

# Random Graphs and the Parity Quantifier

Phokion G. Kolaitis  
UC Santa Cruz  
and IBM Almaden Research Center.  
kolaitis@cs.ucsc.edu

Swastik Kopparty  
CSAIL, MIT  
swastik@mit.edu

## ABSTRACT

The classical zero-one law for first-order logic on random graphs says that for every first-order property  $\varphi$  in the theory of graphs and every  $p \in (0, 1)$ , the probability that the random graph  $G(n, p)$  satisfies  $\varphi$  approaches either 0 or 1 as  $n$  approaches infinity. It is well known that this law fails to hold for any formalism that can express the parity quantifier: for certain properties, the probability that  $G(n, p)$  satisfies the property need not converge, and for others the limit may be strictly between 0 and 1.

In this work, we capture the limiting behavior of properties definable in first order logic augmented with the parity quantifier,  $\text{FO}[\oplus]$ , over  $G(n, p)$ , thus eluding the above hurdles. Specifically, we establish the following “modular convergence law”:

For every  $\text{FO}[\oplus]$  sentence  $\varphi$ , there are two explicitly computable rational numbers  $a_0, a_1$ , such that for  $i \in \{0, 1\}$ , as  $n$  approaches infinity, the probability that the random graph  $G(2n + i, p)$  satisfies  $\varphi$  approaches  $a_i$ .

Our results also extend appropriately to FO equipped with  $\text{Mod}_q$  quantifiers for prime  $q$ .

In the process of deriving the above theorem, we explore a new question that may be of interest in its own right. Specifically, we study the joint distribution of the subgraph statistics modulo 2 of  $G(n, p)$ : namely, the number of copies, mod 2, of a fixed number of graphs  $F_1, \dots, F_\ell$  of bounded size in  $G(n, p)$ . We first show that every  $\text{FO}[\oplus]$  property  $\varphi$  is almost surely determined by subgraph statistics modulo 2 of the above type. Next, we show that the limiting joint distribution of the subgraph statistics modulo 2 depends only on  $n \bmod 2$ , and we determine this limiting distribution completely. Interestingly, both these steps are based on a common technique using multivariate polynomials over finite fields and, in particular, on a new generalization of the Gowers norm that we introduce.

The first step above is analogous to the Razborov-Smolensky method for lower bounds for AC0 with parity gates, yet

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

STOC'09, May 31–June 2, 2009, Bethesda, Maryland, USA.  
Copyright 2009 ACM 978-1-60558-506-2/09/05 ...\$5.00.

stronger in certain ways. For instance, it allows us to obtain examples of simple graph properties that are exponentially uncorrelated with every  $\text{FO}[\oplus]$  sentence, which is something that is not known for  $\text{AC0}[\oplus]$ .

## Categories and Subject Descriptors

F.4.1 [Theory of Computation]: Mathematical Logic and Formal Languages—*Mathematical Logic*; G.2.2 [Mathematics of Computing]: Discrete Mathematics—*Graph Theory*

## General Terms

Theory

## Keywords

First-order logic, 0-1 Law, Modular convergence, Polynomials, Finite fields, AC0

## 1. INTRODUCTION

For quite a long time, combinatorialists have studied the asymptotic probabilities of properties on classes of finite structures, such as graphs and partial orders. Assume that  $\mathcal{C}$  is a class of finite structures and let  $\text{Pr}_n$ ,  $n \geq 1$ , be a sequence of probability measures on all structures in  $\mathcal{C}$  with  $n$  elements in their domain. If  $Q$  is a property of some structures in  $\mathcal{C}$  (that is, a decision problem on  $\mathcal{C}$ ), then the *asymptotic probability*  $\text{Pr}(Q)$  of  $Q$  on  $\mathcal{C}$  is defined as  $\text{Pr}(Q) = \lim_{n \rightarrow \infty} \text{Pr}_n(Q)$ , provided this limit exists. In this paper, we will be focusing on the case when  $\mathcal{C}$  is the class  $\mathcal{G}$  of all finite graphs, and  $\text{Pr}_n = G(n, p)$  for constant  $p$ ; this is the probability distribution on  $n$ -vertex undirected graphs where between each pair of nodes an edge appears with probability  $p$ , independently of other pairs of nodes. For example, for this case, the asymptotic probabilities  $\text{Pr}(\text{CONNECTIVITY}) = 1$  and  $\text{Pr}(\text{HAMILTONICITY}) = 1$ ; in contrast, if  $\text{Pr}_n = G(n, p(n))$  with  $p(n) = 1/n$ , then  $\text{Pr}(\text{CONNECTIVITY}) = 0$  and  $\text{Pr}(\text{HAMILTONICITY}) = 0$ .

Instead of studying separately one property at a time, it is natural to consider formalisms for specifying properties of finite structures and to investigate the connection between the expressibility of a property in a certain formalism and its asymptotic probability. The first and most celebrated such connection was established by Glebskii et al. [6] and, independently, by Fagin [5], who showed that a *0-1 law* holds for first-order logic<sup>1</sup> FO on the random graph  $G(n, p)$  with  $p$  a constant in  $(0, 1)$ ; this means that if  $Q$  is a property of

<sup>1</sup>Recall that the formulas of first-order logic on graphs

graphs expressible in FO and  $\text{Pr}_n = G(n, p)$  with  $p$  a constant in  $(0, 1)$ , then  $\text{Pr}(Q)$  exists and is either 0 or 1. This result became the catalyst for a series of investigations in several different directions. Specifically, one line of investigation [17, 15] investigated the existence of 0-1 laws for first-order logic FO on the random graph  $G(n, p(n))$  with  $p(n) = n^{-\alpha}$ ,  $0 < \alpha < 1$ . Since first-order logic on finite graphs has limited expressive power (for example, FO cannot express CONNECTIVITY and 2-COLORABILITY), a different line of investigation pursued 0-1 laws for extensions of first-order logic on the random graph  $G(n, p)$  with  $p$  a constant in  $(0, 1)$ . In this vein, it was shown in [2, 10] that the 0-1 law holds for extensions of FO with fixed-point operators, such as least fixed-point logic LFP, which can express CONNECTIVITY and 2-COLORABILITY. As regards to higher-order logics, it is clear that the 0-1 law fails even for existential second-order logic ESO, since it is well known that  $\text{ESO} = \text{NP}$  on finite graphs [4]. In fact, even the *convergence law* fails for ESO, that is, there are ESO-expressible properties  $Q$  of finite graphs such that  $\text{Pr}(Q)$  does *not* exist. For this reason, a separate line of investigation pursued 0-1 laws for syntactically-defined subclasses of NP. Eventually, this investigation produced a complete classification of the quantifier prefixes of ESO for which the 0-1 law holds [10, 11, 13], and provided a unifying account for the asymptotic probabilities of such NP-complete problems as  $k$ -COLORABILITY,  $k \geq 3$ , and SATISFIABILITY.

Let  $L$  be a logic for which the 0-1 law (or even just the convergence law) holds on the random graph  $G(n, p)$  with  $p$  a constant in  $(0, 1)$ . An immediate consequence of this is that  $L$  cannot express any *counting* properties, such as EVEN CARDINALITY (“there is an even number of nodes”), since for each  $n$ ,  $\text{Pr}_{2n}(\text{EVEN CARDINALITY}) = 1$  and  $\text{Pr}_{2n+1}(\text{EVEN CARDINALITY}) = 0$ . In this paper, we turn the tables around and systematically investigate the asymptotic probabilities of properties expressible in extensions of FO with *counting quantifiers*  $\text{Mod}_q^i$ , where  $q$  is a prime number. The most prominent such extension is  $\text{FO}[\oplus]$ , which is the extension of FO with the *parity quantifier*  $\text{Mod}_2^1$ . The syntax of  $\text{FO}[\oplus]$  augments the syntax of FO with the following formation rule: if  $\varphi(y)$  is a  $\text{FO}[\oplus]$ -formula, then  $\oplus y\varphi(y)$  is also a  $\text{FO}[\oplus]$ -formula; this formula is true if the number of  $y$ 's that satisfy  $\varphi(y)$  is odd (analogously,  $\text{Mod}_q^i y\varphi(y)$  is true if the number of  $y$ 's that satisfy  $\varphi(y)$  is congruent to  $i \pmod q$ ). A typical property on graphs expressible in  $\text{FO}[\oplus]$  (but not in FO) is  $\mathcal{P} := \{G : \text{every vertex in } G \text{ has odd degree.}\}$ , since a graph is in  $\mathcal{P}$  if and only if it satisfies the  $\text{FO}[\oplus]$ -sentence  $\forall x \oplus y E(x, y)$ .

Our main result (see Theorem 2.1) is a *modular convergence law* for  $\text{FO}[\oplus]$  on  $G(n, p)$  with  $p$  a constant in  $(0, 1)$ . This law asserts that if  $\varphi$  is a  $\text{FO}[\oplus]$ -sentence, then there are two explicitly computable rational numbers  $a_0, a_1$ , such that, as  $n \rightarrow \infty$ , the probability that the random graph  $G(2n + i, p)$  satisfies  $\varphi$  approaches  $a_i$ , for  $i = 0, 1$ . Moreover,  $a_0$  and  $a_1$  are computable and are of the form  $r/2^s$ , where  $r$  and  $s$  are non-negative integers. We also establish that an analogous modular convergence law holds for ev-

are obtained from atomic formulas  $E(x, y)$  (interpreted as the adjacency relation) and equality formulas  $x = y$  using Boolean combinations, existential quantification, and universal quantification; the quantifiers are interpreted as ranging over the set of vertices of the graph (and not over sets of vertices or sets of edges, etc.).

ery extension  $\text{FO}[\text{Mod}_q]$  of FO with the counting quantifiers  $\{\text{Mod}_q^i : i \in [q - 1]\}$ , where  $q$  is a prime. It should be noted that results in [9] imply that the modular convergence law for  $\text{FO}[\oplus]$  does *not* generalize to extensions of  $\text{FO}[\oplus]$  with fixed-point operators. This is in sharp contrast to the aforementioned 0-1 law for FO which carries over to extensions of FO with fixed-point operators.

## 1.1 Methods

Earlier 0-1 laws have been established by a combination of standard methods and techniques from mathematical logic and random graph theory. In particular, on the side of mathematical logic, the tools used include the compactness theorem, Ehrenfeucht-Fraïssé games, and quantifier elimination. Here, we establish the modular convergence law by combining quantifier elimination with, interestingly, algebraic methods related to multivariate polynomials over finite fields. In what follows in this section, we present an overview of the methods and techniques that we will use.

### 1.1.1 The distribution of subgraph frequencies mod $q$ , polynomials and Gowers norms

Let us briefly indicate the relevance of polynomials to the study of  $\text{FO}[\oplus]$  on random graphs. A natural example of a statement in  $\text{FO}[\oplus]$  is a formula  $\varphi$  such that  $G$  satisfies  $\varphi$  if and only if the number of copies of  $H$  in  $G$  is odd, for some graph  $H$  (where by copy we mean an induced subgraph, for now). Thus understanding the asymptotic probability of  $\varphi$  on  $G(n, p)$  amounts to understanding the distribution of the number of copies (mod 2) of  $H$  in  $G(n, p)$ .

In this spirit, we ask: what is the probability that in  $G(n, 1/2)$  there is an odd number of triangles (where we count *unordered* triplets of vertices  $\{a, b, c\}$  such that  $a, b, c$  are all pairwise adjacent<sup>2</sup>).

We reformulate this question in terms of the following “triangle polynomial”, that takes the adjacency matrix of a graph as input and returns the parity of the number of triangles in the graph;  $P_\Delta : \{0, 1\}^{\binom{n}{2}} \rightarrow \{0, 1\}$ , where

$$P_\Delta((x_e)_{e \in \binom{n}{2}}) = \sum_{\{e_1, e_2, e_3\} \text{ forming a } \Delta} x_{e_1} x_{e_2} x_{e_3},$$

where the arithmetic is mod 2. Note that for the random graph  $G(n, 1/2)$ , each entry of the adjacency matrix is chosen independently and uniformly from  $\{0, 1\}$ . Thus the probability that a random graph  $G \in G(n, 1/2)$  has an odd number of triangles is precisely equal to  $\text{Pr}_{x \in \mathbb{Z}_2^{\binom{n}{2}}} [P_\Delta(x) = 1]$ . Thus we have reduced our problem to studying the distribution of the evaluation of a certain polynomial at a random point, a topic of much study in pseudorandomness and algebraic coding theory, and we may now appeal to tools from these areas.

In Section 3, via the above approach, we show that the probability that  $G(n, 1/2)$  has an odd number of triangles equals  $1/2 \pm 2^{-\Omega(n)}$ . Similarly, for any connected graph  $F \neq K_1$  (the graph consisting of one vertex), the probability that  $G(n, 1/2)$  has an odd number of copies<sup>3</sup> of  $F$  is also  $1/2 \pm 2^{-\Omega(n)}$  (when  $F = K_1$ , there is no randomness in the

<sup>2</sup>Counting the number of *unordered* triples is not expressible in  $\text{FO}[\oplus]$ , we ask this question only for expository purposes (nevertheless, we do give an answer to this question in Section 3).

<sup>3</sup>with a certain precise definition of “copy”.

number of copies of  $F$  in  $G(n, 1/2)$ ). In fact, we show that for any collection of distinct connected graphs  $F_1, \dots, F_\ell$  ( $\neq K_1$ ), the joint distribution of the number of copies mod 2 of  $F_1, \dots, F_\ell$  in  $G(n, 1/2)$  is  $2^{-\Omega(n)}$ -close to the uniform distribution on  $\mathbb{Z}_2^\ell$ , i.e., the events that there are an odd number of  $F_i$  are essentially independent of one another.

Generalizing the above to  $G(n, p)$  and counting mod  $q$  for arbitrary  $p \in (0, 1)$  and arbitrary integers  $q$  motivates the study of new kinds of questions about polynomials, that we believe are interesting in their own right. For  $G(n, p)$  with arbitrary  $p$ , we need to study the distribution of  $P(x)$ , for certain polynomials  $P$ , when  $x \in \mathbb{Z}_2^m$  is distributed according to the  $p$ -biased measure. Even more interestingly, for the study of  $\text{FO}[\text{Mod}_q]$ , where we are interested in the distribution of the number of triangles mod  $q$ , one needs to understand the distribution of  $P(x)$  ( $P$  is now a polynomial over  $\mathbb{Z}_q$ ) where  $x$  is chosen uniformly from  $\{0, 1\}^m \subseteq \mathbb{Z}_q^m$  (as opposed to  $x$  being chosen uniformly from all of  $\mathbb{Z}_q^m$ , which is traditionally studied). In Section 4, we develop all the relevant polynomial machinery in order to answer these questions. This involves generalizing some classical results of Babai, Nisan and Szegedy [1] on correlations of polynomials. The key technical innovation here is our definition of a  $\mu$ -Gowers norm (where  $\mu$  is a measure on  $\mathbb{Z}_q^m$ ) that measures the correlation, under  $\mu$ , of a given function with low-degree polynomials (letting  $\mu$  be the uniform measure, we recover the standard Gowers norm). After generalizing several results about the standard Gowers norm to the  $\mu$ -Gowers norm case, we can then use a technique of Viola and Wigderson [19] to establish the generalization of [1] that we need.

### 1.1.2 Quantifier Elimination

Although we studied the distribution of subgraph frequencies mod  $q$  as an attempt to determine the limiting behavior of only a special family of  $\text{FO}[\text{Mod}_q]$  properties, it turns out that this case, along with the techniques developed to handle it, play a central role in the proof of the full modular convergence law. In fact, we reduce the modular convergence law for general  $\text{FO}[\text{Mod}_q]$  properties to the above case. We show that for any  $\text{FO}[\text{Mod}_q]$  sentence  $\varphi$ , with high probability over  $G \in G(n, p)$ , the truth of  $\varphi$  on  $G$  is determined by the number of copies in  $G$ , mod  $q$ , of each small subgraph. Then by the results described earlier on the equidistribution of these numbers (except for the number of  $K_1$ , which depends only on  $n \bmod q$ ), the full modular convergence law for  $\text{FO}[\text{Mod}_q]$  follows.

In Section 5, we establish such a reduction using the method of elimination of quantifiers. To execute this, we need to analyze  $\text{FO}[\text{Mod}_q]$  formulas which may contain free variables (i.e., not every variable used is quantified). Specifically, we show that for every  $\text{FO}[\text{Mod}_q]$  formula  $\varphi(\alpha_1, \dots, \alpha_k)$ , with high probability over  $G \in G(n, p)$ , it holds that for all vertices  $w_1, \dots, w_k$  of  $G$ , the truth of  $\varphi(w_1, \dots, w_k)$  is entirely determined by the following data: (a) which of the  $w_i, w_j$  pairs are adjacent, (b) which of the  $w_i, w_j$  pairs are equal to one another, and (c) the number of copies “rooted” at  $w_1, \dots, w_k$ , mod  $q$ , of each small *labelled graph*. This statement is a generalization of what we needed to prove, but lends itself to inductive proof (*this* is quantifier elimination). This leads us to studying the distribution (via the polynomial approach described earlier) of the number of copies of labelled graphs in  $G$ ; questions of the form, given two spec-

ified vertices  $v, w$  (the “roots”), what is the probability that there are an odd number of paths of length 4 in  $G \in G(n, p)$  from  $v$  to  $w$ ? After developing the necessary results on the distribution of labelled subgraph frequencies, combined with some elementary combinatorics, we can eliminate quantifiers and thus complete the proof of the modular convergence law.

## 1.2 Comparison with $\text{AC0}[\oplus]$

Every  $\text{FO}[\oplus]$  property naturally defines a family of boolean functions  $f_n : \{0, 1\}^{\binom{n}{2}} \rightarrow \{0, 1\}$ , such that a graph  $G$  satisfies  $\varphi$  if and only if  $f_n(A_G) = 1$ , where  $A_G$  is the adjacency matrix of  $G$ . This family of functions is easily seen to be contained in  $\text{AC0}[\oplus]$ , which is  $\text{AC0}$  with parity gates (each  $\forall$  becomes an  $\text{AND}$  gate,  $\exists$  becomes an  $\text{OR}$  gate and  $\oplus$  becomes a parity gate). This may be summarized by saying that  $\text{FO}[\oplus]$  is a highly uniform version of  $\text{AC0}[\oplus]$ .

Currently, all our understanding of the power of  $\text{AC0}[\oplus]$  comes from the Razborov-Smolensky [14, 16] approach to proving circuit lower bounds on  $\text{AC0}[\oplus]$ . At the heart of this approach is the result that for every  $\text{AC0}[\oplus]$  function  $f$ , there is a low-degree polynomial  $P$  such that for  $1 - \epsilon(n)$  fraction of inputs, the evaluations of  $f$  and  $P$  are equal. Note that this result automatically holds for  $\text{FO}[\oplus]$  (since  $\text{FO}[\oplus] \subseteq \text{AC0}[\oplus]$ ).

We show that for the special case when  $f : \{0, 1\}^{\binom{n}{2}} \rightarrow \{0, 1\}$  comes from an  $\text{FO}[\oplus]$  property  $\varphi$ , a significantly improved approximation may be obtained: (i) We show that the degree of  $P$  may be chosen to be a constant depending only on  $\varphi$ , whereas the Razborov-Smolensky approximation required  $P$  to be of  $\text{polylog}(n)$  degree, (ii) The error parameter  $\epsilon(n)$  may be chosen to be exponentially small in  $n$ , whereas the Razborov-Smolensky method only yields  $\epsilon(n) = 2^{-\log^{O(1)} n}$ . (iii): Finally, the polynomial  $P$  can be chosen to be symmetric under the action of  $S_n$  on the  $\binom{n}{2}$  coordinates, while in general, the polynomial produced by the Razborov-Smolensky approach need not be symmetric (due to the randomness involved in the choices).

These strengthened approximation results allow us, using known results about pseudorandomness against low-degree polynomials, to show that (i) there exist explicit pseudorandom generators that fool  $\text{FO}[\oplus]$  sentences, and (ii) there exist explicit functions  $f$  such that for any  $\text{FO}[\oplus]$  formula  $\varphi$ , the probability over  $G \in G(n, p)$  that  $f(G) = \varphi(G)$  is at most  $\frac{1}{2} + 2^{-\Omega(n)}$ . The first result follows from the pseudorandom generators against low-degree polynomials due to Bogdanov-Viola [3], Lovett [12] and Viola [18]. The second result follows from the result of Babai, Nisan and Szegedy [1], and our generalization of it, giving explicit functions that are uncorrelated with low degree polynomials.

Obtaining similar results for  $\text{AC0}[\oplus]$  is one of the primary goals of modern day “low-level” complexity theory.

### Organization of this paper:

In the next section, we formally state our main results and some of its corollaries. In Section 3, we determine the distribution of the number of copies mod  $q$  of all small connected graphs in  $G(n, p)$ . In Section 4, we introduce the  $\mu$ -Gowers Norm and use it to prove a lemma used in Section 3. In Section 5, we give a sketch of the proof of the full modular convergence law via quantifier elimination. We conclude with some open questions.

## 2. THE MODULAR CONVERGENCE LAW

We now state our main theorem.

**THEOREM 2.1.** *Let  $q$  be a prime. Then for every  $\varphi \in \text{FO}[\text{Mod}_q]$ , there exist rationals  $a_0, \dots, a_{q-1}$  such that for every  $p \in (0, 1)$  and every  $i \in \{0, 1, \dots, q-1\}$ ,*

$$\lim_{\substack{n \rightarrow \infty \\ n \equiv i \pmod q}} \Pr_{G \in G(n,p)} [G \text{ satisfies } \varphi] = a_i.$$

**REMARK** The proof of Theorem 2.1 also yields:

- Given a formula  $\varphi$ , the numbers  $a_0, \dots, a_{q-1}$  may be effectively determined.
- Each  $a_i$  is of the form  $r/q^s$ , where  $r, s$  are nonnegative integers.
- For every sequence of numbers  $b_0, \dots, b_{q-1} \in [0, 1]$ , each of the form  $r/q^s$ , there is a formula  $\varphi \in \text{FO}[\text{Mod}_q]$  such that for each  $i$ , the number  $a_i$  given by the theorem equals  $b_i$ .

Before we describe the main steps in the proof of Theorem 2.1, we make a few definitions.

For graphs  $F = (V_F, E_F)$  and  $G = (V_G, E_G)$ , an (injective) homomorphism from  $F$  to  $G$  is an (injective) map  $\chi : V_F \rightarrow V_G$  that maps edges to edges, i.e., for each  $(u, v) \in E_F$ , we have  $(\chi(u), \chi(v)) \in E_G$ . Note that we do not require that  $\chi$  maps non-edges to non-edges. We denote by  $[F](G)$  the number of injective homomorphisms from  $F$  to  $G$ , and we denote by  $[F]_q(G)$  this number mod  $q$ . We let  $\text{aut}(F) := [F](F)$  be the number of automorphisms of  $F$ .

The following lemma (whose proof is omitted in this version), shows that for some graphs  $F$ , as  $G$  varies, the number  $[F](G)$  cannot be arbitrary.

**LEMMA 2.2.** *Let  $F$  be a connected graph and  $G$  be any graph. Then  $\text{aut}(F) \mid [F](G)$ .*

For the rest of this section, let  $q$  be a fixed prime. Let  $\text{Conn}^a$  be the set of connected graphs on at most  $a$  vertices. For any graph  $G$ , let the *subgraph frequency vector*  $\text{freq}_G^a \in \mathbb{Z}_q^{\text{Conn}^a}$  be the vector such that its value in coordinate  $F$  ( $F \in \text{Conn}^a$ ) equals  $[F]_q(G)$ , the number of injective homomorphisms from  $F$  to  $G$  mod  $q$ . Let  $\text{FFreq}(a)$ , the set of *feasible frequency vectors*, be the subset of  $\mathbb{Z}_q^{\text{Conn}^a}$  consisting of all  $f$  such that  $f_F = 0$  whenever  $q \mid \text{aut}(F)$ . By lemma 2.2, for every  $G$  and  $a$ ,  $\text{freq}_G^a \in \text{FFreq}(a)$ , i.e., the subgraph frequency vector is always feasible.

We can now state the two main technical results that underlie Theorem 2.1.

The first states that on almost all graphs  $G$ , every  $\text{FO}[\text{Mod}_q]$  formula can be expressed in terms of the subgraph frequencies,  $[F]_q(G)$ , over all small connected graphs  $F$ .

**THEOREM 2.3.** (SUBGRAPH FREQUENCIES mod  $q$  DETERMINE  $\text{FO}[\text{Mod}_q]$  FORMULAE). *For every sentence  $\varphi \in \text{FO}[\text{Mod}_q]$  of quantifier depth  $t$ , there exists an integer  $c = c(t, q)$  and a function  $\psi : \mathbb{Z}_q^{\text{Conn}^c} \rightarrow \{0, 1\}$  such that for all  $p \in (0, 1)$ ,*

$$\Pr_{G \in G(n,p)} [(G \text{ satisfies } \varphi) \Leftrightarrow (\psi(\text{freq}_G^c) = 1)] \geq 1 - \exp(-n).$$

This result is complemented by the following result, that shows the distribution of subgraph frequencies mod  $q$  in a random graph  $G \in G(n, p)$  is essentially uniform in the space of all feasible frequency vectors, up to the obvious restriction that the number of vertices (namely the frequency of  $K_1$  in  $G$ ) should equal  $n \pmod q$ .

**THEOREM 2.4.** (DISTRIBUTION OF SUBGRAPH FREQUENCIES mod  $q$  DEPENDS ONLY ON  $n \pmod q$ ). *Let  $p \in (0, 1)$ . Let  $G \in G(n, p)$ . Then for any constant  $a$ , the distribution of  $\text{freq}_G^a$  is  $\exp(-n)$ -close to the uniform distribution over the set*

$$\{f \in \text{FFreq}(a) : f_{K_1} \equiv n \pmod q\}.$$

Theorem 2.4 is proved in Section 3 by studying the bias of multivariate polynomials over finite fields via a generalization of the Gowers' norm. Theorem 2.3 is proved in Section 5 using two main ingredients:

1. A generalization of Theorem 2.4 that determines the joint distribution of the frequencies of “labelled subgraphs” with given roots.
2. A variant of quantifier elimination (that may be called quantifier conversion) designed to handle  $\text{Mod}_q$  quantifiers that crucially uses the probabilistic input from the previous ingredient.

**PROOF OF THEOREM 2.1:** Follows by combining Theorem 2.3 and Theorem 2.4.  $\square$

### 2.1 Pseudorandomness against $\text{FO}[\oplus]$

We now point out three simple corollaries of our study of  $\text{FO}[\oplus]$  on random graphs. The second and third ones follow from the first using known results from the literature.

**COROLLARY 2.5.** ( $\text{FO}[\text{Mod}_q]$  IS WELL APPROXIMATED BY LOW-DEGREE POLYNOMIALS) *For every  $\varphi \in \text{FO}[\text{Mod}_q]$ , there is a constant  $d$ , such that for each  $n \in \mathbb{N}$ , there is a degree  $d$  polynomial  $P((X_e)_{e \in \binom{[n]}{2}}) \in \mathbb{Z}_q[(X_e)_{e \in \binom{[n]}{2}}]$ , such that for all  $p \in (0, 1)$ ,*

$$\Pr_{G \in G(n,p)} [\varphi(G) \text{ is true} \Leftrightarrow P(A_G) = 1] \geq 1 - 2^{-\Omega(n)},$$

where  $A_G \in \{0, 1\}^{\binom{[n]}{2}}$  is the adjacency matrix of  $G$ .

**PROOF.** Follows from Theorem 2.3 and the observation that for any graph  $F$  of constant size, there is a polynomial  $Q((X_e)_{e \in \binom{[n]}{2}})$  of constant degree, such that  $Q(A_G) = [F]_q(G)$  for all graphs  $G$ .  $\square$

**COROLLARY 2.6** (PRGS AGAINST  $\text{FO}[\oplus]$ ). *For each  $s \in \mathbb{N}$  and constant  $\epsilon > 0$ , there is a constant  $c \geq 0$  such that for each  $n$ , there is a family  $\mathcal{F}$  of  $\Theta(n^c)$  graphs on  $n$  vertices, computable in time  $\text{poly}(n^c)$ , such that for all  $\varphi \in \text{FO}[\oplus]$  of size at most  $s$ , and for all  $p \in (0, 1)$ ,*

$$\left| \Pr_{G \in \mathcal{F}} [G \text{ satisfies } \varphi] - \Pr_{G \in G(n,p)} [G \text{ satisfies } \varphi] \right| < \epsilon.$$

**PROOF.** For  $p = 1/2$ , this follows from the previous corollary and the result of Viola [18] (building on results of Bogdanov-Viola [3] and Lovett [12]) constructing a pseudorandom generator fooling low-degree polynomials under the uniform distribution. For general  $p$ , note that the same family  $\mathcal{F}$  from the  $p = 1/2$  case works, since the polynomial in the previous corollary is independent of  $p$ .  $\square$

**COROLLARY 2.7.** (EXPLICIT FUNCTIONS EXPONENTIALLY HARD FOR  $\text{FO}[\oplus]$ ) *There is an explicit function  $f : \{0, 1\}^{\binom{n}{2}} \rightarrow \{0, 1\}$  such that for every  $\text{FO}[\oplus]$  property  $\varphi$ ,*

$$\Pr_{G \in \mathcal{G}(n,p)} [(G \text{ satisfies } \varphi) \Leftrightarrow (f(A_G) = 1)] < \frac{1}{2} + 2^{-\Omega(n)}.$$

**PROOF.** Follows from Corollary 2.5, and the result of Babai, Nisan, Szegedy [1] (for  $p = 1/2$ ) and its generalization, Lemma 4.1 (for general  $p$ ), constructing functions exponentially uncorrelated with low degree polynomials under the  $p$ -biased measure. It actually follows from our proofs that, one may even choose a function  $f$  that is a graph property (namely, invariant under the action of  $S_n$  on the coordinates).  $\square$

### 3. SUBGRAPH FREQUENCIES

In this section, we prove Theorem 2.4 on the distribution of subgraph frequencies in  $G(n, p)$ .

We first make a few definitions. If  $F$  is a connected graph and  $G$  is any graph, a *copy* of  $F$  in  $G$  is a set  $E \subseteq E_G$  such that there exists an injective homomorphism  $\chi$  from  $F$  to  $G$  such that  $E = \chi(E_F) := \{(\chi(v), \chi(w)) \mid (v, w) \in E_F\}$ . We denote the set of copies of  $F$  in  $G$  by  $\text{Cop}(F, G)$ , the cardinality of  $\text{Cop}(F, G)$  by  $\langle F \rangle(G)$ , and this number mod  $q$  by  $\langle F \rangle_q(G)$ . We have the following basic relation (whose proof is omitted).

**LEMMA 3.1.** *If  $F$  is an connected graph with  $|E_F| \geq 1$ , then*

$$[F](G) = \text{aut}(F) \cdot \langle F \rangle(G).$$

We can now state the general equidistribution theorem from which Theorem 2.4 will follow easily (We use the notation  $\Omega_{q,p,d}(n)$  to denote the expression  $\Omega(n)$ , where the implied constant depends only on  $q, p$  and  $d$ ). Note that this theorem holds for arbitrary integers  $q$ , not necessarily prime.

**THEOREM 3.2** (EQUIDISTRIBUTION OF GRAPH COPIES). *Let  $q > 1$  be an integer and let  $p \in (0, 1)$ . Let  $F_1, \dots, F_\ell \in \text{Conn}^a$  be distinct graphs with  $1 \leq |E_{F_i}| \leq d$ .*

*Let  $G$  be a random graph distributed according to  $G(n, p)$ . Then the distribution of  $(\langle F_1 \rangle_q(G), \dots, \langle F_\ell \rangle_q(G))$  on  $\mathbb{Z}_q^\ell$  is  $2^{-\Omega_{q,p,d}(n)+\ell}$ -close to uniform in statistical distance.*

Using this theorem, we complete the proof of Theorem 2.4.

**PROOF OF THEOREM 2.4:** Let  $F_1, \dots, F_\ell$  be an enumeration of the elements of  $\text{Conn}^a$  except for  $K_1$ . By Theorem 3.2, the distribution of  $g = (\langle F_i \rangle_q(G))_{i=1}^\ell$  is  $2^{-\Omega(n)}$  close to uniform over  $\mathbb{Z}_q^\ell$ . Given the vector  $g$ , we may compute the vector  $\text{freq}_G^a$  by:

- $(\text{freq}_G^a)_{K_1} = n \pmod q$ .
- For  $F \in \text{Conn}^a \setminus \{K_1\}$ ,  $(\text{freq}_G^a)_F = g_F \cdot \text{aut}(F)$  (by Lemma 3.1).

This implies that the distribution of  $\text{freq}_G^a$  is  $2^{-\Omega(n)}$ -close to uniformly distributed over  $\{f \in \text{FFreq}(a) : f_{K_1} = n \pmod q\}$ .  $\square$

### 3.1 Preliminary lemmas

As indicated in the introduction, the distribution of subgraph frequencies is most naturally studied via the distribution of values of certain polynomials. The following lemma, which is used in the proof of Theorem 3.2 (and repeatedly throughout the proof of the modular convergence law), gives a simple sufficient criterion for the distribution of values of a polynomial to be “unbiased”, i.e., for the distribution of its values to be nearly uniform. The proof appears in Section 4.

**LEMMA 3.3.** *Let  $q > 1$  be an integer and let  $p \in (0, 1)$ . Let  $\mathcal{F} \subseteq 2^{[m]}$ . Let  $d > 0$  be an integer. Let  $Q(Z_1, \dots, Z_m) \in \mathbb{Z}_q[Z_1, \dots, Z_m]$  be a polynomial of the form*

$$\sum_{S \in \mathcal{F}} a_S \prod_{i \in S} Z_i + Q'(\mathbf{Z}),$$

where  $\deg(Q') < d$ . Suppose there exist  $\mathcal{E} = \{E_1, \dots, E_r\} \subseteq \mathcal{F}$  such that:

- $|E_j| = d$  for each  $j$ ,
- $a_{E_j} \neq 0$  for each  $j$ .
- $E_j \cap E_{j'} = \emptyset$  for each  $j, j'$ ,
- For each  $S \in \mathcal{F} \setminus \mathcal{E}$ ,  $|S \cap (\cup_j E_j)| < d$ .

Let  $\mathbf{z} = (z_1, \dots, z_m) \in \mathbb{Z}_q^m$  be the random variable where, independently for each  $i$ , we have  $\Pr[z_i = 1] = p$  and  $\Pr[z_i = 0] = 1 - p$ . Then,

$$|\mathbb{E}[\omega^{Q(\mathbf{z})}]] \leq 2^{-\Omega_{q,p,d}(r)},$$

where  $\omega \in \mathbb{C}$  is a primitive  $q^{\text{th}}$ -root of unity.

The lemma below is a useful tool for showing that a distribution on  $\mathbb{Z}_q^\ell$  is close to uniform.

**LEMMA 3.4** (VAZIRANI XOR LEMMA). *Let  $q > 1$  be an integer and let  $\omega \in \mathbb{C}$  be a primitive  $q^{\text{th}}$ -root of unity. Let  $\mathbf{X} = (X_1, \dots, X_\ell)$  be a random variable over  $\mathbb{Z}_q^\ell$ . Suppose that for every nonzero  $c \in \mathbb{Z}_q^\ell$ ,*

$$|\mathbb{E}[\omega^{\sum_{i \in [\ell]} c_i X_i}]] \leq \epsilon.$$

Then  $\mathbf{X}$  is  $q^\ell \cdot \epsilon$ -close to uniformly distributed over  $\mathbb{Z}_q^\ell$ .

### 3.2 Proof of the equidistribution theorem

**PROOF OF THEOREM 3.2:** By the Vazirani XOR Lemma (Lemma 3.4), it suffices to show that for each nonzero  $c \in \mathbb{Z}_q^\ell$ , we have  $|\mathbb{E}[\omega^R]| \leq 2^{-\Omega_{q,p,d}(n)}$ , where  $R := \sum_{i \in [\ell]} c_i \langle F_i \rangle_q(G)$ , and  $\omega \in \mathbb{C}$  is a primitive  $q^{\text{th}}$ -root of unity.

We will show this by appealing to Lemma 3.3. Let  $m = \binom{n}{2}$ . Let  $\mathbf{z} \in \{0, 1\}^{\binom{n}{2}}$ , be the random variable where, for each  $e \in \binom{[n]}{2}$ ,  $z_e = 1$  if and only if  $e$  is present in  $G$ . Thus, independently for each  $e$ ,  $\Pr[z_e = 1] = p$ .

We may now express  $R$  in terms of the  $z_e$ . Let  $K_n$  denote the complete graph on  $n$  vertices, and associate its vertices with the vertices of  $G$ . Thus  $\text{Cop}(F_i, K_n)$  is the set of  $E$  that could potentially arise as copies of  $F_i$  in  $G$ . Then we may write,

$$\begin{aligned} R &= \sum_{i \in [\ell]} c_i \langle F_i \rangle_q(G) = \sum_{i \in [\ell]} c_i \sum_{E \in \text{Cop}(F_i, K_n)} \prod_{e \in E} z_e \\ &= \sum_{E \in \mathcal{F}} c_E \prod_{e \in E} z_e, \end{aligned}$$

where  $\mathcal{F} \subseteq 2^{\binom{[n]}{2}}$  is the set  $\bigcup_{i:c_i \neq 0} \text{Cop}(F_i, K_n)$ , and for  $E \in \mathcal{F}$ ,  $c_E = c_i$  for the unique  $i$  satisfying  $E \in \text{Cop}(F_i, K_n)$  (note that since the  $F_i$  are nonisomorphic connected graphs, the  $\text{Cop}(F_i, K_n)$  are pairwise disjoint).

Let  $Q(\mathbf{Z}) \in \mathbb{Z}_q[\mathbf{Z}]$ , where  $\mathbf{Z} = (Z_e)_{e \in \binom{[n]}{2}}$  be the polynomial  $\sum_{E \in \mathcal{F}} c_E \prod_{e \in E} Z_e$ . Then  $R = Q(\mathbf{z})$ . We wish to show that

$$|\mathbb{E}[\omega^{Q(\mathbf{z})}]] \leq 2^{-\Omega_{q,p,d}(n)}. \quad (1)$$

We do this by demonstrating that the polynomial  $Q(\mathbf{Z})$  satisfies the hypotheses of Lemma 3.3.

Let  $d^* = \max_{i:c_i \neq 0} |E_{F_i}|$ . Let  $i_0 \in [\ell]$  be such that  $c_{i_0} \neq 0$  and  $|E_{F_{i_0}}| = d^*$ . Let  $\chi_1, \chi_2, \dots, \chi_r \in \text{Inj}(F_{i_0}, K_n)$  be a collection of homomorphisms such that for all distinct  $j, j' \in [r]$ , we have  $\chi_j(V_{F_{i_0}}) \cap \chi_{j'}(V_{F_{i_0}}) = \emptyset$ . Such a collection can be chosen greedily so that  $r = \Omega(\frac{n}{d})$ . Let  $E_j \in \text{Cop}(F_{i_0}, K_n)$  be given by  $\chi_j(E_{F_{i_0}})$ . Let  $\mathcal{E}$  be the family of sets  $\{E_1, \dots, E_r\} \subseteq \mathcal{F}$ . We observe the following properties of the  $E_j$ :

1. For each  $j \in [r]$ ,  $|E_j| = d^*$ .
2. For each  $j \in [r]$ ,  $c_{E_j} = c_{i_0} \neq 0$ .
3. For distinct  $j, j' \in [r]$ ,  $E_j \cap E_{j'} = \emptyset$ .
4. For every  $S \in \mathcal{F} \setminus \mathcal{E}$ ,  $|S \cap (\cup_j E_j)| < d^*$ . To see this, take any  $S \in \mathcal{F} \setminus \mathcal{E}$  and suppose  $|S \cap (\cup_j E_j)| \geq d^*$ . Let  $i' \in [\ell]$  be such that  $c_{i'} \neq 0$  and  $S \in \text{Cop}(F_{i'}, K_n)$ . Let  $\chi \in \text{Inj}(F_{i'}, K_n)$  with  $\chi(E_{F_{i'}}) = S$ . By choice of  $d^*$ , we know that  $|S| \leq d^*$ . Therefore, the only way that  $|S \cap (\cup_j E_j)|$  can be  $\geq d^*$  is if (1)  $|S| = d^*$ , and (2)  $S \cap (\cup_j E_j) = S$ , or in other words,  $S \subseteq (\cup_j E_j)$ . However, since the  $\chi_j(V_{F_{i_0}})$  are all pairwise disjoint, this implies that  $S \subseteq E_j$  for some  $j$ . But since  $|E_j| = |S|$ , we have  $S = E_j$ , contradicting our choice of  $S$ . Therefore,  $|S \cap (\cup_j E_j)| < d^*$  for any  $S \in \mathcal{F} \setminus \mathcal{E}$ .

It now follows that  $Q(\mathbf{Z})$ ,  $\mathcal{F}$  and  $\mathcal{E}$  satisfy the hypothesis of Lemma 3.3. Consequently, (recalling that  $r = \Omega(n/d)$  and  $d^* \leq d$ ) Equation (1) follows, completing the proof of the theorem.  $\square$

**REMARK** We just determined the joint distribution of the number of injective homomorphisms, mod  $q$ , from all small connected graphs to  $G(n, p)$ . This information can be used in conjunction with Lemma 5.6 to determine the joint distribution of the number of injective homomorphisms, mod  $q$ , from *all* small graphs to  $G(n, p)$ .

## 4. THE BIAS OF POLYNOMIALS

We now state and prove some useful lemmas about multivariate polynomials over finite fields and the distribution of their values. We will be especially interested in criteria for polynomials to be “unbiased”, Our main goal in this section will be to give a full proof of Lemma 3.3.

The following lemma, proved in the next subsection, shows that “Generalized Inner Product” polynomials are uncorrelated with polynomials of lower degree. It generalizes a result of Babai Nisan and Szegedy [1] (which dealt with the case  $q = 2$  and  $p = 1/2$ ). Our generalization is enabled by the use of the “ $\mu$ -Gowers norm”.

**LEMMA 4.1.** *Let  $q > 1$  be an integer and let  $p \in (0, 1)$ . Let  $E_1, \dots, E_r$  be pairwise disjoint subsets of  $[m]$  each of cardinality  $d$ . Let  $Q(Z_1, \dots, Z_m) \in \mathbb{Z}_q[Z_1, \dots, Z_m]$  be a polynomial of the form*

$$\left( \sum_{j=1}^r a_j \prod_{i \in E_j} Z_i \right) + R(\mathbf{Z}),$$

where each  $a_j \neq 0$  and  $\deg(R(\mathbf{Z})) < d$ . Let  $\mathbf{z} = (z_1, \dots, z_m) \in \mathbb{Z}_q^m$  be the random variable where, independently for each  $i$ , we have  $\Pr[z_i = 1] = p$  and  $\Pr[z_i = 0] = 1 - p$ . Then,

$$|\mathbb{E}[\omega^{Q(\mathbf{z})}]] \leq 2^{-\Omega_{q,p,d}(r)}.$$

Given Lemma 4.1, we may prove Lemma 3.3 below. Informally, it says that polynomials that have “Generalized Inner Product” polynomials embedded in them are unbiased.

**LEMMA 3.3 (RESTATED)** *Let  $q > 1$  be an integer and let  $p \in (0, 1)$ . Let  $\mathcal{F} \subseteq 2^{[m]}$ . Let  $d > 0$  be an integer. Let  $Q(Z_1, \dots, Z_m) \in \mathbb{Z}_q[Z_1, \dots, Z_m]$  be a polynomial of the form*

$$\sum_{S \in \mathcal{F}} a_S \prod_{i \in S} Z_i + Q'(\mathbf{Z}),$$

where  $\deg(Q') < d$ . Suppose there exist  $\mathcal{E} = \{E_1, \dots, E_r\} \subseteq \mathcal{F}$  such that:

- $|E_j| = d$  for each  $j$ ,
- $a_{E_j} \neq 0$  for each  $j$ .
- $E_j \cap E_{j'} = \emptyset$  for each  $j, j'$ ,
- For any  $S \in \mathcal{F} \setminus \mathcal{E}$ ,  $|S \cap (\cup_j E_j)| < d$ .

Let  $\mathbf{z} = (z_1, \dots, z_m) \in \mathbb{Z}_q^m$  be the random variable where, independently for each  $i$ , we have  $\Pr[z_i = 1] = p$  and  $\Pr[z_i = 0] = 1 - p$ . Then,

$$|\mathbb{E}[\omega^{Q(\mathbf{z})}]] \leq 2^{-\Omega_{q,p,d}(r)}.$$

**PROOF.** Let  $U = \cup_{j=1}^r E_j$ . Fix any  $x \in \{0, 1\}^{[m] \setminus U}$ , and let  $Q_x(\mathbf{Y}) \in \mathbb{Z}_q[(Y_i)_{i \in U}]$  be the polynomial

$$\sum_{S \in \mathcal{F}} a_S \left( \prod_{j \in S \cap ([m] \setminus U)} x_j \right) \left( \prod_{i \in S \cap U} Y_i \right) + Q'(x, \mathbf{Y})$$

so that  $Q_x(y) = Q(x, y)$  for each  $y \in \mathbb{Z}_q^U$ . Notice that the degree (in  $\mathbf{Y}$ ) of the term corresponding to  $S \in \mathcal{F}$  is  $|S \cap U|$ . By assumption, unless  $S = E_j$  for some  $j$ , we must have  $|S \cap U| < d$ .

Therefore the polynomial  $Q_x(\mathbf{Y})$  is of the form:

$$\sum_{j=1}^r a_{E_j} \prod_{i \in E_j} Y_i + R(\mathbf{Y}),$$

where  $\deg(R(\mathbf{Y})) < d$ . By Lemma 4.1,

$$|\mathbb{E}[\omega^{Q_x(\mathbf{y})}]] < 2^{-\Omega_{q,p,d}(r)},$$

where  $\mathbf{y} \in \{0, 1\}^U$  with each  $y_i = 1$  independently with probability  $p$ .

As  $Q_x(y) = Q(x, y)$ , we get

$$|\mathbb{E}[\omega^{Q(\mathbf{z}^x)}]] < 2^{-\Omega_{q,p,d}(r)},$$

where  $\mathbf{z}^x \in \mathbb{Z}_q^n$  is the random variable  $\mathbf{z}$  conditioned on the event  $z_j = x_j$  for every  $j \in [m] \setminus U$ . Now, the distribution of  $\mathbf{z}$  is a convex combination of the distributions of  $\mathbf{z}^x$  as  $x$  varies over  $\{0, 1\}^{[m] \setminus U}$ . This allows us to deduce that

$$|\mathbb{E}[\omega^{Q(\mathbf{z})}]| \leq 2^{-\Omega_{q,p,d}(r)},$$

as desired.  $\square$

## 4.1 $\mu$ -Gowers Norms

Before proving Lemma 4.1, we need to introduce a generalization of the Gowers norm and develop some of its basic properties.

Let  $H$  be an abelian group and let  $\mu$  be a probability distribution on  $H$ . For each  $d \geq 0$ , define a probability distribution  $\mu^{(d)}$  on  $H^{d+1}$  inductively by  $\mu^{(0)} = \mu$ , and, for  $d \geq 1$ , let  $\mu^{(d)}(x, t_1, \dots, t_d)$  equal

$$\frac{\mu^{(d-1)}(x, t_1, \dots, t_{d-1}) \cdot \mu^{(d-1)}(x + t_d, t_1, \dots, t_{d-1})}{\sum_{z \in H} \mu^{(d-1)}(z, t_1, \dots, t_{d-1})}.$$

Equivalently, to sample  $(x, t_1, \dots, t_d)$  from  $\mu^{(d)}$ , first take a sample  $(x, t_1, \dots, t_{d-1})$  from  $\mu^{(d-1)}$ , then take a sample  $(y, t'_1, \dots, t'_{d-1})$  from  $\mu^{(d-1)}$  conditioned on  $t'_i = t_i$  for each  $i \in [d-1]$ , and finally set  $t_d = y - x$  (our sample is then  $(x, t_1, \dots, t_{d-1}, t_d)$ ).

For a function  $f : H \rightarrow \mathbb{C}$  and  $\mathbf{t} \in H^d$ , we define its  $d^{\text{th}}$ -derivative in direction  $\mathbf{t}$  to be the function  $D_{\mathbf{t}}f : H \rightarrow \mathbb{C}$  given by

$$D_{\mathbf{t}}f(x) = \prod_{S \subseteq [d]} f(x + \sum_{i \in S} t_i)^{a^{|S|}},$$

where  $a^{\circ S}$  equals the complex conjugate  $\bar{a}$  if  $|S|$  is odd, and  $a^{\circ S}$  equals  $a$  otherwise. From the definition it immediately follows that  $D_{(\mathbf{t}, u)}f(x) = D_{\mathbf{t}}f(x) \overline{D_{\mathbf{t}}f(x + u)}$  (where  $(\mathbf{t}, u)$  denotes the vector  $(t_1, \dots, t_d, u) \in H^{d+1}$ ).

We now define the  $\mu$ -Gowers norm.

**DEFINITION 4.2** ( $\mu$ -GOWERS NORM). *If  $\mu$  is a distribution on  $H$ , and  $f : H \rightarrow \mathbb{C}$ , we define its  $(d, \mu)$ -Gowers norm by*

$$\|f\|_{U^d, \mu} = \left| \mathbb{E}_{(x, \mathbf{t}) \sim \mu^{(d)}} [(D_{\mathbf{t}}f)(x)] \right|^{\frac{1}{2^d}}.$$

When  $H$  is of the form  $\mathbb{Z}_q^m$ , then the  $(d, \mu)$ -Gowers norm of a function is supposed to estimate the correlation, under  $\mu$ , of that function with polynomials of degree  $d-1$ . Intuitively, this happens because the Gowers norm of  $f$  measures how often the  $d^{\text{th}}$  derivative of  $f$  vanishes. Notice that the value of this norm depends on the values of  $f$  at a random  $2^d$ -tuple of points arranged in a ‘‘cube’’, where the marginal distribution of each vertex of the cube is precisely  $\mu$ .

The next few lemmas enumerate some of the useful properties that  $\mu$ -Gowers norms enjoy.

**LEMMA 4.3.** *Let  $f : H \rightarrow \mathbb{C}$ . Then,*

$$|\mathbb{E}_{x \sim \mu} [f(x)]| \leq \|f\|_{U^d, \mu}.$$

**PROOF.** We prove that for every  $d$ ,  $\|f\|_{U^d, \mu} \leq \|f\|_{U^{d+1}, \mu}$ . The lemma follows by noting that  $\|f\|_{U^0, \mu} = |\mathbb{E}_{x \sim \mu} [f(x)]|$ .

The proof proceeds (following Gowers [7] and Green-Tao [8])

via the Cauchy-Schwarz inequality,

$$\begin{aligned} \|f\|_{U^d, \mu}^{2^{d+1}} &= \left| \mathbb{E}_{(x, \mathbf{t}) \sim \mu^{(d)}} [D_{\mathbf{t}}f(x)] \right|^2 \\ &\leq \mathbb{E}_{\mathbf{t}} \left[ \mathbb{E}_x [D_{\mathbf{t}}f(x)]^2 \right] && \text{by Cauchy-Schwarz} \\ &= \mathbb{E}_{\mathbf{t}} \mathbb{E}_{x, y} \left[ D_{\mathbf{t}}f(x) \overline{D_{\mathbf{t}}f(y)} \right] \end{aligned}$$

where  $y$  is an independent sample of  $x$  given  $\mathbf{t}$

$$\begin{aligned} &= \mathbb{E}_{x, \mathbf{t}, u} \left[ D_{\mathbf{t}}f(x) \overline{D_{\mathbf{t}}f(x + u)} \right] && \text{where } u = y - x \\ &= \mathbb{E}_{(x, \mathbf{t}, u) \sim \mu^{(d+1)}} \left[ D_{\mathbf{t}}f(x) \overline{D_{\mathbf{t}}f(x + u)} \right] \end{aligned}$$

by definition of  $\mu^{(d+1)}$

$$\begin{aligned} &= \mathbb{E}_{(x, \mathbf{t}, u) \sim \mu^{(d+1)}} [D_{\mathbf{t}, u}f(x)] \\ &= \|f\|_{U^{d+1}, \mu}^{2^{d+1}}. \end{aligned}$$

This proves the lemma.  $\square$

**DEFINITION 4.4.** *For each  $i \in [r]$ , let  $g_i : H \rightarrow \mathbb{C}$ . We define  $(\otimes_{i=1}^r g_i) : H^r \rightarrow \mathbb{C}$  by*

$$\left( \otimes_{i=1}^r g_i \right) (x_1, \dots, x_r) = \prod_{i=1}^r g_i(x_i).$$

*For each  $i \in [r]$ , let  $\mu_i$  a probability measure on  $H$ . We define the probability measure  $\otimes_{i=1}^r \mu_i$  on  $H^r$  by*

$$\left( \otimes_{i=1}^r \mu_i \right) (x_1, \dots, x_r) = \prod_{i=1}^r \mu_i(x_i).$$

**LEMMA 4.5.**  $\| \otimes_{i=1}^r g_i \|_{U^d, \otimes_{i=1}^r \mu_i} = \prod_{i=1}^r \|g_i\|_{U^d, \mu_i}$ .

**PROOF.** Follows by expanding both sides and using the fact that  $(\otimes_{i=1}^r \mu_i)^{(d)} = \otimes_{i=1}^r (\mu_i^{(d)})$ .  $\square$

**LEMMA 4.6.** *Let  $q > 1$  be an integer and let  $\omega \in \mathbb{C}$  be a primitive  $q^{\text{th}}$ -root of unity. For all  $f : \mathbb{Z}_q^n \rightarrow \mathbb{C}$ , all probability measures  $\mu$  on  $\mathbb{Z}_q^n$ , and all polynomials  $h \in \mathbb{Z}_q[Y_1, \dots, Y_n]$  of degree  $< d$ ,*

$$\|f\omega^h\|_{U^d, \mu} = \|f\|_{U^d, \mu}.$$

The above lemma follows from the fact that  $(D_{\mathbf{t}}f) = (D_{\mathbf{t}}(f \cdot \omega^h))$ .

**LEMMA 4.7.** *Let  $a \in \mathbb{Z}_q \setminus \{0\}$  and let  $g : \mathbb{Z}_q^d \rightarrow \mathbb{C}$  be given by  $g(y) = \omega^a \prod_{i=1}^d y_i$ . Let  $\mu$  be a probability distribution on  $\mathbb{Z}_q^d$  with  $\text{supp}(\mu) \supseteq \{0, 1\}^d$ . Then  $\|g\|_{U^d, \mu} < 1 - \epsilon$ , where  $\epsilon > 0$  depends only on  $q, d$  and  $\mu$ .*

We now put together the above ingredients.

**LEMMA 4.8.** *Let  $f : (\mathbb{Z}_q^d)^r \rightarrow \mathbb{C}$  be given by*

$$f(x_1, \dots, x_r) = \omega^{\sum_{j=1}^r a_j \prod_{i=1}^d x_{ij}},$$

where  $a_j \in \mathbb{Z}_q \setminus \{0\}$  for all  $j \in [r]$ . Let  $\mu$  be a probability distribution on  $\mathbb{Z}_q^d$  with  $\text{supp}(\mu) \supseteq \{0, 1\}^d$ . Then for all polynomials  $h \in \mathbb{Z}_q[(Y_{ij})_{i \in [d], j \in [r]}]$ , with  $\text{deg}(h) < d$ , we have

$$\left| \mathbb{E}_{x \sim \mu^{\otimes r}} [f(x) \omega^{h(x)}] \right| \leq c^r,$$

where  $c < 1$  depends only on  $q, d$  and  $\mu$ .

PROOF. Let  $g_j : \mathbb{Z}_q^d \rightarrow \mathbb{C}$  be given by  $g_j(y) = \omega^{a_j \prod_{i=1}^d y_i}$  (as in in Lemma 4.7), and take  $c = 1 - \epsilon$  from that Lemma. Notice that  $f = \otimes_{j=1}^r g_j$ . Therefore by Lemma 4.5, we have

$$\|f\|_{U^d, \mu^{\otimes r}} = \prod_{j=1}^r \|g_j\|_{U^d, \mu} < c^r.$$

As the degree of  $h$  is at most  $d - 1$ , Lemma 4.6 implies that

$$\|f\omega^h\|_{U^d, \mu^{\otimes r}} = \|f\|_{U^d, \mu^{\otimes r}} = c^r.$$

Lemma 4.3 now implies that

$$\left| \mathbb{E}_{x \sim \mu^{\otimes r}} \left[ f(x) \omega^{h(x)} \right] \right| \leq c^r,$$

as desired.  $\square$

We can now complete the proof of Lemma 4.1.

PROOF OF LEMMA 4.1: By fixing the variables  $Z_i$  for  $i \notin \cup_j E_j$ , and then averaging over all such fixings, it suffices to consider the case  $[m] = \cup_j E_j$ . Then the polynomial  $Q(Z_1, \dots, Z_m) = \left( \sum_{j=1}^r a_j \prod_{i \in E_j} Z_i \right) + R(Z)$  can be rewritten in the form (after renaming the variables):

$$\sum_{j=1}^r a_j \prod_{i=1}^d X_{ij} + h(\mathbf{X}),$$

where  $\deg(h) < d$ . Let  $\mu$  be the  $p$ -biased probability measure on  $\{0, 1\}^d \subseteq \mathbb{Z}_q^d$ . Lemma 4.8 now implies that

$$\left| \mathbb{E}_{x \sim \mu^{\otimes r}} \left[ \omega^{Q(x)} \right] \right| \leq 2^{-\Omega_{q,p,d}(r)},$$

as desired.  $\square$

## 5. QUANTIFIER ELIMINATION

In this section, we sketch a proof of Theorem 2.3. The full proof is deferred to the full version of the paper. We begin by introducing some notions that will be needed to state our main technical theorem, Theorem 5.2.

### 5.1 Labelled graphs, types and frequency vectors

We will need to generalize some concepts related to graphs and homomorphisms to *labelled graphs*. In the full version of the paper, we prove various equidistribution theorems similar to Theorem 3.2, for the number of copies of small labelled graphs in  $G(n, p)$ .

DEFINITION 5.1. *An  $I$ -labelled graph is a graph  $F = (V_F, E_F)$  where some vertices are labelled by elements of  $I$ , such that (a) for each  $i \in I$ , there is exactly one vertex labelled  $i$ . We denote this vertex  $F(i)$ , and (b) the graph induced on the set of labelled vertices is an independent set. We denote the set of labelled vertices of  $F$  by  $\mathcal{L}(F)$ .*

A *homomorphism* from an  $I$ -labelled graph  $F$  to a pair  $(G, \mathbf{w})$ , where  $G$  is a graph and  $\mathbf{w} \in V_G^I$ , is defined to be a homomorphism  $\chi \in \text{Hom}(F, G)$  such that for each  $i \in I$ ,  $\chi$  maps  $F(i)$  to  $w_i$ . A homomorphism from  $F$  to  $(G, \mathbf{w})$  is called *injective* if for all distinct  $v, w \in V_F$ , such that  $\{v, w\} \not\subseteq \mathcal{L}(F)$ , we have  $\chi(v) \neq \chi(w)$ . We define  $\text{Hom}(F, (G, \mathbf{w}))$  (respectively  $\text{Inj}(F, (G, \mathbf{w}))$ ) to be the set of homomorphisms (respectively injective homomorphisms)

from  $F$  to  $(G, \mathbf{w})$ . We define  $[F](G, \mathbf{w})$  to be the cardinality of  $\text{Inj}(F, (G, \mathbf{w}))$ , and define  $[F]_q(G, \mathbf{w}) = [F](G, \mathbf{w}) \bmod q$ .

For  $F$  an  $I$ -labelled graph, we say  $F$  is *label-connected* if  $F \setminus \mathcal{L}(F)$  is connected. Define  $\text{Conn}_I^t$  to be the set of all  $I$ -labelled label-connected graphs with at most  $t$  unlabelled vertices.

REMARK We will often deal with  $[k]$ -labelled graphs. By abuse of notation we will refer to them as  $k$ -labelled graphs. If  $\mathbf{w} \in V^{[k]}$  and  $u \in V$ , when we refer to the tuple  $(\mathbf{w}, v)$  we mean the  $[k + 1]$ -tuple whose first  $k$  coordinates are given by  $\mathbf{w}$  and whose  $k + 1$ st coordinate is  $v$ . Similarly  $\text{Conn}_k^t$  denotes  $\text{Conn}_{[k]}^t$ .

For a graph  $G$  and vertices  $w_1, \dots, w_k \in V_G$ , we define the *type* of  $\mathbf{w} = (w_1, \dots, w_k)$ , denoted  $\text{type}_G(\mathbf{w})$ , to be the aggregate of the equality and adjacency data, namely:  $\Pi_\tau = \{(i, j) \in [k]^2 \mid w_i = w_j\}$ , and  $E_\tau = \{(i, j) \in [k]^2 \mid w_i \text{ adjacent to } w_j\}$ . The collection of all possible types of  $k$  vertices is denoted  $\text{Types}(k)$ .

For a graph  $G$ , a tuple  $\mathbf{w} \in V_G^k$  and a positive integer  $a$ , we define the *labelled subgraph frequency vector at  $\mathbf{w}$*   $\text{freq}_G^a(\mathbf{w}) \in \mathbb{Z}_q^{\text{Conn}_k^a}$  to be the vector such that for each  $F \in \text{Conn}_k^a$ ,

$$(\text{freq}_G^a(\mathbf{w}))_F = [F]_q(G, \mathbf{w}).$$

### 5.2 Sketch of the main quantifier elimination step

The proof of Theorem 2.3 will be via a more general theorem amenable to inductive proof, Theorem 5.2. Just as Theorem 2.3 states that for almost all  $G \in G(n, p)$ , the truth of any  $\text{FO}[\text{Mod}_q]$  sentence on  $G$  is determined by subgraph frequencies,  $\text{freq}_G^c$ , Theorem 5.2 states that for almost all graphs  $G \in G(n, p)$ , for any  $w_1, \dots, w_k \in V_G$  the truth of any  $\text{FO}[\text{Mod}_q]$  formula  $\varphi(w_1, \dots, w_k)$  on  $G$  is determined by the type of  $\mathbf{w}$ ,  $\text{type}_G(\mathbf{w})$ , and the labelled subgraph frequencies at  $\mathbf{w}$ ,  $\text{freq}_G^c(\mathbf{w})$ .

THEOREM 5.2. *For all primes  $q$  and integers  $k, t > 0$ , there is a constant  $c = c(k, t, q)$  such that for every  $\text{FO}[\text{Mod}_q]$  formula  $\varphi(\alpha_1, \dots, \alpha_k)$  with quantifier depth  $t$ , there is a function  $\psi : \text{Types}(k) \times \mathbb{Z}_q^{\text{Conn}_k^c} \rightarrow \{0, 1\}$  such that for all  $p \in (0, 1)$ ,*

$$\Pr_{G \in G(n, p)} \left[ \begin{array}{l} \forall w_1, \dots, w_k \in V_G \\ ((G \text{ satisfies } \varphi(w_1, \dots, w_k)) \Leftrightarrow \psi(\text{type}_G(\mathbf{w}), \text{freq}_G^c(\mathbf{w})) = 1) \\ \geq 1 - 2^{-\Omega(n)}. \end{array} \right]$$

Putting  $k = 0$ , we recover Theorem 2.3.

This theorem is proved by induction on the structure of the formula  $\varphi$ . When the formula  $\varphi$  has no quantifiers, then the truth of  $\varphi(\mathbf{w})$  on  $G$  is completely determined by  $\text{type}_G(\mathbf{w})$ . The case where  $\varphi$  is of the form  $\varphi_1(\alpha_1, \dots, \alpha_k) \wedge \varphi_2(\alpha_1, \dots, \alpha_k)$  is easily handled via the induction hypothesis. The case where  $\varphi(\alpha_1, \dots, \alpha_k) = \neg \varphi_1(\alpha_1, \dots, \alpha_k)$  is similar.

The key cases for us to handle are thus (i)  $\varphi(\alpha_1, \dots, \alpha_k)$  is of the form  $\text{Mod}_q^t \beta, \varphi'(\alpha_1, \dots, \alpha_k, \beta)$ , and (ii)  $\varphi(\alpha_1, \dots, \alpha_k)$  is of the form  $\exists \beta, \varphi'(\alpha_1, \dots, \alpha_k, \beta)$ . We now give a sketch of how these cases may be handled.

For case (i), let  $\psi' : \text{Types}(k + 1) \times \mathbb{Z}_q^{\text{Conn}_{k+1}^c}$  be the function given by the induction hypothesis for the formula  $\varphi'$ .



Thus for most graphs  $G \in G(n, p)$  (namely the ones for which  $\psi'$  is good for  $\varphi'$ ),  $\varphi(w_1/\alpha_1, \dots, w_k/\alpha_k)$  is true if and only the number of vertices  $v \in V_G$  such that  $(\tau'_v, f'_v) := (\text{type}_G(\mathbf{w}, v), \text{freq}_G^b(\mathbf{w}, v))$  such that  $\psi'(\tau'_v, f'_v) = 1$  is congruent to  $i \bmod q$ . Using the following lemma, proved in the next subsection, we show that the number of such vertices  $v$  can be determined solely as a function of  $\text{type}_G(\mathbf{w})$  and  $\text{freq}_G^a(\mathbf{w})$  for suitable  $a$ . This allows us to define  $\psi$ , and the proof of case (i) is complete.

LEMMA 5.3. *Let  $q$  be a prime, let  $k, b > 0$  be integers and let  $a \geq (q-1) \cdot b \cdot |\text{Conn}_{k+1}^b|$ . There is a function*

$$\lambda : \text{Types}(k+1) \times \mathbb{Z}_q^{\text{Conn}_{k+1}^b} \times \text{Types}(k) \times \mathbb{Z}_q^{\text{Conn}_k^a} \rightarrow \mathbb{Z}_q$$

such that for all  $\tau' \in \text{Types}(k+1)$ ,  $f' \in \mathbb{Z}_q^{\text{Conn}_{k+1}^b}$ ,  $\tau \in \text{Types}(k)$ ,  $f \in \mathbb{Z}_q^{\text{Conn}_k^a}$ , it holds that for every graph  $G$ , and every  $w_1, \dots, w_k \in V_G$  with  $\text{type}_G(\mathbf{w}) = \tau$  and  $\text{freq}_G^a(\mathbf{w}) = f$ , the cardinality of the set

$$\{v \in V_G : \text{type}_G(\mathbf{w}, v) = \tau' \wedge \text{freq}_G^b(\mathbf{w}, v) = f'\}$$

is congruent to  $\lambda(\tau', f', \tau, f) \bmod q$ .

Case (ii) is the most technically involved case. As before, we get a function  $\psi'$  corresponding to  $\varphi'$  by the induction hypothesis. We show that one can define  $\psi$  essentially as follows: define  $\psi(\tau, f) = 1$  if there exists some  $(\tau', f') \in \text{Types}(k+1) \times \mathbb{Z}_q^{\text{Conn}_{k+1}^b}$  that “extends”  $(\tau, f)$  for which  $\psi'(\tau', f') = 1$ ; otherwise  $\psi(\tau, f) = 0$ . Informally, we show that if it is conceivable that there is a vertex  $v$  such that  $\varphi'(\mathbf{w}, v)$  is true, then  $\varphi(\mathbf{w})$  is almost surely true. Proving this statement requires us to develop a number of probabilistic results on the distribution of labelled subgraph frequencies (generalizing Theorem 2.4). Thus we handle the case of the  $\exists$  quantifier. This completes the proof of Theorem 5.2.

### 5.3 Counting extensions

In this subsection we prove Lemma 5.3.

We begin with a definition. A *partial matching* between two  $I$ -labelled graphs  $F_1, F_2$  is a subset  $\eta \subseteq (V_{F_1} \setminus \mathcal{L}(F_1)) \times (V_{F_2} \setminus \mathcal{L}(F_2))$  that is one-to-one. For two graphs  $F_1, F_2$ , let  $\text{PMatch}(F_1, F_2)$  be the set of all partial matchings between them.

DEFINITION 5.4. *Let  $F_1$  and  $F_2$  be two  $I$ -labelled graphs, and let  $\eta \in \text{PMatch}(F_1, F_2)$ . Define  $F_1 \vee_\eta F_2$  to be the graph obtained by first taking the disjoint union of  $F_1$  and  $F_2$ , identifying pairs of vertices with the same label, and then identifying the vertices in each pair of  $\eta$  (and removing duplicate edges). We omit the subscript when  $\eta = \emptyset$ .*

We have the following simple identity (whose proof is omitted).

LEMMA 5.5. *For any  $I$ -labelled graphs  $F_1, F_2$ , any graph  $G$  and any  $\mathbf{w} \in V_G^I$ :*

$$[F_1](G, \mathbf{w}) \cdot [F_2](G, \mathbf{w}) = \sum_{\eta \in \text{PMatch}(F_1, F_2)} [F_1 \vee_\eta F_2](G, \mathbf{w}). \quad (2)$$

LEMMA 5.6. (CONNECTED SUBGRAPH FREQUENCIES DETERMINE ALL SUBGRAPH FREQUENCIES) *For every  $k$ -labelled*

*graph  $F'$  with  $|V_{F'}| = t$ , there is a polynomial  $\delta_{F'}(\mathbf{X}) \in \mathbb{Z}[(X_F)_{F \in \text{Conn}_k^t}]$  such that for all graphs  $G$  and  $\mathbf{w} \in V_G^k$ ,*

$$[F'](G, \mathbf{w}) = \delta_{F'}(x),$$

where  $x \in \mathbb{Z}^{\text{Conn}_k^t}$  is given by  $x_F = [F](G, \mathbf{w})$ .

PROOF. By induction on the number of connected components of  $F' \setminus \mathcal{L}(F')$ . If  $F'$  is label-connected, then we take  $\delta_{F'}(\mathbf{X}) = X_{F'}$ .

Now suppose  $F'$  is label-disconnected. Write  $F' = F_1 \vee F_2$  where  $F_1$  and  $F_2$  are both  $k$ -labelled graphs, and  $F_1 \setminus \mathcal{L}(F_1)$  and  $F_2 \setminus \mathcal{L}(F_2)$  have fewer connected components.

By equation (2), for all  $G$  and  $\mathbf{w}$ ,

$$\begin{aligned} [F_1 \vee F_2](G, \mathbf{w}) &= [F_1](G, \mathbf{w}) \cdot [F_2](G, \mathbf{w}) \\ &- \sum_{\emptyset \neq \eta \in \text{PMatch}(F_1, F_2)} [F_1 \vee_\eta F_2](G, \mathbf{w}). \end{aligned}$$

Observe that for any  $\eta \neq \emptyset$ , each graph  $F_1 \vee_\eta F_2$  has at least one fewer label-connected component than  $F_1 \vee F_2 = F'$ . Thus, by induction hypothesis, we may take

$$\delta_{F'}(\mathbf{X}) = \delta_{F_1}(\mathbf{X}) \cdot \delta_{F_2}(\mathbf{X}) - \sum_{\emptyset \neq \eta \in \text{PMatch}(F_1, F_2)} \delta_{F_1 \vee_\eta F_2}(\mathbf{X}).$$

This completes the proof of the lemma.  $\square$

We can now prove Lemma 5.3.

PROOF OF LEMMA 5.3: We describe the function  $\lambda(\tau', f', \tau, f)$  explicitly. If  $\tau'$  does not “extend”  $\tau$ , then we set  $\lambda(\tau', f', \tau, f) = 0$  (the notion of extend, which can be defined precisely, roughly says that the adjacency and equality information in  $\tau'$  restricted to the first  $k$  vertices, is consistent with the adjacency and equality information in  $\tau$ ).

Now assume  $\tau'$  extends  $\tau$ . We take cases on whether the  $k+1^{\text{st}}$  vertex in  $\tau'$  equals any of the first  $k$  vertices or not.

**Case 1:** For all  $j \in [k]$ ,  $(j, k+1) \notin \Pi_{\tau'}$ . In this case, there is an  $I \subseteq [k]$  such that  $\text{type}_G(w_1, \dots, w_k, v) = \tau'$  if and only if  $v \notin \{w_1, \dots, w_k\}$  and  $(v, w_i) \in E_G \Leftrightarrow i \in I$ .

For vertices  $u, u'$  in  $V_G$ , let  $x_{uu'} \in \{0, 1\}$  equal 1 if and only if  $(u, u') \in E_G$ . Then the number (mod  $q$ ) of  $v$  with  $\text{type}_G(\mathbf{w}, v) = \tau'$  and  $\text{freq}_G(\mathbf{w}, v) = f'$  can be compactly expressed as (and *this* is where we use the primality of  $q$ ):

$$\sum_{v \in V_G \setminus \{w_1, \dots, w_k\}} \left( \prod_{i \in I} x_{vw_i} \right) \cdot \left( \prod_{j \in [k] \setminus I} (1 - x_{vw_j}) \right) \cdot \prod_{F \in \text{Conn}_{k+1}^b} \left( 1 - ([F]_q(G, \mathbf{w}, v) - f'_F)^{q-1} \right).$$

Expanding, the expression  $\prod_{i \in I} x_{vw_i} \prod_{j \in [k] \setminus I} (1 - x_{vw_j})$  may be rewritten in the form  $\sum_{S \subseteq [k]} b_S \prod_{i \in S} x_{vw_i}$ . Lemma 5.5 implies that the expression

$$\prod_{F \in \text{Conn}_{k+1}^b} \left( 1 - ([F]_q(G, \mathbf{w}, v) - f'_F)^{q-1} \right)$$

may be rewritten in the form  $\sum_j c_j [F_j]_q(G, \mathbf{w}, v)$ , where each  $F_j$  is a  $k+1$ -labelled graph with at most  $|\text{Conn}_{k+1}^b| \cdot b \cdot (q-1) \leq a$  unlabelled vertices.

Thus we may rewrite the expression for  $\lambda(\tau', f', \tau, f)$  as:

$$\begin{aligned} & \sum_{v \in [n] \setminus \{w_1, \dots, w_k\}} \left( \sum_S b_S \prod_{i \in S} x_{vw_i} \right) \left( \sum_j c_j [F_j]_q(G, \mathbf{w}, v) \right) \\ &= \sum_{S,j} b_S c_j \sum_{v \in [n] \setminus \{w_1, \dots, w_k\}} \left( \left( \prod_{i \in S} x_{vw_i} \right) [F_j]_q(G, \mathbf{w}, v) \right) \\ &= \sum_{S,j} b_S c_j [F'_{S,j}]_q(G, \mathbf{w}), \end{aligned}$$

where  $F'_{S,j}$  is the  $k$ -labelled graph obtained from  $F_j$  by

- (a) For each  $i \in S$ , adding an edge between the vertex labelled  $k+1$  and the vertex labelled  $i$ , and
- (b) Removing the label from the vertex labelled  $k+1$ .

Finally, note that by Lemma 5.6,  $[F'_{S,j}]_q(G, \mathbf{w})$  is determined by  $\text{freq}_G^a(\mathbf{w})$ .

**Case 2:** There is some  $j \in [k]$  such that  $(j, k+1) \in \Pi_{\tau'}$ . This case is much easier to handle. Then there is only one  $v \in V_G$  such that  $\text{type}_G(\mathbf{w}, v) = \tau'$  (namely,  $w_j$ ).

Then  $\lambda(\tau', f', \tau, f) = 1$  if and only if for all  $F' \in \text{Conn}_{k+1}^b$ ,  $f'_{F'} = f_F$ , where  $F \in \text{Conn}_k^b$  is the graph obtained by identifying the vertex labelled  $k+1$  with the vertex labelled  $j$ , and labelling this new vertex  $j$ . Otherwise  $\lambda(\tau', f', \tau, f) = 0$ .

This completes the definition of our desired function  $\lambda$ .  $\square$

## 6. CONCLUDING REMARKS

The results presented here constitute the first systematic investigation of the asymptotic probabilities of properties expressible in first-order logic with counting quantifiers. Moreover, these results have been established by combining, for the first time, algebraic methods related to multivariate polynomials over finite fields with the method of quantifier elimination from mathematical logic.

We conclude with two open problems:

1. What is the complexity of computing the numbers  $a_0, \dots, a_{q-1}$  in Theorem 2.1? We know that computing these numbers is PSPACE-hard (it is already PSPACE-hard to tell if the asymptotic probability of a FO sentence is 0 or 1). Our proof shows that they may be computed in time  $2^{2^{2^{\dots}}}$  of height proportional to the quantifier depth of the formula. It is likely that a more careful analysis of our approximation of  $\text{FO}[\text{Mod}_q]$  by polynomials can yield better upper bounds.
2. Is there a modular convergence law for  $\text{FO}[\text{Mod}_m]$  for arbitrary  $m$ ? We encounter the same obstacles that prevent the Razborov-Smolensky approach from generalizing to  $\text{AC}_0[\text{Mod}_6]$ . Perhaps an answer to the above question will give us some hints for  $\text{AC}_0[\text{Mod}_6]$ ?

## Acknowledgements

Swastik Kopparty is very grateful to Eli Ben-Sasson, Danny Gutfreund and Alex Samorodnitsky for encouragement and stimulating discussions. We would also like to thank Miki Ajtai, Ron Fagin, Prasad Raghavendra, Ben Rossman, Shubhangi Saraf and Madhu Sudan for valuable discussions.

## 7. REFERENCES

- [1] Babai, Nisan, and Szegedy. Multipart protocols and logspace-hard pseudorandom sequences. In *STOC: ACM Symposium on Theory of Computing (STOC)*, 1989.
- [2] A. Blass, Y. Gurevich, and D. Kozen. A zero-one law for logic with a fixed point operator. *Information and Control*, 67:70–90, 1985.
- [3] A. Bogdanov and E. Viola. Pseudorandom bits for polynomials. In *FOCS*, pages 41–51, 2007.
- [4] R. Fagin. Generalized first-order spectra and polynomial-time recognizable sets. In R. M. Karp, editor, *Complexity of Computation, SIAM-AMS Proceedings, Vol. 7*, pages 43–73, 1974.
- [5] R. Fagin. Probabilities on finite models. *Journal of Symbolic Logic*, 41:50–58, 1976.
- [6] Y. V. Glebskii, D. I. Kogan, M. I. Liogonki, and V. A. Talanov. Range and degree of realizability of formulas in the restricted predicate calculus. *Cybernetics*, 5:142–154, 1969.
- [7] W. T. Gowers. A new proof of Szemerédi’s theorem. *Geom. Funct. Anal.*, 11(3):465–588, 2001.
- [8] B. Green and T. Tao. The primes contain arbitrarily long arithmetic progressions. *Ann. of Math. (2)*, 167(2):481–547, 2008.
- [9] L. Hella, P. Kolaitis, and K. Luosto. Almost everywhere equivalence of logics in finite model theory. *Bulletin of Symbolic Logic*, 2(4):422–443, 1996.
- [10] P. G. Kolaitis and M. Y. Vardi. The decision problem for the probabilities of higher-order properties. In *Proc. 19th ACM Symp. on Theory of Computing*, pages 425–435, 1987.
- [11] P. G. Kolaitis and M. Y. Vardi. 0-1 laws and decision problems for fragments of second-order logic. *Information and Computation*, 87:302–338, 1990.
- [12] S. Lovett. Unconditional pseudorandom generators for low degree polynomials. In *STOC*, pages 557–562, 2008.
- [13] L. Pacholski and W. Szwoz. The 0-1 law fails for the class of existential second-order Gödel sentences with equality. In *Proc. 30th IEEE Symp. on Foundations of Computer Science*, pages 280–285, 1989.
- [14] Razborov. Lower bounds on the size of bounded depth circuits over a complete basis with logical addition. *MATHNASUSSR: Mathematical Notes of the Academy of Sciences of the USSR*, 41, 1987.
- [15] S. Shelah and J. Spencer. Zero-one laws for sparse random graphs. *J. Amer. Math. Soc.*, 1:97–115, 1988.
- [16] R. Smolensky. Algebraic methods in the theory of lower bounds for boolean circuit complexity. In *STOC*, pages 77–82, 1987.
- [17] J. Spencer and S. Shelah. Threshold spectra for random graphs. In *Proc. 19th ACM Symp. on Theory of Computing*, pages 421–424, 1987.
- [18] E. Viola. The sum of  $d$  small-bias generators fools polynomials of degree  $d$ . In *IEEE Conference on Computational Complexity*, pages 124–127, 2008.
- [19] E. Viola and A. Wigderson. Norms, xor lemmas, and lower bounds for  $\text{gf}(2)$  polynomials and multipart protocols. In *22th IEEE Conference on Computational Complexity (CCC)*, 2007.