W_f(|mdq_2|y) I(mdq_2) \right| > \frac{1}{\tau(q_2)}$$

for some $d | q_2$. By splitting the integral defining $I(x)$ smoothly into δ^{-1} integrals and applying Lemma 4.5, taking into account our previous remark on the support of η , we have

$$\delta^{-O(1)} \frac{|\lambda_\pi(q_2/d)|}{\sqrt{dq_2}} \left| \sum_{m \neq 0} \frac{\lambda_\pi(m)}{\sqrt{|m|}} e(mdq_2 \alpha) W_f(|mdq_2|y) \psi(mdq_2 u) \right| > \frac{1}{\tau(q_2)}$$

with $u = \max(|\epsilon N|, |\omega| N^2)$ and ψ a smooth function with

$$\psi^{(j)}(x) \ll_{j,k} \delta^{-j} (1 + |x|)^{-k} \quad k, j \geq 0.$$

The function ψ could depend on some of the parameters, but the bounds on the derivatives are absolute. Thus, if π is discrete, splitting into dyadic intervals and using (3.49), Proposition 3.1 (with its remark) and (3.37) we have

$$\delta^{-O(1)} (q_2/d)^\theta (dq_2)^{-1/2} \min(1, yu^{-1})^{1/2} > \tau(q_2)^{-1}$$

so that $q_2 < \delta^{-O(1)}$ and $u < y\delta^{-O(1)}$. Now, there exists a natural number $q_1 \leq 1/\sqrt{\delta y}$ such that $\|q_1 \alpha\| \ll \sqrt{\delta y}$; picking $q = q_1 q_2$ one can

check that the condition (4.2) is satisfied. If π is continuous, again by Proposition 3.1 and (3.52)—with the remarks after them—we have

$$\delta^{-O(1)}(dq_2)^{-1/2}\sqrt{dq_2y} \min\left(\frac{1}{q^*dq_2(y+u)}, \frac{1}{\|q^*dq_2\alpha\|}\right) > \tau(q_2)^{-2}$$

for some $q^* \ll (dq_2(y+u))^{-1/2}$. Therefore for some $d^* \mid q_2$

$$q^*d^*q_2 < (\delta/\tau(q_2))^{-O(1)}y^{-\frac{1}{2}}(1+uy^{-1})^{-1}$$

and $\|q^*d^*q_2\alpha\| < (\delta/\tau(q_2))^{-O(1)}y^{1/2}$. These conditions imply (4.2) with $q = q^*d^*q_2$. \square

We now deal with the intermediate case. In the proof we simply use bounds for the Fourier coefficients.

Proposition 4.10 (Intermediate case). *Let $0 < \delta < 1/2$, $N \geq 1$, $T = sN \geq 1$ and $g \in G$. Assume $s\delta^{-O(1)} < T^{1/4}$. The sequence $\Gamma gh(s)^n$, $n \leq N$ is δ -equidistributed w.r.t. the probability measure carried by $\Gamma gh(t)$, $t \in [0, T]$ unless either both $y_T < \delta^{-O(1)}s/T$ and $|W_T|y_T < \delta^{-O(1)}s^2T^{2\theta+o(1)}$ or both $y_T\delta^{-O(1)} > s^{-2}T^{-2\theta-o(1)}$ and $|W_T|y_T\delta^{-O(1)} > (T/s)^{2-o(1)}y_T^{2\theta}$.*

Proof. Let us assume for simplicity that $W_T > 0$. If $\Gamma gh(s)^n$, $n \leq N$ is not δ -equidistributed w.r.t. the measure carried by $\Gamma gh(t)$, $t \in [0, T]$, splitting the orbit (2.5) and using (4.1) we have that

$$|\mathbb{E}_{n/N \in I} f(h(x_{sn})a(y_*))\eta(n/N)| \gg \delta,$$

for some $y_T < y_* < y_T\delta^{-O(1)}$, f bounded with $\int_0^1 f(t+iy_*)dt = 0$ and derivatives bounded by $\delta^{-O(1)}$, $\eta \in C_0^\infty$ with derivatives bounded by $\delta^{-O(1)}$ and supported on an interval I of length larger than δ^2 , at distance at least δ^2 from zero and such that $\delta^{-O(1)}|1 \pm sn/W| > 1$, with

$$x_t = \alpha + yW(1 \pm t/W)^{-1},$$

$W = W_T$, $y = y_T$. Now, using the Fourier expansion of $f(t+iy_*)$, (3.43) and Proposition 3.1 we have

$$(4.6) \quad \delta^{-O(1)} \sum_{\delta^{-O(1)}y^{-\epsilon} < |m| < \delta^{-O(1)}y^{-1}} |\tilde{f}(m, y_*)| |\mathbb{E}_n \eta(\frac{n}{N}) e(mx_{sn})| > 1.$$

Then, we have (4.6) with the restriction that $|m|$ is either smaller or larger than $\delta^{-O(1)}T^\epsilon W/T$. In the former possibility, by Poisson summation in the inner sum, Lemmas 4.5 and 4.8 and Proposition 3.1 we deduce that $y\delta^{-O(1)} > s^{-2}T^{-2\theta-o(1)}$ and $Wy\delta^{-O(1)} > (T/s)^{2-o(1)}y^{2\theta}$. In the latter, repeating the same steps we get that Wy is at most

$\delta^{-O(1)}s^2T^{2\theta+o(1)}$, and using Cauchy's inequality and (3.42) instead of Proposition 3.1 we have

$$\delta^{-O(1)} \sum_{W/T < m < \delta^{-O(1)}y^{-1}} \left| \sum_{j \neq 0} b_j \eta_j^* \left(\frac{1 - \sqrt{\frac{my}{j/s}}}{T/W} \right) \frac{e(2W\sqrt{myj/s})}{\sqrt{(j/s)T^2/W}} \right|^2 > 1$$

with $b_j = e(-Wj/s)$. Now, introducing a smooth function in the outer sum, expanding the square and changing the order of summation gives that either (from the diagonal terms) $\delta^{-O(1)}sT/yW > T^2W^{-1}$ or

$$\delta^{-O(1)} \left| \sum_m \rho\left(\frac{my/c - 1}{a}\right) \psi(my/c) e(\epsilon c T \sqrt{my/c}) \right| > Ts^{-2}(\epsilon c)^{-1+o(1)}$$

for some $\psi, \rho \in C_0^\infty$ with support and derivatives bounded by $\delta^{-O(1)}$ and at distance δ^2 from zero, $a < \delta^{-O(1)}$ and $\epsilon, c < \delta^{-O(1)}$, $\epsilon c \gg 1/s$.

The first possibility implies that $y < \delta^{-O(1)}s/T$. The second possibility implies, writing $\rho(x) = \int \hat{\rho}(\theta) e(x\theta) d\theta$, that for some θ we have

$$\delta^{-O(1)} \left| \sum_m \psi(my/c) e(m\theta) e(\epsilon c T \sqrt{my/c}) \right| > Ts^{-2}(\epsilon c)^{-1+o(1)}$$

which applying Poisson summation and Lemmas 4.5 and 4.8 gives that $\delta^{-O(1)}s > T^{1/4}$, a contradiction. \square

Now we treat the case in which the piece of horocycle is far from a closed one. The previous result worked for $y_T > s/T$, that is almost for the whole range. To handle the remaining range, we are going to split the piece of horocycle of length T in pieces of length ϵT , $\epsilon < 1$. Then we shall put those pieces in the fundamental domain $D_{X, \epsilon T}$, and then use the Fourier expansion (as in Proposition 4.10). The advantage is that now we can get an extra cancellation from the fact that different pieces are essentially uncorrelated, and we can control this independence arithmetically. The ideas of the method come from the theory of exponential sums; it is essentially a modification of the one in [23].

Proposition 4.11 (Far from a closed horocycle). *Let $0 < \delta < 1/2$, $N \geq 1$, $T = sN \geq 1$ and $g \in G$. Assume $y_T < \delta^{-O(1)}s/T$ and $W_T y_T < \delta^{-O(1)}s^2T^{2\theta+o(1)}$. The sequence $\Gamma gh(s)^n$, $n \leq N$ is δ -equidistributed w.r.t. the volume measure on X unless $T^{1/5}s^{-1}$ is smaller than $\delta^{-O(1)}$.*

Proof. Let us assume for simplicity that $W_T > 0$, and write $y = y_T$, $W = W_T$. If $\Gamma gh(s)^n$, $n \leq N$ is not δ -equidistributed w.r.t. the volume measure on X , then by (2.5) and (4.1) there exists $f \in L_0^2(X)$ bounded and with derivatives bounded by $\delta^{-O(1)}$ and an interval I of length

satisfying $N|I|^{-1} < \delta^{-O(1)}$ such that

$$|\mathbb{E}_{n \in I} f(h(x_{sn})a(y_*))| \gg \delta$$

where $y \ll y_* < \delta^{-O(1)}y$, $\delta^{-O(1)}|1 \pm sn/W| > 1$ for any $n \in I$ and

$$x_t = \alpha + yW(1 \pm t/W)^{-1}.$$

We can take α such that $|\alpha + yW| \leq 1/2$. Now we divide the sum into sums such that the argument in $h(\cdot)$ is near to a Farey fraction a/q up to $q \leq K$, with $K^2y = (sT)^{-2/3}$. This choice for K will make sense later; by now we can say that since $y < \delta^{-O(1)}s/T$, we have $K > \delta^{-O(1)}$ in the range $\delta^{-O(1)}s < T^{1/5}$. Precisely, we are going to take the interval around a/q

$$\left(\frac{a}{q} - \frac{1}{q(q+q')}, \frac{a}{q} + \frac{1}{q(q+q'')}\right)$$

where q' and q'' are the denominators of the Farey fractions (up to K) to the left and right of a/q respectively. The point is that both $q+q'$ and $q+q''$ are comparable to K . In this way, we can write

$$1(x) = \sum_{a,q} \tilde{\eta}_{a,q} \left(\frac{x - a/q}{1/qK}\right)$$

where a, q ranges over integers a, q with $q \leq K$, and $\tilde{\eta}_{a,q}$ are C_0^∞ functions with $\|\tilde{\eta}_{a,q}^{(j)}\|_{L^\infty} \ll_j \delta^{-O(j)}$ with uniformly bounded support, and $\tilde{\eta}_{a,q} = 0$ if a, q are coprimes. In this way, we have

$$|\mathbb{E}_{q \leq K} \mathbb{E}_{a \in qL} \mathbb{E}_{n \in I_{a/q}} f(h(x_{sn})a(y_*)) \tilde{\eta}_{a,q}(qK(x_{sn} - a/q))| \gg \delta$$

with L a subinterval of $[-\delta^{-O(1)}yT, \delta^{-O(1)}yT]$ with $yT/|L| < \delta^{-O(1)}$, and $I_{a/q}$ an interval of size $\asymp 1/syqK \gg \delta^{-O(1)}$ containing all n such that $qK(x_{sn} - a/q)$ is in the support of $\tilde{\eta}_{a,q}$. Taking into account that the fractions with $q < \delta^2K$ give an amount $O(\delta^2)$, we have

$$|\mathbb{E}_{\delta^2K < q < K} \mathbb{E}_{a \in qL} \mathbb{E}_{n \in I_{a/q}} f(h(x_{sn})a(y_*)) \eta_{a,q}^\#(q^2(x_{sn} - a/q))| \gg \delta$$

with $\eta_{a,q}^\#(t) = \tilde{\eta}_{a,q}(tK/q)$. Now, splitting $\delta^2K < q < K$ into $\delta^{-O(1)}$ intervals of equal length, we deduce that

$$\mathbb{E}_{|q-q_0| < \delta^c K} \mathbb{E}_{a \in qL} \mathbb{E}_{n \in I_{a/q}} f(h(x_{sn})a(y_*)) \eta_{a,q}^\#(q^2(x_{sn} - a/q))| \gg \delta$$

for some $q_0 \in [\delta^2K, K]$, $c > 1$ a large constant. This implies that we can change $a \in qL$ by $a \in q_0L$. Proceeding in the same way we have

$$|\mathbb{E}_{|q-q_0| < \delta^c K} \mathbb{E}_{|a-a_0| < \delta^c K} \mathbb{E}_{n \in I_{a/q}} f(h(x_{sn})a(y_*)) \eta_{a,q}^\#(q^2(x_{sn} - a/q))| \gg \delta$$

with $\text{supp } \eta_{a,q} \subset [w_0, w_0 + \delta^c]$ for some fixed $\delta^2 \ll w_0 \ll 1$ and $a_0 \asymp \delta^{-O(1)} yTK$, $\eta_{a,q}$ with the same properties as $\tilde{\eta}_{a,q}$. Defining $v_t = q^2(x_t - a/q)$ we can write

$$|\mathbb{E}_q \mathbb{E}_a \mathbb{E}_n f(h(a/q + v_{sn} q^{-2})a(y_*))\eta_{a,q}(v_{sn})| \gg \delta$$

with a, q, n moving in the previous ranges. But multiplying by $\gamma_{a/q} \in \Gamma$ as in the proof of Lemma 2.10 we have that $\gamma_{a/q} h(a/q + v/q^2)a(y_*)$ equals

$$h\left(-\frac{\bar{a}}{q} - \frac{v}{v^2 + (q^2 y_*)^2}\right) a\left(\frac{q^2 y_*}{v^2 + (q^2 y_*)^2}\right) k\left(-\text{arccot} \frac{v}{q^2 y_*}\right)$$

which, due to the restrictions on the ranges of $a, q, \eta_{a,q}$, is at distance $O(\delta^2)$ from $h(-\bar{a}/q - 1/v + r_0)a(K^2 y_0)$ for some y_0 with $y < y_0 < \delta^{-O(1)} y$ and r_0 , both independent of a, q . Then

$$|\mathbb{E}_q \mathbb{E}_a \mathbb{E}_n f(h(-\frac{\bar{a}}{q} - \frac{1}{v_{sn}} + r_0)a(K^2 y_0))\eta_{a,q}(v_{sn})| \gg \delta$$

The Fourier expansion $f(h(t)a(y)) = \sum_m \tilde{f}(m, y)e(mt)$ gives that

$$\delta^{-O(1)} \left| \sum_m \tilde{f}(m, K^2 y_0) e(mr_0) \mathbb{E}_{a,q} e\left(-\frac{\bar{a}m}{q}\right) \mathbb{E}_n e\left(-\frac{m}{v_{sn}}\right) \eta_{a,q}(v_{sn}) \right| > 1.$$

By Proposition 3.1 we can assume that m is in the range $|m| > \delta^{-O(1)} T^\epsilon$ for some $\epsilon > 0$. Thus, by Cauchy's inequality and (3.43) we have

$$\delta^{-O(1)} \mathbb{E}_m |\mathbb{E}_{a,q} e\left(-\frac{\bar{a}m}{q}\right) \mathbb{E}_n e\left(-\frac{m}{v_{sn}}\right) \eta_{a,q}(v_{sn})|^2 > K^2 y,$$

where m moves in the range $\delta^{-O(1)} T^\epsilon < m < \delta^{-O(1)} (K^2 y)^{-1}$. By Poisson Summation in the n -sum, Lemmas 4.5 and 4.8 and splitting of the obtained sum into intervals of the shape $[J, (1 + \delta)J]$ we have

$$\delta^{-O(1)} \mathbb{E}_m |\mathbb{E}_x b_x e\left(A \frac{m}{F} + B \sqrt{\frac{m}{F}} \eta_x^*\left(\frac{m}{F}\right)\right)|^2 > c^{-1-o(1)} s^{-2}$$

with $x = (a, q, j)$, b_x independent of m and bounded, j in the range $cs < j < (1 + \delta)cs$ for some $c < \delta^{-O(1)}$, $cs \geq 1$, and

$$F = \frac{c}{K^2 y}, \quad A = \left(-\frac{\bar{a}}{q} + \frac{q^{-2}}{a/q - \alpha}\right) F, \quad B = \frac{2Wy}{a/q - \alpha} \frac{K}{q} \sqrt{\frac{j}{cs}} F.$$

We also can assume that we are summing just in the a 's for which \bar{a}/q is contained in an interval of length $1/2$. Introduction of a smooth

factor in the m -sum and expansion of the square followed by Poisson Summation in m and Lemmas 4.5 and 4.8 gives that

$$(4.7) \quad \delta^{-O(1)} \mathbb{E}_{(x,x')} \left(1_{x=x'} + \frac{1_{\bar{a}/q \neq \bar{a}'/q'}}{\sqrt{|A-A'|}} + \frac{1_V}{\sqrt{1+|A-A'|}} \right) > \frac{c^{-o(1)}}{(cs)^2},$$

where $A' = A(x')$, $B' = B(x')$ and V is the subset of (x, x') with $x \neq x'$ for which $\bar{a}/q = \bar{a}'/q'$ and either $|A-A'|$ and $|B-B'|$ are of comparable size or both of them are bounded by $\delta^{-O(1)}T^\epsilon$. This is so due to the restrictions in the ranges of a, q, j and the fact that $a/q - \alpha \asymp Wy$.

The first summand in (4.7) gives a contribution of $\delta^{-O(1)}(KyTKcs)^{-1}$ to the expectation, which equals $\delta^{-O(1)}c^{-1}(sT)^{-1/3}$. This is smaller than the right hand of (4.7) in the range $\delta^{-O(1)}s < T^{1/5}$. In the second summand, again due to the restriction in the ranges of a, q, j , we have

$$|A - A'| \asymp |\bar{a}/q - \bar{a}'/q'|F.$$

Then, by counting we see that the contribution of the terms with $|\bar{a}/q - \bar{a}'/q'| \asymp U$ is $\delta^{-O(1)}U(UF)^{-1/2}$, so the full contribution of the second summand is $\delta^{-O(1)}F^{-1/2}$ which equals the one of the first summand—this is what motivated the election of K —and then is smaller than the right hand of (4.7).

Finally, we are going to show that the contribution of the third summand in (4.7) is even smaller than the other two. For the terms in V we have $q' = q$ and $a' = a + \lambda q$, with $0 \neq \lambda \ll yT$, and then

$$|A - A'| \asymp \frac{|\lambda|}{q^2(yW)^2}F, \quad |B - B'| \asymp \frac{K}{q} \left| \frac{\lambda}{a/q - \alpha} + 1 - \sqrt{j'/j} \right| F.$$

Suppose first that $|1 - \sqrt{j'/j}| \not\asymp |\lambda/(a/q - \alpha)|$. Then $|B - B'|$ is larger than $|\lambda/(a/q - \alpha)|$, so at least $\delta^{-O(1)}|A - A'|$. Moreover if $j' \neq j$ then $|B - B'|$ is at least F/j which is larger than $\delta^{-O(1)}T^\epsilon$ in the range $s\delta^{-O(1)} < T^{1/2-\epsilon}$, so $(x, x') \notin V$. Thus, we can assume $j' = j$, and then from (4.7) we deduce that for λ in some dyadic interval $\delta^{-O(1)}K^{-2}(|\lambda|/yT)(1/cs)$ is larger than $T^{-o(1)}(cs)^{-2}$, so $\delta^{-O(1)}|\lambda|F > T^{1-o(1)}/s$; using this bound we have

$$\delta^{-O(1)}|B - B'| \asymp \delta^{-O(1)}(Wy)^{-1}|\lambda|F \gg (Wy)^{-1}T^{1-o(1)}/s$$

which by $Wy < \delta^{-O(1)}s^2T^{2\theta+o(1)}$ is larger than $\delta^{-O(1)}T^\epsilon$ in the range $\delta^{-O(1)}s < T^{1/5}$, since $\theta < 1/5$. Therefore $(x, x') \notin V$, a contradiction.

Assume now that $|1 - \sqrt{j'/j}| \asymp |\lambda/(a/q - \alpha)|$. Then

$$(4.8) \quad 1 \leq |j' - j| \asymp j|\lambda|(yW)^{-1}$$

so the proportion of j' is $O(|\lambda|/Wy)$. We can write $a = a_1 + \mu q$ with $1 \leq a_1 \leq q$ and $\mu \ll yT$. For fixed λ, a_1, q, j' and j , it is easy to see

that $|A - A'|$ can be comparable to $|B - B'|$ for at most one μ ; in the same conditions, the number of μ for which both $|A - A'|$ and $|B - B'|$ can be at most $\delta^{-O(1)}T^\epsilon$ is $1 + \delta^{-O(1)}T^\epsilon(Wy)^2/(|\lambda|F)$. Applying the bound on Wy from (4.8) we have

$$(Wy)^2(|\lambda|F)^{-1} \ll |\lambda|j^2c^{-1}(sT)^{-2/3} \ll yTs^2(sT)^{-2/3},$$

which, due to the bound $y \ll \delta^{-O(1)}s/T$, is smaller than 1 in the range $s\delta^{-O(1)} < T^{2/7}$. Then, the proportion of μ in V is $\delta^{-O(1)}T^\epsilon/yT$, so the contribution from the third summand is the supremum in $\lambda \ll yT$ of

$$\delta^{-O(1)}T^\epsilon K^{-2} \frac{|\lambda|}{Wy} \frac{1}{yT} \left(\frac{|\lambda|}{K^2(yW)^2} F \right)^{-1/2}$$

which is at most $\delta^{-O(1)}T^\epsilon c^{-1/2}T^{-1/2}$ and then smaller than then right hand of (4.7) in the range $\delta^{-O(1)}s < T^{1/4-\epsilon}$. □

Remarks. (a) It would be possible to improve the range for s a little by treating non-trivially the exponential sums appearing in the proof in the last application of Stationary Phase. This would improve the range also in Theorem 4.12.

(b) We could prove this result also (with a smaller range for s) by proceeding as in the proof of Proposition 4.10, but treating the sums $\sum_n \lambda_\pi(n)e(F(n))$ with van der Corput's Lemma and shifted convolution. One could also prove it in a quicker way (again with s smaller) by relating the discrete orbit to the continuous one as in Lemma 3.1 of [39], and then using Theorem 4.6.

Now we join the previous cases to prove an effective version of Dani's result for sums

Theorem 4.12 (Effective Dani, discrete orbit). *Let $\xi \in X$ fixed. Let $0 < \delta < 1/2$, $N \geq 1$ and $s > 0$. There exists a positive integer $j < \delta^{-O(1)}$ such that the sequence $\xi h(s)^n$, with n in any subinterval of $n \leq N$ of length N/j , is δ -equidistributed w.r.t. a continuous algebraic measure on X unless either $(sN)^{1/5}s^{-1} < \delta^{-O(1)}$ or the conditions in Lemma 2.10 are satisfied. Moreover both j and the measure are the ones in Theorem 4.6—the continuous case.*

Proof. Let $\xi = \Gamma g$. As before, the case $\delta^{-O(1)} > T = sN$ is trivial. Otherwise, notice that if we show that $\xi h(s)^n, n \leq N$ is δ -equidistributed w.r.t the probability measure carried by $\xi h(t), t \in [0, T]$, by Theorem 4.6 we are done. Thus, by Propositions 4.11 and 4.10 we get the result unless both $y_T \delta^{-O(1)} > s^{-2}T^{-2\theta-o(1)}$ and $|W_T y_T| \delta^{-O(1)} >$

$(T/s)^{2-o(1)}y_T^{2\theta}$, with $T = sN$. In this case, we have that

$$y_T > \delta^{-O(1)}/N, \quad |W_T y_T| > \delta^{-O(1)}(y_T T^{3/2} + sT),$$

in the range $\delta^{-O(1)}s < T^{1/5}$, since $\theta < 1/5$, so

$$\Gamma g h(s)^n = h(x + y_T s n + y_T \frac{s^2 n^2}{W_T}) a(y_T) + O(\delta^2)$$

and we can apply Proposition 4.9, which gives that $\xi h(s)^n, n \leq N$ is δ -equidistributed w.r.t. the algebraic measure on the closed horocycle of period y_T^{-1} unless there exists an integer $q < (\delta/\tau(q_2))^{-O(1)}y_T^{-1/2}$ such that

$$\|qx\| + N\|qsy_T\| + (sN)^2 qy_T/|W_T| < y_T^{\frac{1}{2}}(\delta/\tau(q_2))^{-O(1)}.$$

Actually we must have $q > (\delta/\tau(q_2))^{O(1)}y_T^{-1/2}$ because ξ is fixed. But these are the conditions in Lemma 2.10. \square

Remark. The proof actually shows that $\xi h(s)^n, n \leq N$ is δ -equidistributed w.r.t the measure carried by $\xi h(t), t \in [0, sN]$ unless either $(sN)^{1/5}s^{-1} < \delta^{-O(1)}$ or the conditions in Lemma 2.10 are satisfied.

As a corollary of this result we obtain Dani's Theorem for discrete orbits.

Corollary 4.13 (Dani Theorem). *Let $\xi \in X$ with α irrational. Then $\xi h(s)^n, n \leq N$ becomes equidistributed w.r.t the volume measure on X as N goes to infinity.*

5. EFFECTIVE EQUIDISTRIBUTION: NON-CONTINUOUS ALGEBRAIC MEASURES

In the previous section we gave necessary conditions on ξ for the probability measure carried by any large piece of the orbit $\xi h(s)^n, n \leq N$ to be near to some continuous algebraic measure. This is all we need in order to prove our results for primes. However, it remains to answer the following questions: is it possible for that measure to be always near to an—either continuous or not—algebraic measure? Are the conditions on ξ really sharp for the measure to be near to a continuous algebraic measure? In this section we will show that the answer to both questions is essentially “yes”, as long as $(sN)^{1/5}s^{-1} > \delta^{-O(1)}$.

Theorem 4.12 says that, for $(sN)^{1/5}s^{-1} > \delta^{-O(1)}$, if the sequence $\xi h(s)^n, n \leq N$ is not equidistributed w.r.t. a continuous algebraic measure, then it is near to a sequence

$$(5.1) \quad s_n = \frac{A + Bn}{q} + \epsilon_4 \frac{\epsilon_1 + \epsilon_2 \frac{n}{N} + \epsilon_3 (\frac{n}{N})^2 + i\epsilon_4}{q^2}$$

for some integers A, B, q with $(A, B, q) = 1$ and $\epsilon_j \ll (\tau(q_2)\delta^{-1})^{O(1)}$, with $\frac{B}{q} = \frac{a_2}{q_2}$, $(a_2, q_2) = 1$. We also have $\epsilon_4/q^2 \gg (1+s)^{-2-o(1)}$ from the remarks after Lemma 2.10. A key ingredient to understand s_n is going to be to understand its restriction to an interval $L < n \leq L+q$; there it is very near to a sequence

$$(5.2) \quad \frac{A+Bn}{q} + M_2\left(\frac{M_1}{q^2} + i\frac{M_2}{q^2}\right) \quad \text{with } n \bmod q_2,$$

where $M_1, M_2 \ll \tau(q_2)^{O(1)}$. We can describe the points

$$\frac{A+Bn}{q} \quad n \bmod q_2$$

mod 1 in a more convenient way. First, we can write

$$\frac{A+q_1a_2n}{q_1q_2} \quad n \bmod q_2$$

with $q_1 = q/q_2$ and then $(A, q_1) = 1$. Since a_2 is coprime to q_2 we can write

$$\frac{A+q_1n}{q_1q_2} \quad n \bmod q_2.$$

Next, we can write $q_2 = q'_2q''_2$ with q'_2 coprime to q_1 and q_1 a multiple of every prime dividing q''_2 . Since $(q_1, q'_2) = 1$, the equation

$$A + q_1x \equiv 0 \pmod{q'_2}$$

has solution in x , and then we can write

$$\frac{q'_2h + q_1n}{q_1q_2} \quad \text{with } n \bmod q_2$$

with $(h, q_1) = 1$ and then $(h, q_1q''_2) = 1$. Finally, writing $m = vq'_2 + uq''_2$ with u, v integers the points in (5.2) can be described as

$$\frac{u}{q'_2} + \frac{v}{q''_2} + \frac{h}{q_1q''_2} + M_2\left(\frac{M_1}{q^2} + i\frac{M_2}{q^2}\right) \quad u \bmod q'_2, \quad v \bmod q''_2.$$

To further study the behaviour of these points in X it is natural to split them into classes C_l for any l a divisor of q'_2 , each C_l corresponding to the points such that $(u, q'_2) = l$.

For any point g of $SL(2, \mathbb{R})$ of the form;

$$g = \left(\frac{b}{q} + M_2\frac{M_1 + iM_2}{q^2}, 0\right)$$

with $(b, q) = 1$, we can multiply to the left by $\gamma_{b/q} \in SL(2, \mathbb{Z})$, with

$$\gamma_{b/q} = \begin{bmatrix} -\bar{b} & * \\ q & -b \end{bmatrix}$$

obtaining

$$(5.3) \quad g^* = \gamma_{b/q}g = \left(-\frac{\bar{b}}{q} - \frac{M_1/M_2}{M_1^2 + M_2^2} + \frac{i}{M_1^2 + M_2^2}, -2 \operatorname{arccot} \frac{M_1}{M_2}\right).$$

So in our setting it is easy to see that $g \mapsto g^*$ sends the points in C_1 to

$$\left(\frac{u}{q'_2} + \frac{\overline{q'_2^2} q_1 v + \bar{h}}{q_1 q_2''} - \frac{M_1/M_2}{M_1^2 + M_2^2} + \frac{i}{M_1^2 + M_2^2}, -2 \operatorname{arccot} \frac{M_1}{M_2}\right).$$

where \bar{x} is the inverse modulo $q_1 q_2''$. One can rewrite these points as

$$\left(\frac{u}{q'_2} + \frac{\overline{q'_2^2} (q_1 v + \tilde{h})}{q_1 q_2''} - \frac{M_1/M_2}{M_1^2 + M_2^2} + \frac{i}{M_1^2 + M_2^2}, -2 \operatorname{arccot} \frac{M_1}{M_2}\right).$$

with \tilde{h} an inverse of h modulo q_1 . But then this is the same as

$$\left(\frac{n}{q_2} + \frac{\overline{q'_2^2} \tilde{h}}{q_1 q_2''} - \frac{M_1/M_2}{M_1^2 + M_2^2} + \frac{i}{M_1^2 + M_2^2}, -2 \operatorname{arccot} \frac{M_1}{M_2}\right)$$

with $n \bmod q_2$, $(n, q'_2) = 1$. In general, for any $l \mid q'_2$, everything works the same way but dividing q'_2 , M_1 and M_2 by l , and then the points of C_l can be seen as

$$(5.4) \quad \left(\frac{n}{q_2/l} + l^2 z, -2 \operatorname{arccot} \frac{M_1}{M_2}\right),$$

with $n \bmod q_2/l$, $(n, q'_2/l) = 1$ and

$$z = \frac{\overline{q'_2^2} \tilde{h}}{q_1 q_2''} - \frac{M_1/M_2}{M_1^2 + M_2^2} + \frac{i}{M_1^2 + M_2^2}.$$

We are finally prepared to state the main result concerning the near to a closed horocycle case

Proposition 5.1. *Let $0 < \delta < 1/2$. There exists $j < \delta^{-O(1)}$ such that when n is restricted to any subinterval of $n \leq N$ of length N/j the sequence s_n is δ -equidistributed w.r.t. some algebraic measure in X .*

Proof. If $|\epsilon_2| + |\epsilon_3| < \delta^{-O(1)} \epsilon_4$ then s_n and s_m are at distance $O(\delta^{-O(1)} |n - m|/N)$, hence the result follows—each algebraic measure is supported on a point. Thus, from now on we assume $|\epsilon_2| + |\epsilon_3| \gg \delta^{-O(1)} \epsilon_4$.

Let us treat first the case $q_2 \ll \delta^{-O(1)}$; if n is in an interval J of length $\delta^c N$ containing the point N' (we should choose N' such that $t = N'/N$ makes $|\epsilon_1 + \epsilon_2 t + \epsilon_3 t^2| \gg |\epsilon_1| + |\epsilon_2| + |\epsilon_3|$), we select the sequence

$$r_n = \frac{A + Bn}{pq} + \epsilon_4 \frac{\epsilon_1 + \epsilon_2 \frac{N'}{N} + \epsilon_3 \left(\frac{N'}{N}\right)^2 + i\epsilon_4}{(pq)^2} \quad n \bmod pq_2$$

for any prime $p > \delta^{-O(1)}$, which carries an algebraic measure μ . Let us see that s_n in the interval J is δ -equidistributed w.r.t. μ . For that, let us split J into arithmetic progressions $n \equiv n_0 \pmod{q_2}$; for each one we have (due to (5.3))

$$\Gamma s_n = \Gamma(x_0 - \frac{b_n}{\epsilon_4} \frac{q_0^2}{b_n^2 + \epsilon_4^2} + \frac{i q_0^2}{b_n^2 + \epsilon_4^2}, -2 \operatorname{arccot} \frac{b_n}{\epsilon_4})$$

for $b_n = \epsilon_1 + \epsilon_2 \frac{n}{N} + \epsilon_3 (\frac{n}{N})^2$ and x_0 and $q_0 \mid q$ depending on n_0 . Then

$$\Gamma s_n = x_0 - \frac{\epsilon_4^{-1} q_0^2}{b_n^2} + \frac{i q_0^2}{b_{N'}^2} + O(\delta)$$

for most n 's in the arithmetic progression contained in J . Since $\delta^{-O(1)} < |\epsilon_4^{-1}| < \delta^{-O(1)} s^{2+o(1)} < N^{O(1)}$ by classical methods in exponential sums one can prove (see for instance [21])

$$\mathbb{E}_{n \in J, n \equiv n_0 \pmod{q_2}} e(k \frac{\epsilon_4^{-1} q_0^2}{b_n^2}) \ll \delta^{c'}$$

for large c' and any $k < \delta^{-O(1)}$. Then, we have that s_n , $n \in J$, $n \equiv n_0 \pmod{q_2}$ is δ -equidistributed w.r.t. the measure carried by

$$t + i \frac{q_0^2}{b_{N'}^2} \quad t \in [0, 1].$$

On the other hand we have, again by (5.3), that r_n with $n \in J$, $n \equiv n_0 \pmod{q_2}$ is δ -equidistributed w.r.t. the measure carried by

$$(x'_0 + \frac{m}{p} - \frac{b_{N'}}{\epsilon_4} \frac{q_0^2}{b_{N'}^2 + \epsilon_4^2} + i \frac{q_0^2}{b_{N'}^2 + \epsilon_4^2}, -2 \operatorname{arccot} \frac{b_{N'}}{\epsilon_4}) \quad m \pmod{p}$$

for some x'_0 depending on n_0 . Since any two consecutive points of this sequence are at distance $p^{-1} \delta^{-O(1)} < \delta^2$, we have that r_n $n \in J$, $n \equiv n_0 \pmod{q_2}$ is δ -equidistributed w.r. t. the measure carried by

$$t + i \frac{q_0^2}{b_{N'}^2} \quad t \in [0, 1],$$

so s_n , $n \in J$ is δ -equidistributed w.r.t. μ .

It remains the case $|\epsilon_2| + |\epsilon_3| > \delta^{-O(1)} \epsilon_4$ and $q_2 > \delta^{-O(1)}$. In this case we can take the algebraic measure μ carried by

$$r_n = \frac{A + Bn}{q} + \epsilon_4 \frac{\epsilon_1 + \epsilon_2 \frac{N'}{N} + \epsilon_3 (\frac{N'}{N})^2 + i \epsilon_4}{q^2} \quad n \pmod{q_2},$$

to approximate s_n , $n \in J$. To see that s_n , $n \in J$ is δ -equidistributed w.r.t. μ , we shall show the same restricting n to any subinterval of J

of length q_2 ; since $q < \delta^{-O(1)}y^{-\frac{1}{2}}$ and $y^{-1} \ll sN$ we have that $q < \delta N$; then s_n restricted to a subinterval of J of length q_2 satisfies

$$s_n = \frac{A + Bn}{q} + \epsilon_4 \frac{\epsilon_1 + \epsilon_2 \frac{N''}{N} + \epsilon_3 \left(\frac{N''}{N}\right)^2 + i\epsilon_4}{q^2} + O(\delta) \quad n \bmod q_2,$$

with $N'' \in J$. Now we restrict n further to the set C_l , for some $l \mid q_2'$, of n 's satisfying $(n, q_2') = l$; it is necessary to look just to the l 's with $l < \tau(q_2')\delta^{-O(1)}$, because the others give a negligible contribution. Due to (5.4) the restriction of r_n to C_l is

$$(5.5) \quad \left(\frac{n}{q_2/l} + l^2 x_1 - \frac{b_{N'}}{\epsilon_4} \frac{l^2}{b_{N'}^2 + \epsilon_4^2} + i \frac{l^2}{b_{N'}^2 + \epsilon_4^2}, -2 \operatorname{arccot} \frac{b_{N'}}{\epsilon_4} \right).$$

with $n \bmod q_2/l$, $(n, q_2'/l) = 1$, and $x_1 = \frac{\overline{q_2'}^2 \tilde{h}}{q_1 q_2'}$. Now, it is easy to show that for any $x, \epsilon > 0$ and $j, d \in \mathbb{N}$ we have

$$\sum_{x < n \leq x + \epsilon j d, (n, d) = 1} 1 = \epsilon \sum_{n \leq j d, (n, d) = 1} 1 + O(\tau(d)).$$

Applying it for $d = q_2'/l$, $j = q_2/q_2'$ and $\epsilon = \delta l^2 / (b_{N'}^2 + \epsilon_4^2)$, since $\epsilon j \phi(d) \gg \delta^{-O(1)}\tau(d)$, we have that the sequence in (5.5) is δ -equidistributed w.r.t. the measure carried by

$$\left(t + i \frac{l^2}{b_{N'}^2 + \epsilon_4^2}, -2 \operatorname{arccot} \frac{b_{N'}}{\epsilon_4} \right) \quad t \in [0, 1].$$

One can proceed in the same way for the restriction of s_n , obtaining that it is δ -equidistributed w.r.t. the measure carried by

$$\left(t + i \frac{l^2}{b_{N''}^2 + \epsilon_4^2}, -2 \operatorname{arccot} \frac{b_{N''}}{\epsilon_4} \right) \quad t \in [0, 1].$$

But, since both N' and N'' are in J we have that both measures are similar, and then s_n , $n \in J$ is δ -equidistributed w.r.t. μ . \square

As a corollary of Theorem 4.12 and Proposition 5.1 we obtain

Theorem 5.2 (Effective equidistribution). *Let $\xi \in X$ fixed. Let $0 < \delta < 1/2$, $N > 1$ and $s > 0$. There exists a positive integer $j < \delta^{-O(1)}$ such that the sequence $\xi h(s)^n$, with n in any subinterval of $n \leq N$ of length N/j , is δ -equidistributed w.r.t. an algebraic measure on X , unless $(sN)^{1/5} s^{-1} < \delta^{-O(1)}$.*

Remarks. (a) As shown in the proofs of this section, in the quantitative setting the non-continuous algebraic measures are much more complex than continuous algebraic ones. Therefore, to be near to a continuous algebraic measure is a condition that is much stronger than to be near to a general algebraic measure.

(b) This result gives control over pieces of orbits of the discrete horocycle flow for s not very large. It is possible to show that this control fails for $s > N^{3+\epsilon}$; in fact, taking $x = 0$, $y = q^{-2}$, $s = Aq^{-2}$ and $W^{-1} = A^{-2}q^{-3}$ for A, q natural numbers, $q < N^\epsilon$, $A > N^{3+\epsilon}$ we have that $gh(s)^n$ is very near to the periodic sequence

$$\frac{n^2}{q} + \frac{i}{q^2} \quad n \bmod q.$$

One can show that for certain q 's the measure carried by this sequence is not near to any algebraic measure. The same happens for Theorem 4.12. Perhaps both results remain true for $s < N^{O(1)}$ by substituting algebraic measures by “polynomial algebraic measures”—and changing the conditions in the statement of Theorem 4.12—, meaning any measure carried by a periodic sequence $\Gamma h(p(n))a(y)$ with p a polynomial of degree $O(1)$.

6. LARGE CLOSURE OF PRIME ORBITS

In this section we are going to deduce Theorem 1.1 from an upper bound for the sum

$$\sum_{p < T} f(xu^p).$$

In order to get such a bound we make use of sieve theory, in particular the following special case of a Selberg's result (see [21, Theorem 6.4])

Lemma 6.1 (Upper bound Sieve). *Let $(a_n)_{n \leq T}$ a sequence of non-negative numbers. For any $d \in \mathbb{N}$ write*

$$\sum_{\substack{n \leq T \\ n \equiv 0 \pmod{d}}} a_n = \frac{1}{d}A + r_d.$$

Then, for any $1 < D < T$ we have

$$\sum_{\sqrt{T} < p \leq T} a_p \leq \frac{A}{\log \sqrt{D}} + \sum_{d < D} \tau_3(d)|r_d|,$$

with $\tau_3(d) = \sum_{d_1 d_2 d_3 = d} 1$.

So the task now is reduced to have tight control over the sums

$$\sum_{m \leq T/d} f(x(u^d)^m)$$

for most of the d 's in a range $1 \leq d \leq D$ with D as big as possible. This can be handled by the Theorem 4.12 whenever $D < T^{1/5}$. The precise result that follows is

Theorem 6.2. *Let $\xi \in X$ with α irrational, $f \geq 0$ and $s > 0$. Then*

$$(6.1) \quad \mathbb{E}_{p < T} f(\xi h(s)^p) \leq 10 \int_X f d\mu_G + o_T(1) \|f\|_{\text{Lip}}.$$

Proof. Let $\|f\|_{\text{Lip}} = 1$. Let us begin by dealing with the case in which the conditions in Lemma 2.10 (i) are satisfied for ξ , s , $N = T$ and $\delta = (\log T)^{-A}$, A a large constant. Take $q < \delta^{-O(1)}$ the smallest integer satisfying the conditions in the lemma. It is easy to check that q has to go to infinity with T ; moreover by Lemma 2.10 (ii) and the remarks after the lemma one sees that $\xi h(s)^n$ is at distance δ from a q_2 -periodic sequence, with $q_2 \mid q \lfloor q \frac{s}{R} \rfloor^2 \ll q^3 < \delta^{-O(1)}$, when n is restricted to any subinterval J such that $|J| = \delta^c T$ for some constant c . In this situation we can apply the Siegel-Walfisz theorem [9, page 133] for primes in arithmetic progressions to deduce that

$$\mathbb{E}_{p < T} f(\xi h(s)^p) = \mathbb{E}_{n < T, (n, q_2) = 1} f(\xi h(s)^n) + O((\log T)^{-1}).$$

The conditions of Lemma 2.10 are not satisfied for the parameters $s_* = sd$, $N_* = N/d$, $\delta_* = q^{-\epsilon}$, if $d < q^\epsilon$ and $\epsilon > 0$ is sufficiently small; we can then apply Theorem 4.12 to deduce that

$$(6.2) \quad \mathbb{E}_{n < T, n \equiv 0 \pmod{d}} f(\xi h(s)^n) = \int f d\mu + O(q^{-\epsilon}),$$

for any $d < q^\epsilon$, where $d\mu$ is the average of the algebraic measures appearing in the statement and then it is independent of d . Therefore, using the identity $1_{m=1} = \sum_{d|m} \mu(d)$ and (6.2) for $d < q^\epsilon$ we have

$$\mathbb{E}_{n < T, (n, q_2) = 1} f(\xi h(s)^n) = \mathbb{E}_{n < T} f(\xi h(s)^n) + O(\tau(q_2)q^{-\epsilon})$$

and then the result follows from Corollary 4.13.

Now let us suppose that the conditions in Lemma 2.10 are not satisfied. Then, applying Theorem 4.12 we have

$$\mathbb{E}_{n < T} f(\xi h(s)^n) = \int f d\mu + O((\log T)^{-A}).$$

with μ the average of the algebraic measures there. Let us suppose that for any $D \leq D_0 = T^{1/5}$ and for any d in $D < d < 2D$ but at most $O(\delta D)$ exceptions, the conditions in Lemma 2.10 are not satisfied for the parameters $s_* = sd$ and $N_* = N/d$. For any of the non-exceptional d we have (again by Theorem 4.12)

$$\mathbb{E}_{n < T, n \equiv 0 \pmod{d}} f(\xi h(s)^n) = \int f d\mu + O((\log T)^{-A}).$$

Therefore, applying Lemma 6.1 we get

$$\mathbb{E}_{p < T} f(\xi h(s)^p) \leq 10 \mathbb{E}_{n < T} f(\xi h(s)^n) + O((\log T)^{O(1)-A})$$

and we are done by Corollary 4.13. So, we have finished unless

$$\|dq_d \frac{s}{R_d}\| = D\tau_d^{O(1)}\delta^{-O(1)}T^{-1}, \quad \left\| \left[dq_d \frac{s}{R_d} \right] \alpha_d \right\| = D\tau_d^{O(1)}\delta^{-O(1)}(sT^2)^{-1}$$

for more than δD d 's in $D < d < 2D$ for some D , where $R_d = R(\gamma_d g_0)$, $\alpha_d = \alpha(\gamma_d g_0)$, $\Gamma g_0 = \xi$, $q_d \ll \tau_d^{O(1)}\delta^{-O(1)}$, and $\tau_d = \tau(\widetilde{(q_d)_2})$. Let us see that this cannot be true for any $D < T^{1-\epsilon}$, $\epsilon > 0$.

Let us assume it is true. We know that $\tau_d < L = D^{o(1)}$. Then we can split $[1, L]$ into at most $O(\log \log L)$ intervals of the shape $[t, 2t^2]$, and there will be at least $\delta D(\log \log L)^{-1} \gg \delta^2 D$ d 's with τ_d in one of them; let us consider now just those d 's; for any of them we have $\tau_* \leq \tau_d \leq \tau_*^2$ for some fixed $\tau_* < L$.

This implies that $\gamma_d = \gamma$ and $q_d = q$ for a set \mathcal{A} of more than D/M d 's, with $M = \tau_*^{O(1)}\delta^{-O(1)}$. So for them

$$\|dq \frac{s}{R}\| = DMT^{-1} \quad \left\| \left[dq \frac{s}{R} \right] \alpha \right\| = DM(sT^2)^{-1},$$

and then effective equidistribution in the torus (see Lemma 3.2 in [13]) implies that

$$(6.3) \quad \|hq \frac{s}{R}\| < MT^{-1}$$

for some $h < M$. We can assume h is coprime to $[hqs/R]$. Now, since $M^{-2} < DMT^{-1}$ it is easy to check that $h \mid d$ for any $d \in \mathcal{A}$. But then $d = \lambda h$, we have $[dqs/R] = \lambda[hqs/R]$ and

$$\|\lambda[hqs/R]\alpha\| < DM(sT^2)^{-1}$$

for more than D/M λ 's with $\lambda \ll D$, which again by effective equidistribution in the torus gives

$$\|j[hqs/R]\alpha\| < M(sT^2)^{-1}$$

for some $j < M$. Now $j[hqs/R] = [jhqs/R]$ by (6.3), so choosing $q_* = jq < M$ we have

$$(6.4) \quad \|q_* \frac{s}{R}\| = MT^{-1} \quad \left\| \left[q_* \frac{s}{R} \right] \alpha \right\| = M(sT^2)^{-1}.$$

Now, it is easy to check that for any $d \in \mathcal{A}$ we have

$$\frac{[[dq \frac{s}{R}]\alpha]}{[dq \frac{s}{R}]} = \frac{[[q_* \frac{s}{R}]\alpha]}{[q_* \frac{s}{R}]},$$

which implies that $\widetilde{(q_d)_2} = k \leq [q_* \frac{s}{R}] \ll q_*$. Since $\tau_* < \tau(k) < \tau_*^2$ we arrive at

$$k \ll q_* < \tau_*^{O(1)}\delta^{-O(1)} < \tau(k)^{O(1)}\delta^{-O(1)}$$

so $k < \delta^{-O(1)}$ and then $\tau_* < \delta^{-O(1)}$. But then (6.4) means that the conditions in Lemma 2.10 are satisfied for the original sequence, which is a contradiction. \square

Remark: Theorem 4.12 was not strictly necessary to prove this result (and then Theorem 1.1); one could proceed in a more direct way, taking advantage of the extra average in d in Lemma 6.1. Anyway, it seems difficult to get a much better level in Lemma 6.1 that way. We did not do it that way because we think Theorem 4.12 is interesting by itself.

Finally, let us see that Theorem 1.1 follows from Theorem 6.2. Let $x \in X$ generic, and ν_x an accumulation point for the sequence $(\pi_{x,N})_N$ in $C^*(X^*)$, where X^* is the one-point compactification of X . Take $f \in C(X^*)$, $f \geq 0$. By approximation, we can assume that f as finite Lipschitz norm in X , so that Theorem 6.2 gives

$$\int_X f d\nu_x \leq 10 \int_X f d\mu_G,$$

and we are done.

7. DENSITY OF THE HECKE ORBIT

In this section we are going to prove a stronger result (Theorem 1.3) for the special orbit $H_N h(p)$, $p \leq N$, from which we can deduce in particular that it becomes dense in X when $N \rightarrow \infty$. This will be possible because we can get a good level of distribution for linear sums, and above all because we can handle bilinear sums up to a considerable level. We input those bounds into the following special case of a sieve result from [10]

Lemma 7.1 (Asymptotic sieve). *Let $\{a_n\}_{n \in \mathbb{N}}$ be a sequence of non-negative numbers such that $a_n \ll \tau(n)$ and $a_n = A + c_n$ for some constant $A \geq 0$ and a sequence c_n satisfying the “Type I condition of level α ”*

$$(7.1) \quad \mathbb{E}_{D < d < 2D} |\mathbb{E}_{n \leq x/d} c_{dn}| \ll (\log x)^{-3}$$

for any $D < x^{\alpha-\epsilon}$ and also the “Type II condition of level γ ”

$$(7.2) \quad \mathbb{E}_{D < d_1 < d_2 < 2D} |\mathbb{E}_{n \leq \min(x/d_1, x/d_2)} c_{d_1 n} \overline{c_{d_2 n}}| \ll (\log x)^{-22},$$

for any D with $x^{(\log \log x)^{-3}} \leq D \leq x^{\gamma-\epsilon}$, for any fixed $\epsilon > 0$. Then we have

$$(7.3) \quad |\mathbb{E}_{p < x} a_p - A| \leq c(\alpha, \gamma)A + O(\epsilon),$$

with $c(\alpha, \gamma)$ an explicit decreasing function, such that $c(1/2, 1/3) = 0$, $c(1/2, 1/5) < 4/5$ and $c(1/2, \gamma) = 1$ for some $\gamma \in (1/6, 1/5)$.

Moreover, in the summations we can assume that d is square free and d_1, d_2 are primes.

Now, we are going to check the Type I condition.

Proposition 7.2 (Bound for Type I sums). *Let f with $\int_X f = 0$. We have that*

$$\mathbb{E}_{D < d < 2D} |\mathbb{E}_{n \leq N/d} f(H_N h(dn))| \ll (\log N)^{-3} \|f\|_{\text{Lip}}$$

for any $D < N^{1/2-\theta-\epsilon}$, for any $\epsilon > 0$.

Proof. Let us assume $\|f\|_{\text{Lip}} = 1$, and suppose the result is false. Then, we have

$$\mathbb{E}_{D < d < 2D} |\mathbb{E}_n f(\Gamma H_N h(dn)) \eta(\frac{Dn}{N})| \gg \delta$$

for $\delta = (\log N)^{-A}$, A a large constant, with $\eta \in C_0^\infty(-2, 2)$ with $\|\eta^{(j)}\|_{L^\infty}, \|f\|_{W^j} \ll_j \delta^{-j}$. We have the Iwasawa parametrization

$$H_N h(dn) = h(dnN^{-1})a(N^{-1}),$$

and then by the Fourier expansion $f(h(x)a(y)) = \sum_m \tilde{f}(m, y)e(mx)$ and Poisson Summation we have

$$\mathbb{E}_{D < d < 2D} \sum_m |\tilde{f}(m, N^{-1})| |\hat{\eta}(\frac{N}{D} \{ \frac{dm}{N} \})| \gg \delta$$

so that from Proposition 3.1 we have

$$\delta^{-O(1)} D^{-1} N^{-1/2+\theta_2} \sum_{j \ll \delta^{-1} DN} \tau(j)^2 |\hat{\eta}(\frac{N}{D} \{ \frac{j}{N} \})| \gg 1$$

which taking into account the decay of $\hat{\eta}$ gives a contradiction. \square

Let us go with the Type II condition

Proposition 7.3 (Bound for Type II sums). *Let f_1, f_2 be continuous functions in $\Gamma \backslash G$ with $\|f_1\|_{\text{Lip}} = \|f_2\|_{\text{Lip}} = 1$ and $\int_X f_1 = 0$. Let $N^{(\log \log N)^{-3}} < D < N^{(1-2\theta)/(5+2\theta)-\epsilon}$ and $D < d_1 < d_2 < 2D$, with d_1, d_2 primes. Then we have*

$$\mathbb{E}_{n \leq \min(N/d_1, N/d_2)} f_1(H_N h(d_1 n)) f_2(H_N h(d_2 n)) \ll (\log N)^{-22}.$$

Proof. As in (4.1) we can replace the Lipschitz norm by Sobolev norms and hence if the result is false, we have

$$|\mathbb{E}_{n < N/D} f_1(H_N h(d_1 n)) f_2(H_N h(d_2 n)) \eta(\frac{Dn}{N})| \gg \delta$$

for $\delta = (\log N)^{-A}$, A a large constant, with $\eta \in C_0^\infty(0, 1)$ with $\|\eta^{(j)}\|_{L^\infty} \ll \delta^{-j}$ and $\|Df_i\|_{L^\infty} \ll_{\text{ord}D} \delta^{-\text{ord}D}$. By the Fourier expansion of f_i ,

$f_i(h(x)a(y)) = \sum_m \tilde{f}_i(m, y)e(mx)$, the bounds in Proposition 3.1, Poisson summation in n and integration by parts we have that either

$$\left| \sum_{d_1 m_1 + d_2 m_2 = 0} \tilde{f}_1(m_1, N^{-1}) \tilde{f}_2(m_2, N^{-1}) \right| \gg \delta$$

or

$$(7.4) \quad \left| \sum_{d_1 m_1 + d_2 m_2 = k} \tilde{f}_1(m_1, N^{-1}) \tilde{f}_2(m_2, N^{-1}) \right| \gg \delta D^{-2}$$

for some $k \neq 0$, $k \ll \delta^{-O(1)} DN$. If the first possibility is true, we get

$$\left| \sum_j \tilde{f}_1(d_1 j, N^{-1}) \tilde{f}_2(d_2 j, N^{-1}) \right| \gg \delta.$$

But from Proposition 3.1 we have $\tilde{f}_i(0, N^{-1}) \ll \delta^{-O(1)} N^{-\frac{1}{2}}$, and from the spectral expansion (3.29) of f_i , the multiplicativity of Hecke eigenvalues (4.5) and Parseval (3.50) we have

$$(7.5) \quad \sum_{j \neq 0} \tilde{f}_1(d_1 j, N^{-1}) \tilde{f}_2(d_2 j, N^{-1}) \ll (d_1 d_2)^{\theta - \frac{1}{2}},$$

which gives a contradiction.

Let us now assume that (7.4) is true. This is a shifted convolution sum, and we can proceed as in [3, 4]. We can translate it as

$$\left| \int_0^1 f_1(h(d_1 x)a\left(\frac{1}{N}\right)) f_2(h(d_2 x)a\left(\frac{1}{N}\right)) e(-kx) dx \right| \gg \delta D^{-2},$$

and further as

$$\left| \int f_*(h(x)a\left(\frac{1}{DN}\right)) e(-kx) dx \right| \gg \delta D^{-2},$$

with $f_* = f_{d_1} f_{d_2}$ and

$$f_d(g) = f(a(d) g a(D/d)).$$

f_d is a $\Gamma_0(d)$ -invariant function and since $D/d \asymp 1$ we have $\|Df_d\|_{L^\infty} \ll \|Df\|_{L^\infty}$; thus, f_* can be seen as a function in $\Gamma_0(d_1 d_2) \backslash G$ with Sobolev norms $\|f_*\|_{W^j} \ll_j (d_1 d_2)^{\frac{1}{2} + \epsilon} \delta^{-j}$. Therefore, from Proposition 3.1 we get the bound

$$\tilde{f}_*(k, (DN)^{-1}) \ll \delta^{-O(1)} (DN)^{-\frac{1}{2} + \theta} (d_1 d_2)^{\frac{1}{2} + \epsilon}$$

which gives a contradiction in our range for D . \square

Now, assuming that $\theta = 0$, due to Propositions 7.2 and 7.3 we can apply Lemma 7.1 with $\alpha = 1/2$ and $\gamma = 1/5$ for $a_n = f(\Gamma H_N h(n))$, and then

Theorem 7.4. *Assuming $\theta = 0$, for any non-negative f we have*

$$|\mathbb{E}_{p < N} f(H_N h(p)) - \int_X f d\mu_G| \leq \frac{4}{5} \int_X f d\mu_G + o(1) \|f\|_{\text{Lip}}$$

From this result we can deduce Theorem 1.3 as we did in the previous section with Theorem 1.1.

We end by discussing some improvements on the levels of distribution in Propositions 7.2 and 7.3. For the type I sums we will establish a level of $1/2$ matching what Proposition 7.2 gives with $\theta = 0$. For the type II sums, Proposition 7.3 with $\theta = 7/64$ yields a level of $25/167 = 0.1457\dots$. We establish a level of $3/19 = 0.1578\dots$ which is a small improvement but still falls short of the magic number c ($\frac{1}{6} < c < \frac{1}{5}$) which would make Theorem 7.4 unconditional.

Proposition 7.5. *The result in Proposition 7.2 is true for any $D < N^{1/2-\epsilon}$, for any $\epsilon > 0$.*

Proof. In what follows we assume that f is orthogonal to the Eisenstein series, because for them we have $\theta = 0$; moreover, for simplicity let us consider f K -invariant. The sums in Proposition 7.2 are (assuming we have a smooth sum, as we can)

$$I = \frac{1}{N} \sum_{d \gtrsim D} \left| \sum_n \eta\left(\frac{Dn}{N}\right) f(H_N h(dn)) \right|.$$

We can also assume that f is orthogonal to the space of Eisenstein series (because for them we know that $\theta = 0$), and then we can write the Fourier expansion

$$f(h(x)a\left(\frac{1}{N}\right)) = \sum_k \tilde{f}\left(k, \frac{1}{N}\right) e(kx) = O(\delta) + \sum_k^* \tilde{f}\left(k, \frac{1}{N}\right) e(kx)$$

with $\delta = N^{-1/\log \log N}$ and the sum in k restricted to $|k| < \delta^{-O(1)}N$ and $(k, N) < \delta^{-O(1)}$. This is done by using the spectral expansion, the multiplicativity of $\lambda_\pi(k)$ and the bounds (3.37) and (3.49). Thus it is enough to bound

$$I_* = \frac{1}{N} \sum_{d \gtrsim D} \left| \sum_k^* \tilde{f}\left(k, \frac{1}{N}\right) \sum_n \eta\left(\frac{Dn}{N}\right) e\left(\frac{dkn}{N}\right) \right|.$$

Applying Poisson Summation in n yields

$$(7.6) \quad I_* \ll \frac{1}{D} \sum_{d \gtrsim D} \sum_{k, \nu}^* |\tilde{f}(k, N^{-1})| |\widehat{\eta}\left(\frac{dk - \nu N}{D}\right)| \ll \frac{1}{D} \sum_k^* c_k |\tilde{f}\left(k, \frac{1}{N}\right)|,$$

with c_k the number of $d, t \ll D$ such that $dk \equiv t \pmod{N}$. For $D \leq N^{\frac{1}{2}-\epsilon}$ which we assume is in force, one can check that for each k as above

$c_k \ll N^\epsilon$ as follows: if $t = 0$ there are no solutions for the congruence since $(k, N) < \delta^{-O(1)}$; otherwise, for any pair of solutions (d_1, t_1) and (d_2, t_2) we have that $d_1 t_2 \equiv d_2 t_1 \pmod{N}$, so that $d_1 t_2 = d_2 t_1$. This means that any solution is of the form $(d, t) = \lambda(d_*, t_*)$, for some (d_*, t_*) fixed; since λ divides N we are done. On the other hand, we have $\sum_k^* c_k \ll_\epsilon D^2 N^\epsilon$, so that $\sum_k^* c_k^m \ll_{\epsilon, m} D^2 N^\epsilon$.

From (3.15) and noting that $q = 1$ we have

$$|\tilde{f}(k, N^{-1})| \ll_\epsilon \delta + N^{-\frac{1}{2}+\epsilon} \sum_{|t_j| \ll \delta^{-O(1)}} |\lambda_j(k)|$$

and hence

$$(7.7) \quad I_* \ll \delta + \frac{N^{-\frac{1}{2}+\epsilon}}{D} \sum_j \sum_k^* c_k |\lambda_j(k)| \ll \frac{N^{-\frac{1}{2}+\epsilon}}{D} \sum_j (D^2 N^\epsilon)^{\frac{3}{4}} \left(\sum_k^* |\lambda_j(k)|^4 \right)^{\frac{1}{4}}.$$

It is known (see [24] and [26]) that for $x \geq 1$,

$$(7.8) \quad \sum_{|k| \leq x} |\lambda_j(k)|^4 \ll_\epsilon \lambda_j^\epsilon x,$$

hence

$$(7.9) \quad I_* \ll \delta + \delta^{-O(1)} \frac{N^{-\frac{1}{2}+\epsilon}}{D} N^{\frac{1}{4}} D^{\frac{3}{2}} = \delta + \delta^{-O(1)} N^{-\frac{1}{4}+\epsilon} D^{\frac{1}{2}}.$$

This gives a level of distribution of $1/2$ for these type I sums. \square

Proposition 7.6. *The result in Proposition 7.3 is true in the range for $N^{(\log \log N)^{-3}} < D < N^{3/19-\epsilon}$ for any $\epsilon > 0$.*

Proof. We continue with the setting described in the proof of Proposition 7.5. For the type II sums $y^{-1} = ND$, $q = D^2$ and we could try to use the improvement in Proposition 3.1 in the range $\sqrt{q}y^{-\theta} > y^{-\frac{1}{4}}$, i.e. $D(ND)^\theta > (ND)^{\frac{1}{4}}$ or $D > N^{\frac{1-4\theta}{3+4\theta}}$. For $\theta = 7/64$ this is $D > N^{9/55}$. However $\frac{9}{55} > \frac{25}{167}$ so that the improvement in this range does not give an improvement of the level $25/167$. Instead we again exploit the average over k and again we use [26]; for $x \geq 1$

$$(7.10) \quad \sum_{|m| \leq x} |\lambda_j(m)|^8 \ll_\epsilon (\lambda_j q)^\epsilon x.$$

For $d_1, d_2 \asymp D$

$$\begin{aligned}
(7.11) \quad II &= \frac{D}{N} \sum_n \eta\left(\frac{Dn}{N}\right) f_1(H_N h(d_1 n)) f_2(H_N h(d_2 n)) \\
&= \frac{D}{N} \sum_n \eta\left(\frac{Dn}{N}\right) f_*\left(\frac{n}{N}, \frac{1}{ND}\right) \\
&\ll D^{-\epsilon} + \sum_{|k| \ll ND^{1+O(\epsilon)}} \tilde{f}_*(k, (ND)^{-1}) \sum_{|\nu| \ll D^{1+\epsilon}} \hat{\eta}\left(\frac{k - \nu N}{D}\right) \\
&\ll D^{-\epsilon} + \sum_{|t|, |\nu| \leq D^{1+O(\epsilon)}} |\tilde{f}_*(t + \nu N, (ND)^{-1})|.
\end{aligned}$$

From (3.15) with $|k| \ll ND^{1+\epsilon}$ and for $k = 0$ due to (7.5) we have that

$$|\tilde{f}_*(k, (ND)^{-1})| \ll D^{-\epsilon} + \frac{(ND)^{-\frac{1}{2}}}{\sqrt{q}} \sum_{|t_j| \ll D^{O(\epsilon)}} |\lambda_j(k)| |\langle f_*, \phi_j \rangle|.$$

Hence

$$(7.12) \quad \sum_{|t|, |\nu| \leq D^{1+\epsilon}} |\tilde{f}_*(t + \nu N, (ND)^{-1})| \leq \frac{(ND)^{-\frac{1}{2}}}{\sqrt{q}} \sum_{|t_j| \ll D^\epsilon} |\langle f_*, \phi_j \rangle| \sum_{|t|, |\nu| \leq D^{1+\epsilon}} |\lambda_j(t + \nu N)|.$$

The inner sum may be estimated by Holder,

$$\sum_{|t|, |\nu| \leq D^{1+\epsilon}} |\lambda_j(t + \nu N)| \leq \left(\sum_{|m| \leq ND^{1+\epsilon}} |\lambda_j(m)|^8 \right)^{\frac{1}{8}} (D^{2+2\epsilon})^{\frac{7}{8}}$$

which by (7.10) is

$$(7.13) \quad \ll \lambda_j^\epsilon (ND)^{\frac{1}{8}} D^{\frac{7}{4}} D^{O(\epsilon)}.$$

Substituting this into (7.12) gives (recall that $q = D^2$)

$$\begin{aligned}
(7.14) \quad II &\ll D^{-\epsilon} + \frac{(ND)^{-\frac{1}{2}+\epsilon}}{\sqrt{q}} (ND)^{\frac{1}{8}} D^{\frac{7}{4}} \left(\sum_{|t_j| \ll D^\epsilon} |\langle f_*, \phi_j \rangle|^2 \right)^{\frac{1}{2}} q^{\frac{1}{2}} \\
&\ll D^{-\epsilon} + (ND)^{-\frac{1}{2}} (ND)^{\frac{1}{8}} D^{\frac{7}{4}} (D^2)^{\frac{1}{2}} N^{O(\epsilon)} = D^{-\epsilon} + N^{-\frac{3}{8}} D^{\frac{19}{8}} N^{O(\epsilon)}.
\end{aligned}$$

This gives a level of distribution of $3/19$ for the type II sums. \square

ACKNOWLEDGMENTS

We would like to thank N. Pitt for interesting discussions about his paper [29]. We also acknowledge the referees of an earlier version of this paper for pointing to some errors that needed correcting.

A. Ubis was supported by a Postdoctoral Fellowship from the Spanish Government during part of the writing of the paper; he also thanks both the Department of Mathematics of Princeton University and the Institute for Advanced Study for providing excellent working conditions. P. Sarnak and A. Ubis were supported in part by NSF grants, and A. Ubis by a MINCYT grant.

REFERENCES

- [1] C. B. Balogh, *Asymptotic expansions of the modified Bessel function of the third kind of imaginary order*, SIAM J. Appl. Math. **15** (1967) 1315–1323.
- [2] V. Blomer, *Rankin-Selberg L -functions on the critical line*, Manuscripta Math. **117**, 111–113 (2005).
- [3] V. Blomer and G. Harcos, *The spectral decomposition of shifted convolution sums*, Duke Math. J. Volume **144**, Number 2 (2008), 321–339.
- [4] V. Blomer and G. Harcos, *Twisted L -functions over number fields and Hilbert’s eleventh problem*, Geom. Funct. Anal. **20** (2010), no. 1, 1–52.
- [5] J. Bourgain, *Pointwise ergodic theorems for arithmetic sets (With an appendix by H. Furstenberg, Y. Katznelson and D. Ornstein)*, Inst. Hautes Études Sci. Publ. Math. **69** (1989), 5–45.
- [6] R. W. Bruggeman and Y. Motohashi, *A new approach to the spectral theory of the fourth moment of the Riemann zeta-function*, J. Reine Angew. Math. **579** (2005), 75–114.
- [7] M. Burger, *Horocycle flow in geometrically finite surfaces*, Duke Math. J. **61**, No. 3 (1990), 779–803.
- [8] S. Dani, *On uniformly distributed orbit of certain horocycle flows*, Ergodic Theory and Dynamical Systems (1982), **2**, 139–158.
- [9] H. Davenport, *Multiplicative number theory*, Third edition. Revised and with a preface by Hugh L. Montgomery. Graduate Texts in Mathematics, **74**. Springer-Verlag, New York, 2000. xiv+177 pp.
- [10] W. Duke, J. Friedlander and H. Iwaniec, *Equidistribution of roots of a quadratic congruence to prime moduli*, Ann. of Math. (2) **141** (1995), 423–441.
- [11] J. Friedlander and H. Iwaniec, *Opera de cribro*, American Mathematical Society Colloquium Publications, **57**. American Mathematical Society, Providence, RI, 2010. xx+527 pp.
- [12] A. Gorodnik, *Open problems in dynamics and related fields*, Journal of Modern Dynamics **1**, no.1, 1–35 (2007).
- [13] B. J. Green and T. Tao, *The quantitative behaviour of polynomial orbits on nilmanifolds*, arXiv:0709.3562v4 [math.NT].
- [14] B. J. Green and T. Tao, *The Möbius function is strongly orthogonal to nilsequences*, arXiv:0807.1736v2 [math.NT].
- [15] J. L. Hafner, *Some remarks on odd Maass Wave Forms*, Math. Z. **196** (1987), no. 1, 129–132.

- [16] G. A. Hedlund, *Fuchsian groups and transitive horocycles*, Duke Math. J. **2** (1936), no. 3, 530–542.
- [17] J. Hoffstein and P. Lockhart, *Coefficients of Maass forms and the Siegel zero*, With an appendix by Dorian Goldfeld, Hoffstein and Daniel Lieman. Ann. of Math. (2) **140** (1994), no. 1, 161–181.
- [18] L. Hörmander, *The analysis of partial linear differential operators I*, second edition, Springer Study Edition, Springer-Verlag, Berlin, 1990, xii+440 pp.
- [19] H. Iwaniec, *Nonholomorphic modular forms and their applications*, Modular forms (Durham, 1983), 157–196, Ellis Horwood Ser. Math. Appl.: Statist. Oper. Res., Horwood, Chichester, 1984.
- [20] H. Iwaniec, *Small eigenvalues of Laplacian for $\Gamma_0(N)$* . Acta Arith. **56** (1990), **1**, 65–82.
- [21] H. Iwaniec and E. Kowalski, *Analytic Number Theory*, American Mathematical Society Colloquium publications, **53**. American Mathematical Society, Providence, RI, 2004. xii+615 pp.
- [22] H. Iwaniec, W. Luo and P. Sarnak, *Low lying zeros of families of L-functions*, Inst. Hautes Études Sci. Publ. Math. **91** (2000), 55–131 (2001).
- [23] M. Jutila, *Lectures on a method in the theory of exponential sums*, Tata Institute of Fundamental Research Lectures on Mathematics and Physics, 80. Published for the Tata Institute of Fundamental Research, Bombay; by Springer-Verlag, Berlin, 1987. viii+134 pp.
- [24] H. H. Kim, *Factoriality for the exterior square of GL_4 and the symmetric fourth of GL_2* , with “Appendix 1” by D. Ramakrishnan and “Appendix 2” by H. Kim and P. Sarnak, J. Amer. Math. Soc. **16** (2003), 139–183.
- [25] H. D. Kloosterman, *On the representation of numbers in the form $ax^2 + by^2 + cz^2 + dt^2$* , Acta Math. **49** (1927), no. 3–4, 407–464.
- [26] X. Li, *Upper bounds on L-functions at the edge of the critical strip*. Int. Math. Res. Not. IMRN 2010, **4**, 727–755.
- [27] Ju. V. Linnik and B. F. Skubenko, *On the asymptotic behavior of integral matrices of third order*, Dokl. Akad. Nauk SSSR **14** 1992 1007–1008.
- [28] J. Marklof, A. Strömbergsson, *Equidistribution of Kronecker sequences along closed horocycles*, GAFA **13** (2003), 1239–1280.
- [29] N. Pitt, *On an analogue of Titchmarsh’s divisor problem for holomorphic cusp forms*, preprint 2011.
- [30] M. Ratner, *Horocycle flows, joinings and rigidity of products*, Ann. of Math. (2) **118** (1983), no. 2, 277–313.
- [31] M. Ratner, *Raghunathan’s topological conjecture and distributions of unipotent flows*, Duke Math. J. **63** (1991), no. 1, 235–280.
- [32] P. Sarnak, *Asymptotic behavior of periodic orbits of the horocycle flow and Einsenstein series*, Comm. on Pure and Appl. Math., **34** (1981), 719–739.
- [33] P. Sarnak, *Diophantine problems and linear groups*, Proceedings of the International Congress of Mathematicians, Vol. I, II (Kyoto, 1990), 459–471, Math. Soc. Japan, Tokyo, 1991.
- [34] P. Sarnak, *Some applications of modular forms*, Cambridge Tracts in Mathematics, **99**. Cambridge University Press, Cambridge, 1990. x+111p.
- [35] P. Sarnak, *Estimates for Rankin-Selberg L-functions and Quantum Unique Ergodicity*, Journal of Functional Analysis **184**, **2**, 419–453 (2001).

- [36] A. Selberg, *On the estimation of Fourier coefficients of modular forms*, 1965 Proc. Sympos. Pure Math., Vol. VIII pp. 1–15 Amer. Math. Soc., Providence, R.I.
- [37] A. Strömbergsson, *On the uniform equidistribution of long closed horocycles*, Duke Math. J. **123** (2004), 507–547.
- [38] A. Strömbergsson, *On the deviation of ergodic averages for horocycle flows*. Preprint, available at <http://www.math.uu.se/~astrombe/papers/iha.pdf>
- [39] A. Venkatesh, *Sparse equidistribution problems, period bounds, and subconvexity*, Ann. of Math. (2) **172** (2010), no. 2, 989–1094.

INSTITUTE FOR ADVANCED STUDY, EINSTEIN ROAD, PRINCETON NJ 08540,
USA

E-mail address: sarnak@math.princeton.edu

DEPARTAMENTO DE MATEMÁTICAS, UNIVERSIDAD AUTÓNOMA DE MADRID,
MADRID 28049, SPAIN

E-mail address: adrian.ubis@uam.es