

# AVERAGES OF EULER PRODUCTS, DISTRIBUTION OF SINGULAR SERIES AND THE UBIQUITY OF POISSON DISTRIBUTION

EMMANUEL KOWALSKI

ABSTRACT. We discuss in some detail the general problem of computing averages of convergent Euler products, and apply this to examples arising from singular series for the  $k$ -tuple conjecture and more general problems of polynomial representation of primes. We show that the “singular series” for the  $k$ -tuple conjecture have a limiting distribution when taken over  $k$ -tuples with (distinct) entries of growing size. We also give conditional arguments that would imply that the number of twin primes (or more general polynomial prime patterns) in suitable short intervals are asymptotically Poisson distributed.

## 1. INTRODUCTION

Euler products over primes are ubiquitous in analytic number theory, going back to Euler’s proof that there are infinitely many prime numbers based on the behavior of the zeta function  $\zeta(s)$  as  $s \rightarrow 1$ . As defining  $L$ -functions of various types, Euler products are particularly important, and their properties remain very mysterious. In this paper, we consider the issue of the average or statistical behavior of another important class of Euler products, the so-called *singular series*, arising in counting problems for certain “patterns” of primes (singular series also occur in many problems of additive number theory or diophantine geometry, but we do not consider these here).

The first type of prime patterns are the prime  $k$ -tuples, which are the subject of a famous conjecture of Hardy and Littlewood. Let  $k \geq 1$  be an integer and let  $\mathbf{h} = (h_1, \dots, h_k)$  be a  $k$ -tuple of integers with  $h_i \geq 1$  for all  $i$ . Let then

$$\pi(N; \mathbf{h}) = |\{n \leq N \mid n + h_i \text{ is prime for } 1 \leq i \leq k\}|$$

be the counting function for primes represented by this  $k$ -tuple; note that, for instance,  $\mathbf{h} = (1, 3)$  leads to the function counting twin primes up to  $N$ .

For any prime number  $p$ , let  $\nu_p(\mathbf{h})$  denote the cardinality of the set

$$\{h_1, \dots, h_k\} \pmod{p}$$

of the reductions of the  $h_i$  modulo  $p$ . Note that  $1 \leq \nu_p(\mathbf{h}) \leq \min(k, p)$  for all  $p$ , and that if we assume (as we now do) that the  $h_i$ ’s are distinct, then  $\nu_p(\mathbf{h}) = k$  for all sufficiently large  $p$ .

---

2010 *Mathematics Subject Classification.* 11P32, 11N37, 11K65.

*Key words and phrases.* Singular series, prime  $k$ -tuples conjecture, Bateman-Horn conjecture, limiting distribution, Euler product, moments, Poisson distribution.

The *singular series* associated with  $\mathbf{h}$  is defined as the Euler product

$$(1.1) \quad \mathfrak{S}(\mathbf{h}) = \prod_p \left(1 - \frac{\nu_p(\mathbf{h})}{p}\right) \left(1 - \frac{1}{p}\right)^{-k} = \prod_p \left(1 - \frac{\nu_p(\mathbf{h}) - 1}{p-1}\right) \left(1 - \frac{1}{p}\right)^{1-k}$$

which is absolutely convergent (as will be checked again later; here and throughout the paper, as usual,  $p$  is restricted to prime numbers).

The significance of this value is found in the Hardy-Littlewood prime  $k$ -tuple conjecture (originally stated in [HL]), which states that we should have

$$(1.2) \quad \pi(N; \mathbf{h}) = \mathfrak{S}(\mathbf{h}) \frac{N}{(\log N)^k} (1 + o(1)), \quad \text{as } N \rightarrow +\infty,$$

and in particular, if  $\mathfrak{S}(\mathbf{h}) \neq 0$ , there should be infinitely many integers  $n$  such that  $n + h_1, \dots, n + h_k$  are simultaneously prime. Of course, if  $k \geq 2$ , this is still completely open, but let us mention that from sieve methods, it follows that

$$\pi(N; \mathbf{h}) \leq 2^k k! (1 + o(1)) \mathfrak{S}(\mathbf{h}) \frac{N}{(\log N)^k}$$

as  $N \rightarrow +\infty$  (see, e.g., [IK, Th. 6.7] or [HR, Ch. 4, Th. 5.3]), showing that the singular series does arise naturally. Also some other previously inaccessible additive problems with primes, related to counting arithmetic progressions (of fixed length) of primes are currently being attacked with striking success by B. Green and T. Tao (see [GT]).

More generally, one considers polynomial prime patterns. First, a finite family  $\mathbf{f} = (f_1, \dots, f_m)$  of polynomials in  $\mathbf{Z}[X]$  of degrees  $\deg(f_j) \geq 1$  is said to be *primitive* if the  $f_j$  are distinct, and each  $f_j$  is irreducible, has positive leading coefficient, and the gcd of its coefficients is 1.

If  $\mathbf{f}$  is primitive, we say that an integer  $n \geq 1$  is an  $\mathbf{f}$ -prime seed if  $f_1(n), \dots, f_m(n)$  are all (positive) primes. Then we denote by

$$\pi(N; \mathbf{f}) = |\{n \leq N \mid n \text{ is an } \mathbf{f}\text{-prime seed}\}|$$

for  $N \geq 1$  the counting function for those prime seeds. Moreover, let

$$\text{peg}(\mathbf{f}) = \prod_{j=1}^m \deg(f_j).$$

A generalization of the  $k$ -tuple conjecture, due to Bateman and Horn [BH],<sup>1</sup> states that

$$(1.3) \quad \pi(N; \mathbf{f}) \sim \frac{1}{\text{peg}(\mathbf{f})} \mathfrak{S}(\mathbf{f}) \frac{N}{(\log N)^m}, \quad \text{as } N \rightarrow +\infty,$$

if  $\mathfrak{S}(\mathbf{f}) \neq 0$ , where<sup>2</sup>

$$(1.4) \quad \mathfrak{S}(\mathbf{f}) = \prod_p \left(1 - \frac{\nu_p(\mathbf{f})}{p}\right) \left(1 - \frac{1}{p}\right)^{-m},$$

with  $\nu_p(\mathbf{f})$  being now the number of  $x \in \mathbf{Z}/p\mathbf{Z}$  such that  $f_j(x) = 0$  for some  $j$ ,  $1 \leq j \leq m$ .

<sup>1</sup> The qualitative version of which is due to Schinzel [S].

<sup>2</sup> Here, except in the special case where all  $f_j$  are linear, the singular series  $\mathfrak{S}(\mathbf{f})$  is *not* absolutely convergent (see below for more details on this; the problem is that  $\nu_p(\mathbf{f})$  is only equal to  $m$  on average over  $p$ , and not for all  $p$  large enough, except if each  $f_j$  is linear); the product is thus *defined* as the limit of partial products over primes  $p \leq y$ .

The Hardy-Littlewood conjecture for a  $k$ -tuple  $\mathbf{h}$  is equivalent with this conjecture for the primitive family

$$\mathbf{f} = (X + h_1, \dots, X + h_k)$$

for which  $\nu_p(\mathbf{h})$  as defined previously does coincide with  $\nu_p(\mathbf{f})$ .

Our goal is to study various averages of singular series, for which there is undoubted arithmetic interest. A result of Gallagher [Ga] states that

$$(1.5) \quad \lim_{h \rightarrow +\infty} \frac{1}{h^k} \sum_{|\mathbf{h}| \leq h}^* \mathfrak{S}(\mathbf{h}) = 1,$$

for any fixed  $k$ , as  $h \rightarrow +\infty$ , where  $|\mathbf{h}| = \max h_i$  and  $\sum^*$  restricts to  $k$ -tuples with distinct components. This property was used by Gallagher himself to understand the behavior of primes in short intervals (see also the recent work by Montgomery and Soundararajan [MS]), and it is also important the remarkable results of Goldston, Pintz and Yıldırım concerning small gaps between primes (see [GPY] or the survey [K1]).

Our first question is to ask about finer aspects of the distribution of  $\mathfrak{S}(\mathbf{h})$ . To apply the method of moments, we first prove the following:

**Theorem 1.1.** *Let  $k \geq 1$  be fixed. For any complex number  $m \in \mathbf{C}$  with  $\operatorname{Re}(m) \geq 0$ , there exists a complex number  $\mu_k(m)$  such that*

$$\lim_{h \rightarrow +\infty} \frac{1}{h^k} \sum_{|\mathbf{h}| \leq h}^* \mathfrak{S}(\mathbf{h})^m = \mu_k(m).$$

Moreover, for  $m, k \geq 1$  both integers, we have the symmetry property

$$(1.6) \quad \mu_k(m) = \mu_m(k);$$

in addition, we have  $\mu_1(m) = 1$  for all integers  $m \geq 1$ , and hence  $\mu_k(1) = 1$  for all  $k \geq 1$ .

The last statement ( $\mu_k(1) = 1$ ) is of course Gallagher's theorem (1.5); our proof is not intrinsically different, but maybe more enlightening. These results are in fact quite straightforward, and only the final symmetry in  $k$  and  $m$  is maybe surprising. However, its origin is not particularly mysterious: it is a "local" phenomenon, and it can be guessed from (1.2) by a formal computation.

We will also find estimates for the size of the moments which are good enough to imply the existence of a limiting distribution of  $\mathfrak{S}(\mathbf{h})$  for  $k$ -tuples ( $k$  fixed):

**Theorem 1.2.** *Let  $k \geq 1$  be fixed. There exists a probability law  $\nu_k$  on  $\mathbf{R}^+ = [0, +\infty[$  such that  $\mathfrak{S}(\mathbf{h})$ , for  $\mathbf{h}$  with  $|\mathbf{h}| \leq h$  and  $h \rightarrow +\infty$ , becomes equidistributed with respect to  $\nu_k$ , or equivalently*

$$\lim_{h \rightarrow +\infty} \frac{1}{h^k} \sum_{|\mathbf{h}| \leq h}^* f(\mathfrak{S}(\mathbf{h})) = \int_{\mathbf{R}^+} f(t) d\nu_k(t)$$

for any bounded continuous function on  $\mathbf{R}$ .

The second question we explore is the generalization to other prime patterns of the result of Gallagher (based on (1.5)) that shows that a uniform version of the prime  $k$ -tuple conjecture implies that for a fixed  $\lambda > 0$ , the distribution of

$\pi(x + \lambda \log x) - \pi(x)$  is close to a Poisson distribution of parameter  $\lambda$  as  $x \rightarrow +\infty$ , i.e., it implies that

$$(1.7) \quad \frac{1}{N} |\{n \leq N \mid \pi(n+h) - \pi(n) = m\}| \rightarrow e^{-\lambda} \frac{\lambda^m}{m!}, \quad \text{as } N \rightarrow +\infty,$$

for any integer  $m \geq 0$ . It turns out that, indeed, under a general uniform version of the Bateman-Horn conjecture, for *any* fixed primitive family  $\mathbf{f}$ , the number of  $\mathbf{f}$ -prime seeds in short intervals of “fair” length (i.e., intervals around  $n$  in which (1.3) predicts that, on average, there should be a fixed number of  $\mathbf{f}$ -prime seeds) *always* follows a Poisson distribution. As for the symmetry property of the higher moments for the singular series related to  $k$ -tuple conjecture, this turns out to depend primarily on local identities, but we found this rigidity of patterns to be quite surprising at first sight. Precisely:

**Theorem 1.3.** *Assume that the Bateman-Horn conjecture holds uniformly for all primitive families with non-zero singular series, in the sense that*

$$(1.8) \quad \pi(N; \mathbf{f}) = \frac{1}{\text{peg}(\mathbf{f})} \mathfrak{S}(\mathbf{f}) \frac{N}{(\log N)^m} \left(1 + O\left(\frac{c(\mathbf{f})^\varepsilon}{\log N}\right)\right)$$

*holds for all primitive families  $\mathbf{f}$ , all  $\varepsilon > 0$ , and all  $N \geq 2$ , where*

$$c(\mathbf{f}) = \sum_{1 \leq j \leq m} H(f_j), \quad H(a_0 + a_1 X + \cdots + a_d X^d) = \max_i |a_i|,$$

*and the implied constant depends at most on the degrees of the elements of  $\mathbf{f}$  and on  $\varepsilon$ .*

*Let  $\mathbf{f}$  be a fixed primitive family with  $\mathfrak{S}(\mathbf{f}) \neq 0$ . For  $N \geq 1$ , let*

$$\delta(N, \mathbf{f}) = \frac{\text{peg}(\mathbf{f})}{\mathfrak{S}(\mathbf{f})} (\log N)^m.$$

*Then for any  $\lambda > 0$  and any integer  $r \geq 0$ , we have*

$$\lim_{N \rightarrow +\infty} \frac{1}{N} |\{n \leq N \mid \pi(n + \lambda \delta(N, \mathbf{f}); \mathbf{f}) - \pi(n; \mathbf{f}) = r\}| = e^{-\lambda} \frac{\lambda^r}{r!}.$$

*In other words, for  $N$  large, the number of  $\mathbf{f}$ -prime seeds in an interval around  $N \geq 1$  of length  $\lambda (\log N)^m$  is asymptotically distributed like a Poisson random variable with mean given by  $\mathfrak{S}(\mathbf{f}) \text{peg}(\mathbf{f})^{-1} \lambda$ .*

The final purpose of this paper is to emphasize the fact that Theorems 1.1 and 1.3 are special cases of the problem of computing the average of some families of values of Euler products, and (because here the Euler products are absolutely convergent or almost so) the outcome is consistent with the heuristic that the  $p$ -factors are *independent random variables*, so the average of the Euler product is the product of “local” averages. All this is a fairly common theme in analytic number theory, but our presentation is maybe more systematic than usual. The works of Granville-Soundararajan [GS] and Cogdell-Michel [CM] also present this point of view very successfully for values of certain families of  $L$ -functions at the edge of the critical strip, and Y. Lamzouri [La] has developed this type of ideas in a quite general context. Although this is not really relevant from the point of view of singular series, we just mention that Euler products built of local averages still make sense inside the critical strip for many families of  $L$ -functions, and are closely related to their distribution (as one can see, e.g., from the work of Bohr and Jessen [1] for the Riemann zeta function). On the critical line, “renormalized” Euler products still

occur in the moment conjectures for  $L$ -functions (see, e.g., [KS]), although other factors (conjecturally linked to Random Matrices) also appear.

In the next section, we state in probabilistic terms a general result on averages of random Euler products. Then we use it to prove Theorem 1.1 and Theorem 1.2 in Sections 3 and 4. In Section 5, we prove Theorem 1.3.

**Notation.** As usual,  $|X|$  denotes the cardinality of a set. By  $f \ll g$  for  $x \in X$ , or  $f = O(g)$  for  $x \in X$ , where  $X$  is an arbitrary set on which  $f$  is defined, we mean synonymously that there exists a constant  $C \geq 0$  such that  $|f(x)| \leq Cg(x)$  for all  $x \in X$ . The “implied constant” is any admissible value of  $C$ . It may depend on the set  $X$  which is always specified or clear in context. On the other hand,  $f \sim g$  as  $x \rightarrow x_0$  means  $f/g \rightarrow 1$  as  $x \rightarrow x_0$ .

We use standard probabilistic terminology: a probability space  $(\Omega, \Sigma, \mathbf{P})$  is a triple made of a set  $\Omega$  with a  $\sigma$ -algebra and a measure  $\mathbf{P}$  on  $\Sigma$  with  $\mathbf{P}(\Omega) = 1$ . A random variable is a measurable function  $\Omega \rightarrow \mathbf{R}$  (or  $\Omega \rightarrow \mathbf{C}$ ), and the expectation  $\mathbf{E}(X)$  on  $\Omega$  is the integral of  $X$  with respect to  $\mathbf{P}$  when defined. The law of  $X$  is the measure  $\nu$  on  $\mathbf{R}$  (or  $\mathbf{C}$ ) defined by  $\nu(A) = \mathbf{P}(X \in A)$ . If  $A \subset \Omega$ , then  $\mathbb{1}_A$  is the characteristic function of  $A$ .

For  $k$ -tuples  $\mathbf{h} = (h_1, \dots, h_k)$ , we recall that  $|\mathbf{h}| = \max(|h_i|)$ . When different values of  $k$  can occur, we sometimes write  $|\mathbf{h}|_k$  to indicate the number of components of  $\mathbf{h}$ , in particular a sum such as

$$\sum_{|\mathbf{h}|_k \leq h} a(\mathbf{h})$$

is a sum over  $k$ -tuples (of positive integers) with components  $\leq h$ .

## 2. A PROBABILISTIC STATEMENT

We assume given a probability space  $(\Omega, \Sigma, \mathbf{P})$ , and two sequences of random variables

$$X_p, Y_p : \Omega \rightarrow \mathbf{C}$$

which are indexed by prime numbers.

We assume that  $(Y_p)$  is an independent sequence; recall that this means that

$$\mathbf{P}(Y_{p_1} \in A_1, \dots, Y_{p_k} \in A_k) = \prod_{1 \leq i \leq k} \mathbf{P}(Y_{p_i} \in A_i)$$

for all choices of finitely many distinct primes  $p_1, \dots, p_k$ , and all measurable sets  $A_i \subset \mathbf{C}$ , and that a consequence is that (when the expectation makes sense), we have

$$\mathbf{E}(Y_{p_1} \cdots Y_{p_k}) = \mathbf{E}(Y_{p_1}) \cdots \mathbf{E}(Y_{p_k}).$$

We now extend the family to all integers by denoting

$$X_q = \prod_{p|q} X_p, \quad Y_q = \prod_{p|q} Y_p,$$

for any squarefree integer  $q \geq 1$ , and  $X_q = Y_q = 0$  if  $q \geq 1$  is not squarefree.

We will consider the behavior of the random Euler products

$$Z_X = \prod_p (1 + X_p), \quad Z_Y = \prod_p (1 + Y_p)$$

and in particular their expectations  $\mathbf{E}(Z_X)$  and  $\mathbf{E}(Z_Y)$ .

For this purpose, we assume that the products converge absolutely (almost surely). More precisely, expand formally

$$\prod_p (1 + X_p) = \sum_{q \geq 1}^b X_q,$$

where  $\sum^b$  restricts the sum to squarefree numbers. Then we assume that

$$(2.1) \quad \sum_{q > x}^b |X_q| \leq R_X(x)$$

where  $R_X(x)$  is an integrable non-negative random variable such that  $R_X(x) \rightarrow 0$  almost surely as  $x \rightarrow +\infty$ . It then follows that  $Z_X$  is almost surely an absolutely convergent infinite product.

We moreover assume that the product

$$(2.2) \quad \prod_p (1 + |\mathbf{E}(Y_p)|)$$

converges (absolutely). By independence of the  $(Y_p)$ , we know that

$$|\mathbf{E}(Y_q)| = \left| \mathbf{E} \left( \prod_{p|q} Y_p \right) \right| = \prod_{p|q} |\mathbf{E}(Y_p)|$$

and so expanding again in series, we obtain that

$$(2.3) \quad \sum_{q \geq 1}^b |\mathbf{E}(Y_q)| = \sum_{q \geq 1}^b \prod_{p|q} |\mathbf{E}(Y_p)| = \prod_p (1 + |\mathbf{E}(Y_p)|) < +\infty.$$

Our goal is to show that if  $(X_p)$  is distributed “more or less” like  $(Y_p)$ , but without being independent, the expectation of  $Z_X$  is close to

$$\prod_p (1 + \mathbf{E}(Y_p)).$$

In particular, we will typically have  $(X_p)$  depend on another parameter (say  $h$ ), in such a way that  $X_{p,h}$  converges in law to  $Y_p$  (which will remain fixed) when  $h \rightarrow +\infty$ , and this will lead to the relation

$$\lim_{h \rightarrow +\infty} \mathbf{E} \left( \prod_p (1 + X_{p,h}) \right) = \prod_p (1 + \mathbf{E}(Y_p))$$

in a number of situations. We interpret this as saying that (when applicable) the average of the Euler product  $Z_X$  is obtained “as if” the factors were independent, and taking the product of the local averages  $1 + \mathbf{E}(Y_p)$  of the “model” random variables defining  $Z_Y$ .

Here is the precise (and almost tautological) “finitary” statement from which applications will be derived.

**Proposition 2.1.** *Let  $(X_p), (Y_p)$  be as above. Then for any choice of the auxiliary parameter  $x > 0$ , we have*

$$\mathbf{E}(Z_X) = \prod_p (1 + \mathbf{E}(Y_p)) + O \left( \mathbf{E}(R_X(x)) + \sum_{q \leq x}^b |\mathbf{E}(X_q - Y_q)| + \sum_{q > x}^b |\mathbf{E}(Y_q)| \right),$$

where the implied constant is absolute, and in fact has modulus at most 1.

*Proof.* This more or less proves itself: for any  $x \geq 1$ , write first

$$\prod_p (1 + X_p) = \sum_{q \geq 1}^b X_q = \sum_{q \leq x}^b X_q + \sum_{q > x}^b X_q,$$

then use (2.1) to estimate the second term, and take the expectation, which leads to

$$\mathbf{E}(Z_X) = \sum_{q \leq x} \mathbf{E}(X_q) + O(\mathbf{E}(R_X(x))).$$

Next, we insert  $Y_q$  by writing  $X_q = Y_q + (X_q - Y_q)$ , getting

$$\mathbf{E}(Z_X) = \sum_{q \leq x}^b \mathbf{E}(Y_q) + \sum_{q \leq x}^b \mathbf{E}(X_q - Y_q) + O(\mathbf{E}(R_X(x)))$$

and then use

$$\sum_{q \leq x}^b \mathbf{E}(Y_q) = \sum_{q \geq 1}^b \mathbf{E}(Y_q) + O\left(\sum_{q > x}^b |\mathbf{E}(Y_q)|\right) = \prod_p (1 + \mathbf{E}(Y_p)) + O\left(\sum_{q > x}^b |\mathbf{E}(Y_q)|\right),$$

to conclude the proof.  $\square$

*Remark 2.2.* Observe that by (2.3), the last term in the remainder tends to zero as  $x \rightarrow +\infty$ . Moreover, if  $R_X(x)$  is dominated by an integrable function as  $x \rightarrow +\infty$ , the assumption that  $R_X(x) \rightarrow 0$  almost surely implies that the first term also tends to zero. Thus to conclude in practical applications, one needs to control the middle term.

In terms of the “extra” parameter  $h$  mentioned before the statement of the proposition, we may typically hope for uniform estimates for  $\mathbf{E}(R_X(x))$ , in terms of  $h$ , say

$$\mathbf{E}(R_X(x)) \ll h^\alpha x^{-\beta}, \quad \alpha, \beta > 0;$$

if we also have a bound of the type

$$(2.4) \quad \mathbf{E}(X_q) = \mathbf{E}(Y_q) + O(q^\gamma h^{-\delta}), \quad \gamma, \delta > 0,$$

(or if this holds on average over  $q < x$ , which may often be easier to prove, as is the case for the error term in the prime number theorem, as shows the Bombieri-Vinogradov theorem), this leads to a remainder term which is

$$\ll h^\alpha x^{-\beta} + x^{1+\gamma} h^{-\delta} + \varepsilon(x)$$

with  $\varepsilon(x) \rightarrow 0$  as  $x \rightarrow +\infty$ , uniformly in  $h$ . Then we can conclude that

$$(2.5) \quad \lim_{h \rightarrow +\infty} \mathbf{E}(Z_X) = \prod_p (1 + \mathbf{E}(Y_p))$$

by choosing  $x$  suitably as a function of  $h$ , *provided* we have

$$\frac{\alpha}{\beta} < \frac{\delta}{\gamma + 1}.$$

We will see this in action concretely in the next sections. Notice that if  $\alpha$  can be chosen arbitrarily small (i.e.,  $R_X(x)$  is bounded almost uniformly in terms of  $h$ ), then this condition can be met.

*Remark 2.3.* If we assume, instead of (2.2), that the product of  $1 + \mathbf{E}(|Y_p|)$  converges, which is stronger, it follows that  $\sum |Y_p| < +\infty$  almost surely (its expectation being finite), and hence the infinite product defining  $Z_Y$  converges absolutely almost surely. Also, since we have

$$\mathbf{E}\left(\prod_{p \leq P} (1 + Y_p)\right) = \prod_{p \leq P} (1 + \mathbf{E}(Y_p))$$

for all  $P$ , we would obtain

$$\mathbf{E}(Z_Y) = \prod_p (1 + \mathbf{E}(Y_p)).$$

provided  $Z_Y$  converges dominatedly, for instance. This formula is also valid if  $Y_p \geq 0$ , by the monotone convergence theorem. It provides an interpretation of the right-hand side of (2.5).

### 3. MOMENTS OF SINGULAR SERIES FOR THE $k$ -TUPLE CONJECTURE

In this section, we prove Theorem 1.1, which includes in particular Gallagher's theorem, in a way which may seem somewhat complicated but which clarifies the result.

We first assume an integer  $k \geq 1$  to be fixed. We rewrite (1.1) as

$$\mathfrak{S}(\mathbf{h}) = \prod_p \left(1 + \frac{p^k - \nu_p(\mathbf{h})p^{k-1} - (p-1)^k}{(p-1)^k}\right).$$

It is therefore natural to define

$$a(p, \nu) = \frac{p^k - \nu p^{k-1} - (p-1)^k}{(p-1)^k}$$

for all primes  $p$  and real numbers  $\nu$ ,  $0 < \nu \leq p$  (omitting the dependency on  $k$ ). We then define  $a_m(p, \nu)$ , for  $m \in \mathbf{C}$  with  $\operatorname{Re}(m) \geq 0$ , by requiring that

$$1 + a_m(p, \nu) = (1 + a(p, \nu))^m,$$

with the convention  $0^m = 0$  if  $\operatorname{Re}(m) = 0$ ; the condition  $\nu \leq p$  implies that  $1 + a(p, \nu) \geq 0$ , so this is well-defined indeed. (If we assume  $\nu < p$ , we may extend this to all  $m \in \mathbf{C}$ ).

We first need a technical lemma.

**Lemma 3.1.** *For  $m \in \mathbf{C}$  with  $\operatorname{Re}(m) \geq 0$ , write  $m^+ = 0$  if  $\operatorname{Re}(m) < 1$ , and  $m^+ = m - 1$  otherwise. For all  $p$  prime and  $\nu$  with  $1 \leq \nu \leq \min(p, k)$ , we have*

$$(3.1) \quad a_m(p, k) \ll \frac{|m|}{p^2} \left(1 + O\left(\frac{1}{p^2}\right)\right)^{m^+},$$

$$(3.2) \quad a_m(p, \nu) \ll \frac{|m|}{p} \left(1 + O\left(\frac{1}{p}\right)\right)^{m^+}, \quad \text{if } 1 \leq \nu < k,$$

where the implied constants depend only on  $k$ .

*Proof.* Notice first that, in the stated range, we have

$$\begin{aligned} a(p, k) &\ll p^{-2}, \\ a(p, \nu) &\ll p^{-1}, \quad \text{if } 1 \leq \nu < k, \end{aligned}$$

where the implied constants depend only on  $k$ , and then write

$$a_m(p, \nu) = (1 + a(p, \nu))^m - 1 = ma(p, \nu) \int_0^1 (1 + ta(p, \nu))^{m-1} dt$$

and estimate directly.  $\square$

We are now going to prove Theorem 1.1. Fix  $h \geq 1$  (though  $h$  will tend to infinity at the end). We first interpret the  $m$ -th moment of the singular series in probabilistic terms, then introduce the source of its limiting value in the framework of the previous section.

Consider the finite set (again, depending on  $k$ )

$$\Omega_1 = \{\mathbf{h} = (h_i) \mid 1 \leq h_i \leq h, h_i \text{ distinct}\},$$

with the normalized counting measure. Denoting  $h_k^* = |\Omega_1|$ , notice that

$$(3.3) \quad h_k^* = h^k(1 + O(h^{-1}))$$

for  $h \geq 1$ , the implied constant depending only on  $k$ . We will denote by  $\mathbf{E}_1$  and  $\mathbf{P}_1$  the expectation and probability for this discrete space. So we have, for instance, that

$$\mathbf{P}_1(\nu_p = \nu) = \frac{1}{h_k^*} |\{\mathbf{h} \in \Omega_1 \mid \nu_p(\mathbf{h}) = \nu\}|.$$

Our goal is to find the limit as  $h \rightarrow +\infty$  of the average

$$\frac{1}{h_k^*} \sum_{\substack{|\mathbf{h}| \leq h \\ h_i \text{ distinct}}} \mathfrak{S}(\mathbf{h})^m = \mathbf{E}_1(\mathfrak{S}(\mathbf{h})^m)$$

(notice that, by (3.3), if the limit exists, it is also the limit of

$$\frac{1}{h^k} \sum_{\substack{|\mathbf{h}| \leq h \\ h_i \text{ distinct}}} \mathfrak{S}(\mathbf{h})^m,$$

as  $h \rightarrow +\infty$ ).

We write  $X_p(\mathbf{h}) = a(p, \nu_p(\mathbf{h}))$  and  $X_p(m, \mathbf{h}) = a_m(p, \nu_p(\mathbf{h}))$ , so that

$$\prod_p (1 + X_p(m, \mathbf{h})) = \mathfrak{S}(\mathbf{h})^m$$

by construction.

Now consider a second space

$$\Omega_2 = \prod_p (\mathbf{Z}/p\mathbf{Z})^k$$

with the product measure of the probability counting measures on each factor. We denote by  $\omega = (\mathbf{h}_p)_p$  the elements of  $\Omega_2$ . To avoid confusion with  $\nu_p$  defined for  $\mathbf{h} \in \Omega_1$ , we introduce the random variables

$$\rho_p : \begin{cases} \Omega_2 \rightarrow \{1, \dots, k\} \\ \omega = (\mathbf{h}_p)_p \mapsto \text{number of distinct } h_i \text{ in } \mathbf{Z}/p\mathbf{Z}, \end{cases}$$

which satisfy  $1 \leq \rho_p \leq \min(k, p)$ .

We can now define “random” singular series using  $\Omega_2$ , writing  $Y_p = a(p, \rho_p)$  and considering the Euler product

$$\prod_p (1 + Y_p),$$

and similarly with  $Y_p(m) = a_m(p, \rho_p)$  and

$$\prod_p (1 + Y_p(m)) = \left( \prod_p (1 + Y_p) \right)^m.$$

We denote by  $\mathbf{P}_2$  and  $\mathbf{E}_2$  the probability and expectation for this space. By construction of  $\Omega_2$ , the random variables  $(\rho_p)$  are independent, and so are the  $(Y_p)$ , and the  $(Y_p(m))$  for a given  $m$ . Note also that the components  $\mathbf{h}_p$  are equidistributed: for any prime  $p$  and any  $a \in (\mathbf{Z}/p\mathbf{Z})^k$ , we have

$$(3.4) \quad \mathbf{P}_2(\mathbf{h}_p = a) = \frac{1}{p^k}.$$

We now use Proposition 2.1 to compare the average  $\mathbf{E}_1(\mathfrak{S}(\mathbf{h})^m)$  with

$$\prod \mathbf{E}_2((1 + Y_p)^m).$$

Although this proposition is phrased with a single probability space  $\Omega$  on which both Euler vectors are defined, this is not a serious issue and the statement remains valid, provided the expectations are suitably subscripted and one writes

$$\left| \mathbf{E}_1(X_q(m)) - \mathbf{E}_2(Y_q(m)) \right|$$

on the right-hand side instead of  $|\mathbf{E}(X_q(m) - Y_q(m))|$ .<sup>3</sup>

We start by estimating the tail  $R(x) = R_{X(m)}(x)$  of the Euler product defining  $\mathfrak{S}(\mathbf{h})^m$ . In keeping with probabilistic conventions, we omit the argument  $\mathbf{h} \in \Omega_1$  in many places. Denoting

$$\Delta(\mathbf{h}) = \left| \prod_{i < j} (h_i - h_j) \right| \geq 1,$$

and noting that  $\nu_p = k$  unless  $p \mid \Delta$ , we have from Lemma 3.1 the bound

$$|X_p(m)| \ll |m| \left( 1 + O\left(\frac{(p, \Delta)}{p^2}\right) \right)^{m^+} (p, \Delta) p^{-2}$$

for some  $C > 0$  (depending only on  $k$ ) and all  $\mathbf{h}$ ,  $m$  (with  $\operatorname{Re}(m) \geq 0$ ) and  $p$ , the implied constant depending only on  $k$  (this justifies, in particular, the convergence of the Euler product  $Z_X$  for every  $\mathbf{h}$ ). Hence, taking the product over  $p \mid q$  for a squarefree integer  $q$ , we get

$$|X_q(m)| \leq (|m|B)^{\omega(q)} (q, \Delta) q^{-2} \prod_{p \mid q} \left( 1 + C \frac{(p, \Delta)}{p^2} \right)^{m^+}$$

<sup>3</sup> We could also simply consider  $\Omega = \Omega_1 \times \Omega_2$  with the product measure, or equivalently (and maybe more elegantly) assume that we start with some space  $\Omega$  and two vectors  $(X_p)$ ,  $(Y_p)$ , distributed according to the prescription of  $\Omega_1$  and  $\Omega_2$  respectively, i.e., with

$$\begin{aligned} \mathbf{P}(X_p = a) &= \frac{1}{h_k^*} |\{\mathbf{h} \in \Omega_1 \mid a(p, \nu_p(\mathbf{h})) = a\}|, \\ \mathbf{P}(Y_p = a) &= \frac{1}{p^k} |\{\mathbf{h} \in (\mathbf{Z}/p\mathbf{Z})^k \mid a(p, \rho_p(\mathbf{h})) = a\}|. \end{aligned}$$

for some constants  $B > 0$  and  $C \geq 0$  depending only on  $k$ . Since  $\Delta$  is bounded by

$$(3.5) \quad |\Delta| \leq (2h)^{k^2},$$

a standard computation with sums of multiplicative functions leads to

$$\sum_{q>x}^b |X_q(m)| \ll x^{-1} (\log 2hx)^D$$

for  $x \geq 2$  and some constant  $D \geq 0$ , depending on  $k$  and  $m$ .

The next step is to justify the analogue of the convergence of (2.2); more precisely, we have

$$(3.6) \quad \prod_p (1 + \mathbf{E}_2(|Y_p(m)|)) < +\infty.$$

Indeed, Lemma 3.1 leads to

$$\mathbf{E}_2(|Y_p(m)|) \ll p^{-2} + p^{-1} \mathbf{P}_2(\rho_p < k) \ll p^{-2}$$

for  $p \geq 2$ , where the implied constant depends on  $k$  and  $m$ , since it is clear that we have

$$(3.7) \quad \mathbf{P}_2(\rho_p < k) \leq \frac{k(k-1)}{2p}$$

for all primes  $p$  and  $k \geq 1$  (write that the event  $\{\rho_p < k\}$  is the union – not necessarily disjoint – of the  $k(k-1)/2$  events  $h_i = h_j$  with  $i \neq j$ , each of which has probability  $1/p$  by uniform distribution (3.4)). By independence, we then also get

$$(3.8) \quad \mathbf{E}_2(|Y_q(m)|) \leq A^{\omega(q)} q^{-2}.$$

for all squarefree integers  $q$  and some constant  $A \geq 1$ , which depends only on  $k$  and  $m$ .

Finally, it remains to estimate  $\mathbf{E}_1(X_q(m)) - \mathbf{E}_2(Y_q(m))$ . We claim that, for any  $a \in \mathbf{C}$ , we have

$$(3.9) \quad \mathbf{P}_1(X_q(m) = a) = \left(1 + O\left(\frac{q}{h}\right)\right) \mathbf{P}_2(Y_q(m) = a) + O\left(\frac{k^{\omega(q)}}{h}\right)$$

where the implied constants depend only on  $k$ . Assuming this, and noting that  $X_q(m)$  and  $Y_q(m)$  take the same finitely many values (at most  $k^{\omega(q)}$  distinct values, which are

$$\ll \frac{F^{\omega(g)}}{q}$$

where the implied constant and  $F$  depend on  $m$  and  $k$ ), it follows that

$$\mathbf{E}_1(X_q(m)) = \left(1 + O\left(\frac{q}{h}\right)\right) \mathbf{E}_2(Y_q(m)) + O\left(\frac{G^{\omega(q)}}{h}\right),$$

where  $G$  depends on  $m$  and  $k$ , leading in turn to

$$\left| \mathbf{E}_1(X_q(m)) - \mathbf{E}_2(Y_q(m)) \right| \ll \frac{q}{h} \mathbf{E}_2(|Y_q(m)|) + \frac{G^{\omega(q)}}{h} \ll \frac{E^{\omega(q)}}{h}$$

(see (3.8)), where the implied constant depends only on  $k$  and  $m$ , as does  $E$ .

Summing over  $q < x$ , it then follows from Proposition 2.1 that

$$\frac{1}{h_k^*} \sum_{\mathbf{h}}^* \mathfrak{S}(\mathbf{h})^m = \mathbf{E}_1 \left( \prod_p (1 + X_p(m)) \right) = \prod_p (1 + \mathbf{E}_2(Y_p(m))) + O\left(xh^{-1}(\log 2hx)^B + x^{-1}(\log 2hx)^D\right)$$

for some  $B$  depending on  $k$  and  $m$ . Choosing for instance  $x = h^{1/2}$  leads to the existence of the  $m$ -th moment of singular series, with limiting value given by

$$(3.10) \quad \mu_k(m) = \prod_p (1 + \mathbf{E}_2(Y_p(m))) = \prod_p \left(1 - \frac{1}{p}\right)^{-km} \left\{ \frac{1}{p^k} \sum_{\mathbf{h} \in (\mathbf{Z}/p\mathbf{Z})^k} \left(1 - \frac{\rho_p(\mathbf{h})}{p}\right)^m \right\}.$$

It only remains to prove (3.9). Note that this is clearly an expression of quantitative equidistribution (or convergence in law) of  $X_q$  to  $Y_q$  as  $h \rightarrow +\infty$ .<sup>4</sup>

The proof is quite simple. First of all, given arbitrary integers  $s_p$  with  $p \mid q$ , we have

$$\begin{aligned} \mathbf{P}_1(\nu_p(\mathbf{h}) = s_p \text{ for } p \mid q) &= \frac{1}{h_k^*} \sum_{\substack{\nu_p(\mathbf{h})=s_p \text{ for } p \mid q \\ |\mathbf{h}| \leq h}}^* 1 \\ &= \frac{1}{h_k^*} \sum_{\substack{\rho_p(\mathbf{h}_p)=s_p \\ \mathbf{h}_p \in (\mathbf{Z}/p\mathbf{Z})^k}} \cdots \sum_{\substack{|\mathbf{h}| \leq h \\ \mathbf{h} \equiv \mathbf{h}_p \pmod{p|q}}}^* 1 \end{aligned}$$

(where there are as many outer sums in the last line as there are primes dividing  $q$ , and the last sum involves summation conditions for all  $p \mid q$ ). This inner sum is

$$(3.11) \quad \sum_{\substack{|\mathbf{h}| \leq h \\ \mathbf{h} \equiv \mathbf{h}_p \pmod{p|q}}}^* 1 = \sum_{\substack{|\mathbf{h}| \leq h \\ \mathbf{h} \equiv \mathbf{h}_p \pmod{p|q}}} 1 + O(h^{k-1})$$

where the implied constant depends on  $k$  (i.e., we now forget the condition on  $\mathbf{h}$  to have distinct components). Lattice-point counting leads to

$$\sum_{\substack{|\mathbf{h}| \leq h \\ \mathbf{h} \equiv \mathbf{h}_p \pmod{p|q}}} 1 = \frac{h^k}{q^k} \left(1 + O\left(\frac{q}{h}\right)\right)$$

where the implied constant depends again only on  $k$ . In view of the equidistribution of  $\mathbf{h}_p$  for  $(\mathbf{h}_p)_p \in \Omega_2$ , we therefore derive from the above the following *quantitative equidistribution* result:

$$(3.12) \quad \mathbf{P}_1(\nu_p(\mathbf{h}) = s_p \text{ for } p \mid q) = \mathbf{P}_2(\rho_p(\mathbf{h}_p) = s_p \text{ for } p \mid q) \left(1 + O\left(\frac{q}{h}\right)\right) + O\left(\frac{1}{h}\right).$$

Now to derive (3.9), we need only observe that  $Y_q(m)$  and  $X_q(m)$  are “identical” functions of  $\rho_p$  and  $\nu_p$  respectively (for  $p \mid q$ ). Hence (3.12) implies (3.9) by summing over all possible values of  $(s_p)_{p|q}$  leading to a given  $a$ , using the fact that there are at most  $k^{\omega(q)}$  such values (the latter being a very rough estimate!).

It remains to prove the symmetry property (1.6) to finish the proof of Theorem 1.1. We note in advance that since  $\mathfrak{S}(\mathbf{h}) = 1$  for all 1-tuple  $\mathbf{h}$ , we have

<sup>4</sup> It can also be interpreted as a form of “sieve axiom”.

$\mu_1(m) = 1$  for all  $m \geq 1$ , and hence  $\mu_k(1) = 1$  for all  $k \geq 1$ , which is Gallagher's result (1.5).

The symmetry turns out to be true "locally", i.e., the  $p$ -factor of the Euler products (3.10) defining  $\mu_k(m)$  and  $\mu_m(k)$  coincide for all  $p$  and integers  $k, m \geq 1$ .

There are different ways to see this, and the following seems to encapsulate the origin of the phenomenon. Given a finite set  $F$  (which will be  $\mathbf{Z}/p\mathbf{Z}$ ), consider the following obviously symmetric expression of  $m$  and  $k$ :

$$\frac{1}{|F|^{m+k}} \sum_{\substack{\mathbf{x} \in F^m, \\ \{x_i\} \cap \{h_j\} = \emptyset}} \sum_{\mathbf{h} \in F^k} 1$$

(which is the probability, for the normalized counting measure on  $F^{k+m}$ , that a pair of a  $k$ -tuple and an  $m$ -tuple, both of elements of  $F$ , do not contain a common element). Then it can be interpreted either as

$$\begin{aligned} \frac{1}{|F|^m} \sum_{\tau=1}^m \sum_{\substack{\mathbf{x} \in F^m \\ \rho(\mathbf{x})=\tau}} \frac{1}{|F|^k} \sum_{\substack{\mathbf{h} \in F^k \\ \{h_j\} \cap \{x_i\} = \emptyset}} 1 &= \frac{1}{|F|^m} \sum_{\tau=1}^m \sum_{\substack{\mathbf{x} \in F^m \\ \rho(\mathbf{x})=\tau}} \left(1 - \frac{\tau}{|F|}\right)^k \\ &= \frac{1}{|F|^m} \sum_{\mathbf{x} \in F^m} \left(1 - \frac{\rho(\mathbf{x})}{|F|}\right)^k \end{aligned}$$

or (by the same computation with  $m$  and  $k$  reversed) as

$$\frac{1}{|F|^k} \sum_{\mathbf{h} \in F^k} \left(1 - \frac{\rho(\mathbf{h})}{|F|}\right)^m,$$

(using  $\rho(\cdot)$  to denote the number of distinct elements in  $F$  of an  $m$ -tuple, then of a  $k$ -tuple).

Applied with  $F = \mathbf{Z}/p\mathbf{Z}$ , up to the symmetric factor  $(1 - 1/p)^{-km}$  in (3.10), the first is the  $p$ -factor for  $\mu_m(k)$ , and the second is the  $p$ -factor for  $\mu_k(m)$ , showing that they are indeed equal.

*Remark 3.2.* Quantitatively, we have proved that

$$\sum_{|\mathbf{h}| \leq h}^* \mathfrak{S}(\mathbf{h})^m = \mu_k(m) h_k^* + O(h^{k-1/2+\varepsilon}),$$

for any  $\varepsilon > 0$ , where the implied constant depends on  $k$  and  $m$ . For  $m = 1$ , Montgomery and Soundararajan [MS, (17), p. 593] have obtained a more refined expansion with contributions of size  $h^{k-1} \log h$  and  $h^{k-1}$ , and error term of size  $h^{k-3/2+\varepsilon}$ .

*Remark 3.3.* The fact that  $\mu_k(1) = 1$  can be used to recover the combinatorial identities used by Gallagher [Ga, p. 7–8] instead of the probabilistic phrasing above. We review this for completeness: in order to prove  $\mu_k(1) = 1$ , it suffices to show that the average of  $a(p, \rho_p)$  is zero. We have

$$\sum_{\mathbf{h} \in (\mathbf{Z}/p\mathbf{Z})^k} a(p, \rho_p(\mathbf{h})) = \sum_{\nu=1}^p a(p, \nu) |\{\mathbf{h} \in (\mathbf{Z}/p\mathbf{Z})^k \mid \rho_p(\mathbf{h}) = \nu\}|$$

and on the other hand, we have

$$|\{\mathbf{h} \in (\mathbf{Z}/p\mathbf{Z})^k \mid \rho_p(\mathbf{h}) = \nu\}| = \binom{p}{\nu} \left\{ \begin{matrix} k \\ \nu \end{matrix} \right\},$$

where  $\left\{ \begin{matrix} k \\ \nu \end{matrix} \right\}$  is the number of surjective maps from a set with  $k$  elements to one with  $\nu$  elements<sup>5</sup>; indeed, a  $k$ -tuple  $\mathbf{h}$  with  $\nu$  distinct values is the same as a map  $\{1, \dots, k\} \rightarrow \mathbf{Z}/p\mathbf{Z}$  with image of cardinality  $\nu$ , i.e., the set of such tuples is the disjoint union of those sets of surjective maps

$$\{1, \dots, k\} \rightarrow I$$

over  $I \subset \mathbf{Z}/p\mathbf{Z}$  with order  $\nu$ .

Therefore, Gallagher's result follows from the identity

$$\sum_{\nu=1}^p a(p, \nu) \binom{p}{\nu} \left\{ \begin{matrix} k \\ \nu \end{matrix} \right\} = 0$$

which is proved in [Ga, p. 7], and which we have therefore reproved. Similarly, the identities

$$\sum_{\nu=1}^p \binom{p}{\nu} \left\{ \begin{matrix} k \\ \nu \end{matrix} \right\} = p^k, \quad \sum_{\nu=1}^p \nu \binom{p}{\nu} \left\{ \begin{matrix} k \\ \nu \end{matrix} \right\} = p^{k+1} - p(p-1)^k,$$

of [Ga, p. 8] can be derived from the proof that the  $p$ -factor for  $\mu_k(1)$  is 1.

*Remark 3.4.* From (1.2), one can guess that  $\mu_k(m) = \mu_m(k)$  for  $m \geq 1$  integer, by computing

$$\sum_{|\mathbf{h}| \leq h} \left( \sum_{n \leq N} \prod_{1 \leq i \leq k} \Lambda(n + h_i) \right)^m = \sum_{|\mathbf{h}|_k \leq h} \sum_{|\mathbf{n}|_m \leq N} \prod_{\substack{1 \leq i \leq k \\ 1 \leq j \leq m}} \Lambda(n_j + h_i)$$

(where  $\mathbf{n}$  is an  $m$ -tuple), which is a symmetric expression in  $\mathbf{n}$  and  $\mathbf{h}$ , except for the ranges of summation, and which should be asymptotic to either  $\mu_k(m)h^k N^m$  or  $\mu_m(k)h^k N^m$  by a uniform  $k$ -tuple conjecture. In fact, the computation we did amounts to doing the same argument locally (i.e., looking on average over  $\mathbf{h}$  at the distribution of integers such that, for a fixed prime  $p$ ,  $n + h_1, \dots, n + h_k$  are not divisible by  $p$ ).

This symmetry  $\mu_k(m) = \mu_m(k)$ , despite the simplicity of its proof, is a very strong property, as pointed out to us by A. Nikeghbali. Indeed, write  $X_k = Z_{Y,k}$ , the random variable given by the random singular series. Since we have

$$\mu_k(m) = \int_{\mathbf{R}^+} t^m d\nu_k(t) = E(X_k^m),$$

the symmetry implies that the sequence  $(E(X_k^m))_k$ , for a *fixed* value of  $m$ , is the sequence of moments of a probability distribution of  $[0, +\infty[$ , which is a highly non-trivial property. We refer to the survey [Si] of the classical theory surrounding the ‘‘moment problems’’, noting that from Theorem 1 of loc. cit. it follows that, for any fixed  $m \geq 1$ , we have

$$\sum_{\substack{0 \leq i \leq N \\ 0 \leq j \leq N}} \alpha_i \bar{\alpha}_j \mu_{i+j}(m) > 0, \quad \sum_{\substack{0 \leq i \leq N \\ 0 \leq j \leq N}} \alpha_i \bar{\alpha}_j \mu_{i+j+1}(m) > 0,$$

<sup>5</sup> This is denoted  $\sigma(k, \nu)$  in [Ga], and it is *not* the standard notation, which would write  $r! \left\{ \begin{matrix} k \\ r \end{matrix} \right\}$  instead.

for any  $N \geq 1$  and any complex numbers  $(\alpha_i) \in \mathbf{C}^N - \{0\}$ .

It would be quite interesting to know what other types of natural sequences of random variables (or probability distributions) satisfy the relation  $E(X_k^m) = E(X_m^k)$ . One fairly general construction is as follows (this was pointed out by A. Nikeghbali and P. Bourgade): just take  $X_n = Z^n$  for  $Z$  a random variable such that all moments of  $Z$  exist, or a bit more generally, take a sequence  $(X_n)$  of positive random variables such that the  $X_n^{1/n}$  are identically distributed. But note that the variables we encountered are not of this type.

**Example 3.5.** Let  $m = 2$ . We find (using the symmetry property) that the mean-square of  $\mathfrak{S}(\mathbf{h})$  is given by

$$\lim_{h \rightarrow +\infty} \frac{1}{h^k} \sum_{|\mathbf{h}| \leq h}^* \mathfrak{S}(\mathbf{h})^2 = \mu_k(2),$$

where

$$\mu_k(2) = \prod_p \left( \left( \left(1 - \frac{1}{p}\right) \left(1 - \frac{2}{p}\right)^k + \frac{1}{p} \left(1 - \frac{1}{p}\right)^k \right) \left(1 - \frac{1}{p}\right)^{-2k} \right).$$

In particular, we find (using `Pari/GP` for instance):

$$\begin{aligned} \mu_2(2) &= 2.300\dots & \mu_3(2) &= 6.03294\dots \\ \mu_4(2) &= 17.562\dots & \mu_5(2) &= 55.255\dots \\ \mu_6(2) &= 184.18\dots \end{aligned}$$

Note that the second (and higher) moments increase quickly with  $k$  (as proved in Proposition 4.1 in the next section). This is explained intuitively by the fact that  $\mathfrak{S}(\mathbf{h})$  is often zero: for instance, the 2-factor of  $\mathfrak{S}(\mathbf{h})$  is zero unless all  $h_i$  are of the same parity, which happens with probability  $2^{1-k}$  only (see Example 4.3 for a more precise estimate). For those, of course, the 2-factor is very large (equal to  $2^{k-1}$ ).

#### 4. GROWTH AND DISTRIBUTION OF MOMENTS OF SINGULAR SERIES

In this section, we will prove Theorem 1.2, using the methods of moments. For this, we consider the problem (which has independent interest) of determining the growth of  $\mu_k(m)$ . We look at the dependency on  $m$  for fixed  $k$ , or equivalently the dependency on  $k$  for fixed  $m$ , by symmetry (as in Example 3.5). The result is that the moments grow just a bit faster than exponentially.

**Proposition 4.1.** *For any fixed  $k \geq 1$ , we have*

$$\log \mu_k(m) = km \log \log 3m + O(m), \quad \text{for } m \geq 1,$$

where the implied constant depends on  $k$ .

*Proof.* We use the formula (3.10), written in the form

$$\mu_k(m) = \prod_p \left(1 - \frac{1}{p}\right)^{-km} \mathbf{E}_2 \left( \left(1 - \frac{\rho_p}{p}\right)^m \right).$$

We will prove first that

$$\log \mu_k(m) \geq km \log \log 3m + O(m),$$

for  $m \geq 1$ , with an implied constant depending on  $k$ , before proving the corresponding upper bound.

We start by checking that all terms in the Euler product are  $\geq 1$ , i.e., for all primes  $p$ , all integers  $k$  and all real numbers  $m \geq 1$ , we have

$$(4.1) \quad \mathbf{E}_2 \left( \left( 1 - \frac{\rho_p}{p} \right)^m \right) \geq \left( 1 - \frac{1}{p} \right)^{mk}.$$

Indeed, by the symmetry between the  $p$ -factor for  $\mu_k(1)$  and for  $\mu_1(k)$ , we have

$$\left( 1 - \frac{1}{p} \right)^k = \mathbf{E}_2 \left( 1 - \frac{\rho_p}{p} \right),$$

while raising to the  $m$ -th power and applying Hölder's inequality gives

$$\left( \mathbf{E}_2 \left( 1 - \frac{\rho_p}{p} \right) \right)^m \leq \mathbf{E}_2 \left( \left( 1 - \frac{\rho_p}{p} \right)^m \right).$$

From this we can bound  $\mu_k(m)$  from below by any subproduct, and we look at

$$\mu_k^*(m) = \prod_{p \leq m} \left( 1 - \frac{1}{p} \right)^{-km} \mathbf{E}_2 \left( \left( 1 - \frac{\rho_p}{p} \right)^m \right).$$

The probability that  $\rho_p$  is 1 is clearly equal to  $p^{-(k-1)}$  (there are only  $p$   $k$ -tuples with this property). Hence we have crude lower bounds

$$\mathbf{E}_2 \left( \left( 1 - \frac{\rho_p}{p} \right)^m \right) \geq \frac{1}{p^{k-1}} \left( 1 - \frac{1}{p} \right)^k$$

and

$$\mu_k(m) \geq \mu_k^*(m) \geq \prod_{p \leq m} \left( 1 + \frac{1}{p-1} \right)^{k(m-1)} \frac{1}{p^{k-1}}.$$

The logarithm of this expression is easily bounded from below as follows:

$$\begin{aligned} \log \mu_k(m) &\geq k(m-1) \sum_{p \leq m} \log \left( 1 + \frac{1}{p-1} \right) - (k-1) \sum_{p \leq m} \log p \\ &= km \log \log 3m + O(m), \end{aligned}$$

for  $m \geq 2$ , the implied constant depending only on  $k$ , by standard estimates, and we can incorporate trivially  $m = 1$  also.

To prove the corresponding upper bound, we split the Euler product (3.10) into two ranges: we write

$$\mu_k(m) = \mu_k^{(1)}(m) \mu_k^{(2)}(m),$$

where  $\mu_k^{(1)}(m)$  is the product over primes  $p < km$  (which includes the range used for the lower bound), while  $\mu_k^{(2)}(m)$  is the product over the other primes  $p \geq km$ . We will show that

$$\log \mu_k^{(1)}(m) \leq km \log \log 3m + O(m), \quad \log \mu_k^{(2)}(m) \ll \frac{m}{\log 2m},$$

with implied constants depending on  $k$ , and this will conclude the proof.

We start with small primes, and simply bound the expectation of  $(1 - \rho/p)^m$  by the trivial bound 1; this leads to

$$\log \mu_k^{(1)}(m) \leq -km \sum_{p < km} \log \left( 1 - \frac{1}{p} \right) = km \log \log 3m + O(m),$$

where the implied constant depends on  $k$ , again by standard estimates.

Next, we estimate  $\mu_k^{(2)}(m)$  more carefully. The logarithm (say  $\mathcal{L}(x)$ ) of the product restricted to  $km \leq p \leq x$  is given by

$$\mathcal{L}(x) = -km \sum_{km \leq p \leq x} \log(1 - p^{-1}) + \sum_{km \leq p \leq x} \log \mathbf{E}_2 \left( \left(1 - \frac{\rho_p}{p}\right)^m \right).$$

Using (3.7), we write first, for  $p \geq km$ , the upper bound

$$\begin{aligned} \mathbf{E}_2 \left( \left(1 - \frac{\rho_p}{p}\right)^m \right) &\leq \left(1 - \frac{k}{p}\right)^m (1 - \mathbf{P}_2(\rho_p < k)) + \mathbf{P}_2(\rho_p < k) \\ &= \left(1 - \frac{k}{p}\right)^m + \mathbf{P}_2(\rho_p < k) \left(1 - \left(1 - \frac{k}{p}\right)^m\right) \\ &\leq \left(1 - \frac{k}{p}\right)^m + \frac{mk^2(k-1)}{2p^2} \\ &\leq 1 - \frac{mk}{p} + \frac{m(m-1)k^2}{2p^2} + \frac{mk^2(k-1)}{2p^2}, \\ &= 1 - \frac{mk}{p} + \frac{m^2k^2}{2p^2} + \frac{mA_k}{2p^2} \end{aligned}$$

(with  $A_k = k^3 - 2k^2$ ) since

$$1 - mx \leq (1 - x)^m \leq 1 - mx + \frac{m(m-1)}{2}x^2 \quad \text{for } 0 \leq x \leq 1, m \geq 1.$$

Moreover, we have  $\log(1 - x) \leq -x - x^2/2$  for  $0 \leq x < 1$ , and hence after some rearranging, we obtain

$$\begin{aligned} \log \mathbf{E}_2 \left( \left(1 - \frac{\rho_p}{p}\right)^m \right) &\leq -\frac{mk}{p} + \frac{m^2k^2}{2p^2} + \frac{mA_k}{2p^2} - \frac{1}{2} \left( \frac{mk}{p} - \frac{m^2k^2}{2p^2} - \frac{mA_k}{2p^2} \right)^2 \\ &= -\frac{mk}{p} + \frac{m^3k^2}{p^3} - \frac{m^4k^4}{8p^4} + \frac{mA_k}{2p^2} - \frac{m^2kA_k}{2p^3} - \frac{m^2A_k^2 - 2m^3k^2A_k}{8p^4}, \end{aligned}$$

the terms involving  $(m^2k^2)/(2p^2)$  having cancelled out.

Summing over  $km \leq p \leq x$ , we can let  $x$  go to infinity in all but the first resulting term since they define convergent series; bounding the tail by

$$\sum_{p > km} \frac{1}{p^\sigma} \ll (km)^{1-\sigma} (\log 2km)^{-1},$$

leads to

$$\sum_{km \leq p \leq x} \log \mathbf{E}_2 \left( \left(1 - \frac{\rho_p}{p}\right)^m \right) \leq -km \sum_{km \leq p \leq x} \frac{1}{p} + O\left(\frac{m}{\log 2m}\right)$$

for all  $m$  and  $x \geq km$ , where the implied constant depends on  $k$ . Finally,

$$\log \mathcal{L}(x) \leq -km \sum_{km < p \leq x} \left( \frac{1}{p} + \log \left(1 - \frac{1}{p}\right) \right) + O\left(\frac{m}{\log 2m}\right),$$

and since  $p^{-1} + \log(1 - p^{-1})$  defines an absolutely convergent series with tail (for  $p > y$ ) decreasing like  $y^{-1}(\log y)^{-1}$ , we obtain the desired bound for

$$\log \mu_k^{(2)}(m) = \lim_{x \rightarrow +\infty} \mathcal{L}(x).$$

□

The existence of a limiting distribution (Theorem 1.2) is an easy consequence of this.

**Corollary 4.2.** *Let  $k \geq 1$  be a fixed integer. As  $h$  goes to infinity, the singular series  $\mathfrak{S}(\mathbf{h})$  for  $\mathbf{h} \in \Omega_1$ , i.e., such that  $|\mathbf{h}| \leq h$ , converges in law to the random singular series*

$$Z_Y = Z_{Y,k} = \prod_p \left(1 - \frac{1}{p}\right)^{-k} \left(1 - \frac{\rho_p}{p}\right)$$

on  $\Omega_2$ . In other words, there exists a probability law  $\nu_k$  on  $[0, +\infty[$ , which is the law of  $Z_Y$ , such that  $\mathfrak{S}(\mathbf{h})$ , for  $|\mathbf{h}| \leq h$ , becomes equidistributed with respect to  $\nu_k$ , or equivalently

$$\lim_{h \rightarrow +\infty} \frac{1}{h^k} \sum_{|\mathbf{h}| \leq h}^* f(\mathfrak{S}(\mathbf{h})) = \int_{\mathbf{R}^+} f(t) d\nu_k(t)$$

for any bounded continuous function on  $\mathbf{R}$ . Moreover we have

$$(4.2) \quad \mu_k(m) = \mathbf{E}_2(Z_Y^m) = \int_{\mathbf{R}^+} t^m d\nu_k(t).$$

*Proof.* First of all, using (3.10), the monotone and dominated convergence theorems and (3.6) imply that we have

$$(4.3) \quad \mu_k(m) = \mathbf{E}_2(Z_Y^m)$$

for all integers  $m \geq 1$ . Now a standard result of probability theory (the ‘‘method of moments’’) states that given a positive random variable  $X$  and a sequence of positive random variables  $(X_n)$ , such that  $\mathbf{E}(X^m) < +\infty$ ,  $\mathbf{E}(X_n^m) < +\infty$  for all  $n$  and  $m$ , the condition

$$\lim_{n \rightarrow +\infty} \mathbf{E}(X_n^m) = \mathbf{E}(X^m)$$

for all  $m \geq 1$  implies the convergence in law of  $X_n$  to  $X$ , if the moments  $\mathbf{E}(X^m)$  do not grow too fast (a sufficient, but not necessary condition). In fact, it is enough that the power series

$$\sum_{m \geq 0} i^m \frac{\mathbf{E}(X^m)}{m!} t^m$$

have a non-zero radius of convergence, which in our case holds (with  $X = Z_Y$ ) by the almost exponential upper bound for  $\mu_k(m)$  in Proposition 4.1. Finally, the formula (4.2) follows from (4.3).  $\square$

**Example 4.3.** As a corollary of Proposition 4.1 and symmetry, we have

$$\log \mu_k(2) = 2k \log \log 3k + O(k)$$

for  $k \geq 1$ .

Combined with the classical lower bound for non-vanishing arising from Cauchy’s inequality, it follows that for every fixed  $k \geq 1$ , we have

$$\liminf_{h \rightarrow +\infty} \frac{1}{h^k} |\{\mathbf{h} \mid |\mathbf{h}| \leq h \text{ and } \mathfrak{S}(\mathbf{h}) \neq 0\}| \geq \frac{\mu_k(1)^2}{\mu_k(2)} \geq \exp(-(2k \log \log 3k + O(k))).$$

This is close to the truth, as one can check by noting that we have in fact<sup>6</sup>

$$\lim_{h \rightarrow +\infty} \frac{1}{h^k} |\{\mathbf{h} \mid |\mathbf{h}| \leq h \text{ and } \mathfrak{S}(\mathbf{h}) \neq 0\}| = \mathbf{P}_2(Z_{Y,k} \neq 0) = \prod_{p \leq k} \mathbf{P}_2(\rho_p < p)$$

using the almost sure absolute convergence of the random Euler product  $Z_{Y,k}$ . We have the bounds

$$\frac{(p-1)^k}{p^k} \leq \mathbf{P}_2(\rho_p < p) \leq \frac{p(p-1)^k}{p^k}$$

(since, for  $p \leq k$ , a  $k$ -tuple will have  $\rho_p < p$  only if it omits at least one value in  $\mathbf{Z}/p\mathbf{Z}$ ; the lower bound follows by looking at those omitting 0, for instance, and the upper one is a union bound over the possible omitted values), from which we get

$$-k \log \log 3k + O(k) \leq \log \mathbf{P}_2(Z_{Y,k} \neq 0) \leq k - k \log \log 3k + O(k),$$

i.e., we have

$$\mathbf{P}_2(Z_{Y,k} \neq 0) = \exp(-k \log \log 3k + O(k)).$$

It follows from this that if we replace the space  $\Omega_1$  of all  $k$ -tuples with distinct entries by the much smaller one

$$\tilde{\Omega}_1 = \{\mathbf{h} \in \Omega_1 \mid \mathfrak{S}(\mathbf{h}) \neq 0\},$$

(which still depends on  $h$ , with cardinality  $\tilde{h}_k$ ), the singular series still has a limiting distribution when interpreted as a random variable on  $\tilde{\Omega}_1$  with  $h \rightarrow +\infty$ : indeed, this is the distribution  $\tilde{\nu}_k$  given by

$$\tilde{\nu}_k(A) = \frac{\nu_k(A \cap ]0, +\infty[)}{\nu_k(]0, +\infty[)},$$

since, for any integer  $m \geq 1$ , we have

$$\frac{1}{\tilde{h}_k} \sum_{\mathbf{h} \in \tilde{\Omega}_1} \mathfrak{S}(\mathbf{h})^m = \frac{h_k^*}{\tilde{h}_k} \mathbf{E}_1(\mathfrak{S}(\mathbf{h})^m) \rightarrow \frac{\mu_k(m)}{\mathbf{P}_2(Z_{Y,k} \neq 0)} = \int_{]0, +\infty[} t^m d\tilde{\nu}_k(t),$$

as  $h \rightarrow +\infty$ .

Of course, those moments do not satisfy the symmetry property enjoyed by  $\mu_k(m)$ .

*Remark 4.4.* Before going on to the second part of this paper, the following question seems natural: are there arithmetic consequences (possibly conditional, similarly to Gallagher's proof of (1.7)) of the existence of  $m$ -th moments of the singular series for  $k$ -tuples?

## 5. POISSON DISTRIBUTION FOR GENERAL PRIME PATTERNS

In this section, we prove Theorem 1.3, essentially by following Gallagher's reduction to averages of Euler products, which turn out to be easily computable after application of Proposition 2.1.

We fix a primitive family of polynomials  $\mathbf{f}$  with  $\mathfrak{S}(\mathbf{f}) \neq 0$  (the reader may want to review the notation in the introduction for what follows). To apply Gallagher's method, we also require some auxiliary families of polynomials, indexed by  $k$ -tuples.

<sup>6</sup> This does not follow directly from convergence in law for  $\mathfrak{S}(\mathbf{h})$ , but from the absolute convergence and local structure of the singular series.

Thus let  $k \geq 1$  be an integer and  $\mathbf{h}$  a  $k$ -tuple of integers. For our fixed primitive  $\mathbf{f}$ , we denote

$$\mathbf{f} \odot \mathbf{h} = (f_j(X + h_i))_{\substack{1 \leq j \leq m, \\ 1 \leq i \leq k}}$$

which is a family of  $km$  integer polynomials.

Technical difficulties will arise because this family may not be primitive, even if the components of  $\mathbf{h}$  are distinct (which is a necessary condition), i.e., we may have an equality

$$f_{j_1}(X + h_{i_1}) = f_{j_2}(X + h_{i_2}),$$

for some  $i_1 \neq i_2, j_1 \neq j_2$ .

For instance, we have  $(X, X + 2) \odot (3, 1) = (X + 3, X + 1, X + 5, X + 3)$  (in the case of twin primes). However, we will show that these degeneracies have no effect for the problem at hand. Moreover,  $\mathbf{f} \odot \mathbf{h}$  is primitive whenever  $\mathbf{h}$  has distinct arguments, in the following quite general situations:

- if  $m = 1$ ;
- if the degrees of the  $f_j$  are distinct;
- if no two among the polynomials  $f_j$  are related by a translation  $X \mapsto X + \alpha$ , for some  $\alpha \in \mathbf{Z}$ .

This means that the reader may well disregard the technical problems in a first reading (for the twin primes, see also Example 5.9 which explains a special reason why the degeneracies have no consequence then). The following lemma is already a first step, and we will need it before proving the full statement.

**Lemma 5.1.** *Let  $\mathbf{f}$  be a primitive family and  $k \geq 1$ . Then for any  $h \geq 1$ , we have*

$$|\{\mathbf{h} \mid |\mathbf{h}|_k \leq h, \quad \mathbf{f} \odot \mathbf{h} \text{ is not primitive}\}| \ll h^{k-1}$$

where the implied constant depends only on  $k$  and  $m$ .

*Proof.* Let  $I$  be the set of  $k$ -tuples  $\mathbf{h}$  with distinct components such that  $\mathbf{f} \odot \mathbf{h}$  is not primitive. If  $\mathbf{h} \in I$ , then there exists at least one relation of the type

$$(5.1) \quad f_{j_1}(X + h_{i_1}) = f_{j_2}(X + h_{i_2}), \quad i_1 \neq i_2, \quad j_1 \neq j_2,$$

hence

$$f_{j_1}(X) = f_{j_2}(X + h_{i_2} - h_{i_1}),$$

so the two polynomials differ by a “shift”. Let  $\mathcal{R}$  be the set of pairs  $(j_1, j_2)$  for which

$$f_{j_1}(X) = f_{j_2}(X + \delta(j_1, j_2))$$

for some integer  $\delta(j_1, j_2) \neq 0$ . Because the polynomials involved are non-constant, this integer is indeed unique. The cardinality of  $\mathcal{R}$  is bounded in terms of  $m$  only, and from the above, any  $k$ -tuple  $\mathbf{h} \in I$  must satisfy at least one relation

$$h_{i_1} - h_{i_2} = \delta(j_1, j_2),$$

for some  $i_1 \neq i_2$  and  $(j_1, j_2) \in \mathcal{R}$ . Each such relation is valid for at most  $h^{k-1}$  among the  $k$ -tuples with  $|\mathbf{h}| \leq h$ .  $\square$

We will deduce Theorem 1.3 from the following (unconditional) result, which is another instance of average of Euler products:

**Proposition 5.2.** *Let  $\mathbf{f} = (f_1, \dots, f_m)$  be a primitive family and  $k \geq 1$  an integer. Then we have*

$$\lim_{h \rightarrow +\infty} \frac{1}{h^k} \sum_{|\mathbf{h}| \leq h}^* \mathfrak{S}(\mathbf{f} \odot \mathbf{h}) = \mathfrak{S}(\mathbf{f})^k,$$

where  $\sum^*$  here restricts the summation to those  $k$ -tuples for which  $\mathbf{f} \odot \mathbf{h}$  is primitive.

*Remark 5.3.* Taking  $\mathbf{f} = (X)$ , with  $\mathfrak{S}(\mathbf{f}) = 1$  and  $\mathbf{f} \odot \mathbf{h} = (X + h_1, \dots, X + h_k)$ , we recover once more Gallagher's result (1.5).

We have the following complementary statement, which is also unconditional (recall that, in many cases, it holds for trivial reasons; it does *not* follow trivially from Lemma 5.1 because although fewer  $k$ -tuples are concerned, the number of prime seeds increases when  $\mathbf{f} \odot \mathbf{h}$  is not primitive).

**Lemma 5.4.** *Let  $\mathbf{f} = (f_1, \dots, f_m)$  be a primitive family with  $\mathfrak{S}(\mathbf{f}) \neq 0$ , and  $k \geq 1$  an integer. Then for any  $N \geq 2$ , if  $h \leq \lambda(\log N)^m$  for some  $\lambda > 0$ , and for any  $\varepsilon > 0$ , we have*

$$\sum_{\substack{|\mathbf{h}|_k \leq h \\ \mathbf{f} \odot \mathbf{h} \text{ not primitive}}}^* \pi(N; \mathbf{f} \odot \mathbf{h}) \ll \frac{N}{(\log N)^{1-\varepsilon}}$$

where  $\sum^*$  restricts the sum to those  $k$ -tuples with distinct entries, and where the implied constant depends only on  $k$ ,  $\mathbf{f}$ ,  $\lambda$  and  $\varepsilon$ .

Here is the proof of the (conditional) Poisson distribution, assuming those two results.

*Proof of Theorem 1.3.* The argument is essentially identical with that of Gallagher, but we reproduce it for completeness, and so that the necessary uniformity in the Bateman-Horn conjecture becomes clear.

Because the Poisson distribution is characterized by its moments, it is enough to prove that for any fixed integer  $k \geq 1$ , we have

$$\frac{1}{N} \sum_{n \leq N} \left( \pi(n + \lambda \delta(N, \mathbf{f}); \mathbf{f}) - \pi(n; \mathbf{f}) \right)^k \rightarrow \mathbf{E}(P_\lambda^k), \quad \text{as } N \rightarrow +\infty,$$

where  $P_\lambda$  is any Poisson random variable with mean  $\lambda$ .

Write  $h = \lambda \delta(N, \mathbf{f})$ . Expanding the left-hand side, we obtain

$$\frac{1}{N} \sum_{n \leq N} \left( \sum_{\substack{n < m_i \leq n+h \\ m_i \text{ } \mathbf{f}\text{-prime seed}}} \cdots \sum 1 \right)$$

where there are  $k$  sums over  $m_1, \dots, m_k$ . Write  $m_i = n + h_i$ , so that  $1 \leq h_i \leq h$ , and the condition becomes that  $f_j(n + h_i)$  is prime for all  $i$  and  $j$ , i.e., that  $n$  be an  $\mathbf{f} \odot \mathbf{h}$ -prime seed. Exchanging the order of summation, we get

$$\frac{1}{N} \sum_{|\mathbf{h}|_k \leq h} \pi(N; \mathbf{f} \odot \mathbf{h}).$$

Before applying (1.8), we need to account for the  $k$ -uples which do not necessarily have distinct components, and for those where  $\mathbf{f} \odot \mathbf{h}$  is *not* primitive.

For this, observe first that  $\pi(N; \mathbf{f} \odot \mathbf{h})$  only depends on the set containing the components of the  $k$ -tuple  $\mathbf{h}$ . This justifies the fact that the reorderings that

follow are permissible. For each  $r$ ,  $1 \leq r \leq k$ , and each  $r$ -tuple  $\mathbf{h}'$  with distinct components, the set of those  $k$ -tuples for which the set of values is given by the set of components of  $\mathbf{h}'$  has cardinality depending only on  $r$  and  $k$ , but independent of  $\mathbf{h}'$ , and in fact it is given by  $\binom{k}{r}$  (one can assume that  $\mathbf{h}' = (1, \dots, r)$ , and obtain a bijection

$$\left\{ \begin{array}{l} \{\text{suitable } k\text{-tuples}\} \\ \mathbf{h} \end{array} \right\} \begin{array}{l} \rightarrow \{\text{surjective maps } \{1, \dots, k\} \rightarrow \{1, \dots, r\}\} \\ \mapsto (f : i \mapsto h_i) \end{array}$$

between the two sets).

Then we can write

$$\frac{1}{N} \sum_{|\mathbf{h}|_k \leq h} \pi(N; \mathbf{f} \odot \mathbf{h}) = \frac{1}{N} \sum_{r=1}^k \frac{1}{r!} \binom{k}{r} \sum_{|\mathbf{h}'|_r \leq h}^* \pi(N; \mathbf{f} \odot \mathbf{h}')$$

where we divide by  $r!$  because we sum over all  $r$ -tuples instead of only ordered ones, and  $\sum^*$  restricts to  $r$ -tuples with distinct entries.

Now, for each  $r$ , we separate the sum over  $r$ -tuples for which  $\mathbf{f} \odot \mathbf{h}'$  is primitive from the other subsum. Applying (1.8) and using the easy bound

$$c(\mathbf{f} \odot \mathbf{h}') \ll c(\mathbf{f}) |\mathbf{h}'|_r^{\max \deg(f_j)},$$

(where the implied constant depends on  $r$  and  $\mathbf{f}$ ) the first sum (still denoted  $\sum^*$ ) is equal to

$$\sum_{r=1}^k \frac{1}{r!} \binom{k}{r} \frac{1}{\text{peg}(\mathbf{f})^r} \frac{1}{(\log N)^{rm}} \sum_{|\mathbf{h}'|_r \leq h}^* \mathfrak{S}(\mathbf{f} \odot \mathbf{h}') \left(1 + O\left(\frac{h^\varepsilon}{\log N}\right)\right),$$

for any  $\varepsilon > 0$ , where the implied constant depends on  $\mathbf{f}$ ,  $k$  and  $\varepsilon$ . Using Proposition 5.2 and the choice of  $h = \lambda \text{peg}(\mathbf{f}) \mathfrak{S}(\mathbf{f})^{-1} (\log N)^m$ , this converges as  $N \rightarrow +\infty$  to the limit

$$\sum_{r=1}^k \frac{\lambda^r}{r!} \binom{k}{r},$$

which is well-known to be the  $k$ -th moment of a Poisson distribution with mean  $\lambda$  (this is checked by Gallagher for instance, see [Ga, §3]). Hence, to conclude the proof, we need only notice that Lemma 5.4 (applied with  $k = r$  for  $1 \leq r \leq k$ ) implies (taking  $\varepsilon = 1/2$  for concreteness) that the complementary sum is bounded by

$$\frac{1}{N} \sum_{r=1}^k \frac{1}{r!} \binom{k}{r} \sum_{\substack{|\mathbf{h}'|_r \leq h \\ \mathbf{f} \odot \mathbf{h}' \text{ not primitive}}} \pi(N; \mathbf{f} \odot \mathbf{h}') \ll (\log N)^{-1/2}$$

for  $N \geq 2$ , where the implied constant depends on  $k$ ,  $\mathbf{f}$  and  $\lambda$ . Hence this second contribution goes to 0 as  $N \rightarrow +\infty$ , as desired.  $\square$

We now prove Proposition 5.2. This is the conjunction of the two following lemmas, where we use the same notation as in Section 3, but change a bit the definition of probability spaces. Precisely,

$$\Omega_2 = \prod_p (\mathbf{Z}/p\mathbf{Z})^k$$

is unchanged, but we let

$$\Omega_1 = \{\mathbf{h} = (h_1, \dots, h_k) \mid 1 \leq h_i \leq h, \mathbf{f} \odot \mathbf{h} \text{ is primitive}\}$$

with the counting probability measure (note that the condition forces  $\mathbf{h}$  to have distinct coordinates). By Lemma 5.1, note that we have

$$(5.2) \quad |\Omega_1| \sim h^k \quad \text{as } h \rightarrow +\infty.$$

The next lemma shows that the average of Euler product involved can be computed as if the components were independent:

**Lemma 5.5.** *Let  $\mathfrak{S}(\mathbf{f}) = (f_1, \dots, f_m)$  be a primitive family with  $\mathfrak{S}(\mathbf{f}) \neq 0$ . Then for any  $k \geq 1$ , we have*

$$\begin{aligned} \lim_{h \rightarrow +\infty} \frac{1}{h^k} \sum_{|\mathbf{h}| \leq h} \mathfrak{S}(\mathbf{f} \odot \mathbf{h}) &= \lim_{h \rightarrow +\infty} \mathbf{E}_1 \left( \prod_p \left(1 - \frac{1}{p}\right)^{-km} \left(1 - \frac{\nu_{p,\mathbf{f}}}{p}\right) \right) \\ &= \prod_p \mathbf{E}_2 \left( \left(1 - \frac{1}{p}\right)^{-km} \left(1 - \frac{\rho_{p,\mathbf{f}}}{p}\right) \right), \end{aligned}$$

where

$$\begin{aligned} \nu_{p,\mathbf{f}}(\mathbf{h}) &= \nu_p(\mathbf{f} \odot \mathbf{h}) \text{ for } \mathbf{h} = (h_1, \dots, h_k) \text{ with } h_i \geq 1, \\ \rho_{p,\mathbf{f}}(\mathbf{h}) &= |\{x \in \mathbf{Z}/p\mathbf{Z} \mid f_j(x + h_i) = 0 \text{ for some } i, j\}| \text{ for } \mathbf{h} \in (\mathbf{Z}/p\mathbf{Z})^r. \end{aligned}$$

The second lemma computes the limit locally:

**Lemma 5.6.** *Let  $\mathbf{f} = (f_1, \dots, f_m)$  be a primitive family. Then for any  $k \geq 1$  and any prime  $p$ , we have*

$$\mathbf{E}_2 \left( \left(1 - \frac{1}{p}\right)^{-km} \left(1 - \frac{\rho_{p,\mathbf{f}}}{p}\right) \right) = \left(1 - \frac{1}{p}\right)^{-km} \left(1 - \frac{\nu_p(\mathbf{f})}{p}\right)^k.$$

Looking at the definition (1.4) of  $\mathfrak{S}(\mathbf{f})$ , both lemmas together prove Proposition 5.2. We start by proving Lemma 5.6 because Lemma 5.5 is certainly plausible enough in view of Section 3, and the reader may be more interested by the final formal flourish.

*Proof of Lemma 5.6.* It suffices to compute

$$\mathbf{E}_2 \left( 1 - \frac{\rho_{p,\mathbf{f}}}{p} \right)$$

since the other factor is the same on both sides. We argue probabilistically, although one can also just expand the various sums (and do the same steps in a different language, as we did when proving the symmetry (1.6)). We can write

$$1 - \frac{\rho_{p,\mathbf{f}}}{p} = \frac{1}{p} |\mathbf{Z}/p\mathbf{Z} - M|$$

where  $M \subset \mathbf{Z}/p\mathbf{Z}$  is the (random) subset of those  $x \in \mathbf{Z}/p\mathbf{Z}$  such that  $f_j(x + h_i) = 0$  for some  $i$  and  $j$ . We write

$$|\mathbf{Z}/p\mathbf{Z} - M| = \sum_{x \in \mathbf{Z}/p\mathbf{Z}} (1 - \chi_M(x))$$

where  $\chi_M(x)$  is the random variable equal to one if  $x \in M$  and zero otherwise. We have

$$1 - \chi_M(x) = \prod_{1 \leq i \leq k} \prod_{1 \leq j \leq m} (1 - \mathbb{1}_{\{f_j(x + h_i) = 0\}}) = \prod_{1 \leq i \leq k} \xi_{\mathbf{f},i}(x),$$

say. Since  $\xi_{\mathbf{f},i}(x)$  only involves the  $i$ -th component of the random  $\mathbf{h} \in \Omega_2$ , the family  $(\xi_{\mathbf{f},i}(x))$  is an independent  $k$ -tuple of random variables. Consequently we derive

$$\begin{aligned} \mathbf{E}_2\left(1 - \frac{\rho_{p,\mathbf{f}}}{p}\right) &= \frac{1}{p} \sum_{x \in \mathbf{Z}/p\mathbf{Z}} \mathbf{E}_2\left(\prod_{1 \leq i \leq k} \xi_{\mathbf{f},i}(x)\right) \\ &= \frac{1}{p} \sum_{x \in \mathbf{Z}/p\mathbf{Z}} \prod_{1 \leq i \leq k} \mathbf{E}_2(\xi_{\mathbf{f},i}(x)). \end{aligned}$$

To conclude we notice that for every  $x$  and  $i$ ,  $\mathbf{h} \mapsto x + h_i$  is identically (uniformly) distributed, so that all  $\xi_{\mathbf{f},i}(x)$  are identically distributed like

$$\xi_{\mathbf{f}} = \xi_{\mathbf{f},1}(0) = \prod_{1 \leq j \leq m} (1 - \mathbb{1}_{\{f_j(h_1)=0\}}).$$

Hence all  $x$  give the same contribution, and we derive that

$$\mathbf{E}_2\left(1 - \frac{\rho_{p,\mathbf{f}}}{p}\right) = \mathbf{E}_2(\xi_{\mathbf{f}})^k = \mathbf{P}_2(f_1(h_1) \cdots f_m(h_1) \neq 0)^k = \left(1 - \frac{\nu_p(\mathbf{f})}{p}\right)^k,$$

since  $h_1$  is uniformly distributed in  $\mathbf{Z}/p\mathbf{Z}$ .  $\square$

To prove Lemma 5.5, we wish to apply Proposition 2.1. A complication is that, if  $\text{peg}(\mathbf{f}) \neq 1$ , the singular series  $\mathfrak{S}(\mathbf{f} \odot \mathbf{h})$  are not defined by absolutely convergent products, and therefore the result is not directly applicable. However, we can bypass this difficulty here without significant work because of the following fact: all the relevant Euler products can be *uniformly* “renormalized” to absolutely convergent ones. This is the content of the next lemma.

**Lemma 5.7.** *Let  $\mathbf{f}$  be a primitive family with  $\mathfrak{S}(\mathbf{f}) \neq 0$ , and let  $k \geq 1$  be an integer. There exist real numbers  $\gamma_p(\mathbf{f}) > 0$ , for all primes  $p$ , such that the product*

$$\prod_p \gamma_p(\mathbf{f})$$

*converges, and such that the following hold:*

(1) *For all prime  $p$ , and all  $k$ -tuple  $\mathbf{h} \in (\mathbf{Z}/p\mathbf{Z})^k$ , we have*

$$\left(1 - \frac{1}{p}\right)^{-km} \left(1 - \frac{\rho_{p,\mathbf{f}}(\mathbf{h})}{p}\right) = \gamma_p(\mathbf{f}) \times (1 + X_{p,\mathbf{f}}(\mathbf{h}))$$

*for some coefficients  $X_{p,\mathbf{f}}(\mathbf{h})$ , and for all  $k$ -tuple of integers  $\mathbf{h}$  such that  $\mathbf{f} \odot \mathbf{h}$  is primitive, the product*

$$(5.3) \quad \prod_p (1 + X_{p,\mathbf{f}}(\mathbf{h}))$$

*is absolutely convergent.*

(2) *We have*

$$\lim_{h \rightarrow +\infty} \frac{1}{h^k} \sum_{|\mathbf{h}|_k \leq h}^* \prod_p (1 + X_{p,\mathbf{f}}(\mathbf{h})) = \prod_p (1 + \mathbf{E}_2(X_{p,\mathbf{f}})),$$

*where the sum is over  $k$ -tuples with  $\mathbf{f} \odot \mathbf{h}$  primitive.*

*Proof.* (1) To define  $\gamma_p(\mathbf{f})$ , let  $\theta_j$ ,  $1 \leq j \leq m$ , be a complex root of the irreducible polynomial  $f_j$ , and let  $K_j = \mathbf{Q}(\theta_j)$  be the extension of  $\mathbf{Q}$  of degree  $\deg(f_j)$  generated by  $\theta_j$ . Then put

$$\gamma_p(\mathbf{f}) = \prod_{1 \leq i \leq m} \left(1 - \frac{1}{p}\right)^{k(r_j(p)-1)}$$

where  $r_j(n)$ , for  $n \geq 1$ , is the number of prime ideals of norm  $n$  in the ring of integers of  $K_j$ . In view of this definition, to check first that the product of  $\gamma_p(\mathbf{f})$  converges, we can do so for each  $f_j$  separately. Then the statement follows, after taking the logarithm of a partial product over  $p \leq X$ , from the well-known asymptotic formula

$$\sum_{p \leq X} \frac{r_j(p)}{p} = \sum_{p \leq X} \frac{1}{p} + c(K_j) + O((\log X)^{-1})$$

for  $X \geq 2$ , where  $c(K_j)$  is a constant depending only on  $K_j$ , and the implied constant also depends only on  $K_j$ .

It therefore remains to prove that the product (5.3) is absolutely convergent for any  $k$ -tuple of integers  $\mathbf{h}$  with  $\mathbf{f} \odot \mathbf{h}$  primitive. To do so, we claim that there exists an integer  $D(\mathbf{h}) \geq 1$  (which may also depend on  $\mathbf{f}$ ) such that, for  $p \nmid D(\mathbf{h})$ , we have

$$(5.4) \quad \rho_{p,\mathbf{f}}(\mathbf{h}) = k \sum_{j=1}^m \nu_p(f_j) = k \sum_{j=1}^m r_j(p).$$

The desired convergence then follows from that of

$$\prod_{p \nmid D(\mathbf{h})} \gamma_p(\mathbf{f})^{-1} \left(1 - \frac{1}{p}\right)^{-km} \left(1 - \frac{\rho_{p,\mathbf{f}}(\mathbf{h})}{p}\right) = \prod_{p \nmid D(\mathbf{h})} \left(1 - \frac{1}{p}\right)^{-\rho_{p,\mathbf{f}}(\mathbf{h})} \left(1 - \frac{\rho_{p,\mathbf{f}}(\mathbf{h})}{p}\right),$$

and the latter is clear since the  $p$ -factor can be written  $1 + O(p^{-2})$ , where the implied constant depends only on  $k$  and  $\mathbf{f}$ .

The existence of  $D(\mathbf{h})$  is easy; first, let

$$D_1(\mathbf{h}) = \left| \prod_{(i,j) \neq (i',j')} \text{Res}(f_j(X + h_i), f_{j'}(X + h_{i'})) \right|,$$

where  $\text{Res}(\cdot, \cdot)$  is the resultant of two polynomials. By compatibility of the resultant with reduction modulo  $p$ , we have  $p \mid D_1(\mathbf{h})$  if and only if, for some  $(i, j) \neq (i', j')$ , there exists a common zero  $x \in \mathbf{Z}/p\mathbf{Z}$  of  $f_j(X + h_i)$  and  $f_{j'}(X + h_{i'})$ . By contraposition, we first obtain

$$\rho_{p,\mathbf{f}}(\mathbf{h}) = k\nu_p(\mathbf{f}) = k \sum_{j=1}^m \nu_p(f_j),$$

for  $p \nmid D_1(\mathbf{h})$  (the sets of zeros modulo  $p$  of the components of  $\mathbf{f} \odot \mathbf{h}$  are then distinct, and obviously there are as many, namely the sum  $\nu_p(\mathbf{f})$  of the  $\nu_p(f_j)$ , for each of the  $k$  shifts  $h_i$ ).

Next, it is a standard fact of algebraic number theory that for each  $j$ , there exists an integer  $\Delta_j \geq 1$  such that  $\nu_p(f_j) = r_j(p)$  for  $p \nmid \Delta_j$ . Thus we can take

$$D(\mathbf{h}) = D_1(\mathbf{h}) \prod_{1 \leq j \leq m} \Delta_j$$

to obtain the second equality in (5.4).

Note that  $D(\mathbf{h})$  is non-zero (hence  $\geq 1$ ) because otherwise, there would exist a common zero  $\theta \in \mathbf{C}$  of  $f_j(X+h_i)$  and  $f_{j'}(X+h_{i'})$ , and because those are irreducible integral primitive<sup>7</sup> polynomials with positive leading coefficient, this is only possible if

$$f_j(X+h_i) = f_{j'}(X+h_{i'}),$$

which is excluded by the assumption that  $\mathbf{f} \odot \mathbf{h}$  be primitive.

Note in passing the estimate

$$D(\mathbf{h}) \ll (2|\mathbf{h}|_k)^{2k^2 m \sum \deg(f_j)}$$

for all  $\mathbf{h}$ , where the implied constant depends only on  $\mathbf{f}$ ; this follows straightforwardly from the determinant expression of the resultant in  $D_1(\mathbf{h})$  (see, e.g., [L, §V.10]).

(2) With the bounds we have proved on  $X_{p,\mathbf{f}}(\mathbf{h})$  (leading to an analogue of Lemma 3.1), and the estimate on  $D(\mathbf{h})$  (analogue of (3.5)), together with Lemma 5.1 to ensure that the equidistribution of  $k$ -tuples modulo squarefree integers  $q$  remains valid (compare with (3.11)), we can pretty much follow the steps of the proof of Theorem 1.1. We also use (5.2) to go from the limit of the expectation on  $\Omega_1$  to summing over  $k$ -tuples normalized by  $1/h^k$  and taking  $h \rightarrow +\infty$ . The details are left to the reader.  $\square$

*Proof of Lemma 5.5.* We have first

$$\begin{aligned} \frac{1}{h^k} \sum_{|\mathbf{h}| \leq h} \mathfrak{S}(\mathbf{f} \odot \mathbf{h}) &= \frac{1}{h^k} \sum_{|\mathbf{h}| \leq h} \left( \prod_p \gamma_p(\mathbf{f}) \right) \prod_p (1 + X_{p,\mathbf{f}}(\mathbf{h})) \\ &\rightarrow \left( \prod_p \gamma_p(\mathbf{f}) \right) \prod_p (1 + \mathbf{E}_2(X_{p,\mathbf{f}})) \quad \text{as } h \rightarrow +\infty, \end{aligned}$$

by the above, and then we can simply write this limit as

$$\begin{aligned} \left( \prod_p \gamma_p(\mathbf{f}) \right) \prod_p (1 + \mathbf{E}_2(X_{p,\mathbf{f}})) &= \prod_p \mathbf{E}_2(\gamma_p(\mathbf{f})(1 + X_{p,\mathbf{f}})) \\ &= \prod_p \mathbf{E}_2\left( \left(1 - \frac{1}{p}\right)^{-km} \left(1 - \frac{\rho_{p,\mathbf{f}}(\mathbf{h})}{p}\right) \right). \end{aligned}$$

$\square$

We conclude with the last remaining part of the proof, namely Lemma 5.4. The following proof can almost certainly be improved, but although the statement becomes fairly clear after checking one or two examples, the author has not found a cleaner way to deal with the apparent possibilities of combinatorial complications. The point is that as  $\mathbf{f} \odot \mathbf{h}$  becomes “less primitive” (i.e., there are less distinct elements among the  $km$  polynomials involved), the number of prime seeds  $\leq N$  should increase (by a power of  $(\log N)$ ), but also the number of  $k$ -tuples with this property diminishes (by a power of  $h \leq \lambda(\log N)^m$ ), and this gain has to compensate for the loss.

*Proof of Lemma 5.4.* We first quote a standard sieve upper-bound for an individual primitive family  $\mathbf{f}$  (with  $m$  elements), which is uniform, and which allows us to

<sup>7</sup> In the sense that the gcd of their coefficients is 1.

prove the lemma unconditionally: for  $N \geq 2$ , for any  $k$ -tuple  $\mathbf{h}$  with distinct elements for which  $\mathbf{f} \odot \mathbf{h}$  contains  $\ell$  distinct components, we have

$$(5.5) \quad \pi(N; \mathbf{f} \odot \mathbf{h}) \ll (\log \log 3|\mathbf{h}|)^{km} \frac{N}{(\log N)^\ell},$$

where the implied constant depends only on  $k$  and  $\mathbf{f}$ . Precisely, (5.5) for  $k$ -tuples follows immediately from, e.g, Th. 2.3 in [HR], and it is easy to adapt this to the case at hand since uniformity is only asked with respect to  $\mathbf{h}$ . Note also that, since the application we give is conditional on much stronger statements like (1.8), we could also apply the latter for this purpose.

Now, as in the proof of Lemma 5.1, we denote by  $I$  the set of  $k$ -tuples  $\mathbf{h}$  with distinct components such that  $\mathbf{f} \odot \mathbf{h}$  is not primitive. Recall  $\mathcal{R}$  is the set of pairs  $(j_1, j_2)$  for which

$$f_{j_1}(X) = f_{j_2}(X + \delta(j_1, j_2))$$

for some (unique) integer  $\delta(j_1, j_2) \neq 0$ .

We continue as follows: for an  $\mathbf{h} \in I$ , let  $\Gamma_{\mathbf{h}}$  be the graph with vertex set  $\{1, \dots, k\}$  and with (unoriented) edges  $(i_1, i_2)$  corresponding to those indices for which the relation

$$(5.6) \quad h_{i_1} - h_{i_2} = \delta(j_1, j_2)$$

holds for some  $(j_1, j_2) \in \mathcal{R}$ ; the proof of Lemma 5.1 shows that there is at least one edge. Because the number of possibilities for  $\Gamma_{\mathbf{h}}$  is clearly bounded in terms of  $k$  only, and we allow a constant depending on  $k$  in our estimate, we may continue by fixing one possible graph  $\Gamma$  and assuming that *all*  $\mathbf{h} \in I$  satisfy  $\Gamma_{\mathbf{h}} = \Gamma$ .

This being done, we first estimate from above the number of  $k$ -tuples which lie in  $I$  (under the above assumption that the graph is fixed!). We claim that

$$(5.7) \quad |\{\mathbf{h} \in I \mid |\mathbf{h}| \leq h\}| \leq h^c$$

where  $c = |\pi_0(\Gamma)|$  is the number of connected components of  $\Gamma$ .

To see this, notice that each connected component  $C$  corresponds to a set of variables which are *independent* of all others, so that  $I$  is the product over the connected components of sets  $I_C$  of  $|C|$ -tuples satisfying the relations (5.6) dictated by  $C$ . Now we have

$$|\{\mathbf{h} \in I_C \mid |\mathbf{h}| \leq h\}| \leq h,$$

because  $C$  is connected: if we fix some vertex  $i_0$  of  $C$ , then for any choice of  $h_{i_0}$ , the value of  $h_i$  is determined by means of the relations (5.6) for all vertices  $i$  of  $C$ , using induction on the length of a path from  $i_0$  to  $i$  (which exists by connectedness).

Taking the product over  $C$  of these individual upper bounds, we obtain the desired estimate (5.7).

We next need to estimate from below the number of distinct elements in the family  $\mathbf{f} \odot \mathbf{h}$  for a fixed  $\mathbf{h} \in I$  (still under the assumption that the graph  $\Gamma_{\mathbf{h}} = \Gamma$  is fixed).

Let again  $C$  be a connected component of the graph  $\Gamma$ . We consider the set (say  $\{\mathbf{f} \odot \mathbf{h}\}_C$ ) of polynomials of the form  $f_j(X + h_i)$ , where  $1 \leq j \leq m$  and  $i$  is a vertex of  $C$ . We claim this set contains at least  $m+1$  distinct polynomials if  $C$  has at least 2 vertices, and  $m$  if  $C$  is a singleton. Indeed, fixing a vertex  $i_0$  of  $C$ , the set contains the polynomials  $f_j(X + h_{i_0})$ , which are distinct since  $\mathbf{f}$  is a primitive family. This already takes care of the case where  $C$  is a singleton, so assume now that  $C$  contains

at least another vertex  $i$ . If all the  $m$  distinct polynomials  $f_j(X + h_i)$  were already in the set  $\{f_j(X + h_{i_0})\}$ , this would define a permutation  $\sigma$  of  $\{1, \dots, m\}$  such that

$$f_j(X + h_i) = f_{\sigma(j)}(X + h_{i_0}), \quad 1 \leq j \leq m.$$

Consider a cycle  $(j_1, \dots, j_\ell)$  of length  $\ell$  in the decomposition of  $\sigma$ ; applying the identity to  $j_1$ ,  $\sigma(j_1) = j_2$ , etc, in turn, we derive the identity

$$f_{j_1}(X) = f_{\sigma^\ell(j_1)}(X + (\ell - 1)(h_{i_0} - h_i)) = f_{j_1}(X + (\ell - 1)(h_{i_0} - h_i)).$$

Since  $f_j$  is non-constant and  $h_{i_0} \neq h_i$ , we deduce that  $\ell = 1$ ; this holding for all cycles in  $\sigma$  would mean that  $\sigma$  is the identity, but then  $f_1(X + h_i) = f_1(X + h_{i_0})$  again contradicts the fact that  $\mathbf{h}$  has distinct components. This means that  $\sigma$  can not exist, and so the set  $\{f_j(X + h_i)\}$  contains at least one polynomial not among the first  $m$  ones, which was our objective.

Next observe that, by the very definition of the graph  $\Gamma$ , the sets  $\{\mathbf{f} \odot \mathbf{h}\}_C$  are disjoint when  $C$  runs over the connected components of  $\Gamma$ , and hence we find that any  $\mathbf{f} \odot \mathbf{h}$  contains at least  $cm + d$  elements, where  $d$  is the number of connected components of  $\Gamma$  which are not singletons. Note that  $d \geq 1$ , because  $\Gamma$  has at least one edge.

We finally estimate the contribution of  $k$ -tuples in  $I$  using (5.5) and (5.7): we obtain

$$\frac{1}{N} \sum_{\substack{\mathbf{h} \in I \\ |\mathbf{h}| \leq h}} \pi(N; \mathbf{f} \odot \mathbf{h}) \ll h^c (\log 2h)^{km} (\log N)^{-cm-d}$$

where the implied constant depends on  $k$  and  $\mathbf{f}$ . If  $h \leq \lambda(\log N)^m$ , as assumed in Lemma 5.4, we obtain

$$\frac{1}{N} \sum_{\substack{\mathbf{h} \in I \\ |\mathbf{h}| \leq h}} \pi(N; \mathbf{f} \odot \mathbf{h}) \ll (\log N)^{-d+\varepsilon}$$

for any  $\varepsilon > 0$ , where the implied constant depends on  $k$ ,  $\lambda$ ,  $\mathbf{f}$  and  $\varepsilon$ . Since  $d \geq 1$ , the lemma is finally proved.  $\square$

*Remark 5.8.* The gain of  $(\log N)^{-1}$  is indeed the best possible in general. Consider for example the primitive family  $\mathbf{f} = (f_1, f_2, f_3) = (X^2 + 7, (X + 2)^2 + 7, (X + 4)^2 + 7)$  for which it is easy to check that  $\mathfrak{S}(\mathbf{f}) \neq 0$  (7 is not a square modulo 3 or 5, and each  $f_j(0)$  is odd). We have relations  $f_1(X + 2) = f_2(X)$ ,  $f_2(X + 2) = f_3(X)$ .

Consider  $k = 2$ . If we look at 2-tuples  $\mathbf{h} = (h_1, h_2)$  for which  $h_2 = h_1 + 2$ , we obtain

$$\begin{aligned} \mathbf{f} \odot \mathbf{h} &= (f_1(X + h_1), f_2(X + h_1), f_3(X + h_1), \\ &\quad f_1(X + h_2), f_2(X + h_2), f_3(X + h_2)) \\ &= (f_1(X + h_1), f_1(X + h_1 + 2), f_1(X + h_1 + 4), \\ &\quad f_1(X + h_2), f_1(X + h_2 + 2), f_1(X + h_2 + 4)) \\ &= (f_1(X + h_2 - 2), f_1(X + h_2), f_1(X + h_2 + 2), \\ &\quad f_1(X + h_2), f_1(X + h_2 + 2), f_1(X + h_2 + 4)), \end{aligned}$$

which contains 4 distinct polynomials. With  $h \asymp \lambda(\log N)^3$ , those 2-tuples with  $|\mathbf{h}| \leq h$  contribute about  $N(\log N)^{3-4}$  to the sum of Lemma 5.4 (under (1.8), of course).

Finally, here are a few examples.

**Example 5.9.** (1) If we take  $\mathbf{f}_1 = (X, X + 2)$ , we obtain that the number of twin primes  $(p, p + 2)$  with  $n < p \leq n + \lambda(\log n)^2$  should be approximately distributed like a Poisson random variable with mean

$$2\lambda \prod_{p \geq 3} \left(1 - \frac{1}{(p-1)^2}\right) \approx 1.320336593 \dots \lambda.$$

Similarly, if we take  $\mathbf{f}_2 = (X, 2X + 1)$ , we find that the number of Germain primes (i.e., primes  $p$  with  $2p + 1$  also prime) with  $n < p \leq n + \lambda(\log n)^2$  should be approximately distributed like a Poisson random variable with mean

$$\lambda \mathfrak{S}(\mathbf{f}_2) = 2\lambda \prod_{p \geq 3} \left(1 - \frac{1}{(p-1)^2}\right).$$

Two further remarks are interesting here. First, the proof of Theorem 1.3 shows that whenever  $\mathbf{f}$  consists of linear polynomials (in particular for those two results), “only” the (uniform) Hardy-Littlewood conjecture is needed. In other words, no assumption is required beyond those of Gallagher’s original result for the primes themselves.

Secondly, if one is interested in the case of twin primes in particular, Lemma 5.4 has a trivial proof from the following coincidence: if  $\mathbf{f} = (X, X + 2)$ ,  $\mathbf{h}$  has distinct entries, and  $\mathbf{f} \odot \mathbf{h}$  is not primitive, then

$$\mathfrak{S}(\mathbf{f} \odot \mathbf{h}) = 0, \quad \pi(N; \mathbf{f} \odot \mathbf{h}) \leq 1.$$

Indeed, if  $\mathbf{f} \odot \mathbf{h}$  is not primitive, we have  $k \geq 2$  and an equality  $h_{i_2} = h_{i_1} + 2$  for some  $i_1, i_2$ . The family  $\mathbf{f} \odot \mathbf{h}$  contains in particular the three polynomials  $X + h_{i_1}$ ,  $X + h_{i_2} = X + h_{i_1} + 2$  and  $X + h_{i_2} + 2 = X + h_{i_1} + 4$ . Hence, to be a prime seed for  $\mathbf{f} \odot \mathbf{h}$ , an integer  $n \geq 1$  must be such that, in particular, the triple  $(n + h_{i_1}, n + h_{i_1} + 2, n + h_{i_1} + 4)$  consists of prime numbers. But those three numbers are distinct modulo 3, showing that  $\nu_3(\mathbf{f} \odot \mathbf{h}) = 3$ , and the only possible case is  $(n, n + 2, n + 4) = (3, 5, 7)$ . (Examples such as  $\mathbf{f} = (X^2 + 7, (X + 2)^2 + 7)$  and  $\mathbf{h} = (3, 1)$  show that this special situation where imprimitive  $k$ -tuples lead to vanishing singular series for  $\mathbf{f} \odot \mathbf{h}$  is indeed a coincidence).

(2) If we take  $\mathbf{f}_3 = (X^2 + 1)$ , and renormalize in an obvious way, we find that the number of primes of the form  $p = n^2 + 1$  in an interval of the form  $N^2 < n \leq (N + \lambda(\log N))^2$  should be approximately distributed like a Poisson random variable with mean

$$\lambda \mathfrak{S}(\mathbf{f}_3) = \frac{4\lambda}{\pi} \prod_{p \equiv 1 \pmod{4}} \left(1 - \frac{1}{(p-1)^2}\right) \prod_{p \equiv 3 \pmod{4}} \left(1 - \frac{1}{p^2 - 1}\right).$$

#### REFERENCES

- [BH] P. T. Bateman and R. A. Horn, *A heuristic asymptotic formula concerning the distribution of prime numbers*, Mathematics of Computation 16 (1962), 363–367.
- [1] [BJ] H. Bohr and B. Jessen, *On the distribution of the values of the Riemann zeta function*, Amer. J. of Math. 58 (1936), 35–44.
- [CM] J. Cogdell and P. Michel, *On the complex moments of symmetric power L-functions at  $s = 1$* , Int. Math. Res. Not. 2004, no. 31, 1561–1617
- [Ga] P.X. Gallagher, *On the distribution of primes in short intervals*, Mathematika 23 (1976), 4–9.
- [GPY] D. Goldston, J. Pintz and C. Y. Yıldırım, *Primes in tuples I*, Annals of Mathematics 170 (2009), 819–862.

- [GS] A. Granville and K. Soundararajan, *The distribution of values of  $L(1, \chi_d)$* , *Geom. Funct. Anal.* 13 (2003), 992–1028.
- [GT] B. Green and T. Tao, *Linear equations in primes*, *Annals of Math.* (to appear).
- [HR] H. Halberstam and H.E. Richert, *Sieve methods*, Academic Press 1974.
- [HL] G.H. Hardy and J.E. Littlewood, *Some problems of 'Partitio Numerorum' III. On the expression of a number as a sum of primes*, *Acta Math.* 44 (1923), 1–70.
- [IK] H. Iwaniec and E. Kowalski, *Analytic number theory*, A.M.S. Coll. Publ. 53, 2004.
- [KS] J. Keating and N. Snaith, *Random matrix theory and  $\zeta(1/2 + it)$* , *Comm. Math. Phys.* 214 (2000), 57–89.
- [K1] E. Kowalski, *Petits écarts entre nombres premiers, d'après Goldston, Pintz et Yıldırım*, *Séminaire Bourbaki*, exp. 959, Astérisque 311 (2007), 177–210.
- [La] Y. Lamzouri, *Distribution of the values of  $L$ -function at the edge of the critical strip*, *Proc. London Math. Soc.* (to appear).
- [L] S. Lang, *Algebra*, 2nd edition, Addison-Wesley, 1984.
- [MS] H.L. Montgomery and K. Soundararajan, *Primes in short intervals*, *Comm. Math. Phys.* 252 (2004), 589–617.
- [S] A. Schinzel, and W. Sierpiński, *Sur certaines hypothèses concernant les nombres premiers. Remarque*, *Acta Arithm.* 4 (1958), 185–208. 1958.
- [Si] B. Simon, *The classical moment problem as a self-adjoint finite difference operator*, *Advances in Math.* 137 (1998), 82–203.

ETH ZÜRICH – D-MATH, RÄMISTRASSE 101, 8092 ZÜRICH, SWITZERLAND  
E-mail address: kowalski@math.ethz.ch