# Parity check matrices and product representations of squares

Assaf Naor
Microsoft Research
anaor@microsoft.com

Jacques Verstraëte
University of Waterloo
jverstra@math.uwaterloo.ca

**Abstract**

Let $N_{\mathbb{F}}(n,k,r)$ denote the maximum number of columns in an $n$-row matrix with entries in a finite field $\mathbb{F}$ in which each column has at most $r$ nonzero entries and every $k$ columns are linearly independent over $\mathbb{F}$. We obtain near-optimal upper bounds for $N_{\mathbb{F}}(n,k,r)$ in the case $k > r$. Namely, we show that $N_{\mathbb{F}}(n,k,r) \ll n^{\frac{r}{2}+\frac{cr}{k}}$ where $c \approx \frac{4}{3}$ for large $k$. Our method is based on a novel reduction of the problem to the extremal problem for cycles in graphs, and yields a fast algorithm for finding short linear dependences. We present additional applications of this method to problems in extremal hypergraph theory and combinatorial number theory.

*"It is certainly odd to have an instruction in an algorithm asking you to play with some numbers to find a subset with product a square . . . Why should we expect to find such a subsequence, and, if it exists, how can we find it efficiently?"*

**Carl Pomerance [32].**

## 1 Introduction

Since the mid-1930s it has been well established that there is a tight connection between combinatorial number theory and extremal combinatorics (see [33, 12] and the references therein). The basic paradigm is that for certain number theoretical problems one can construct a combinatorial object (e.g. a graph or a hypergraph), and prove that it cannot contain certain configurations (e.g. cycles). Thus, in many cases one can use theorems on excluded configurations in extremal combinatorics to bound the size of the combinatorial structure that was constructed, and this translates back to give number theoretical consequences. The present paper develops a novel reduction of this type, and applies it to two problems in algebraic combinatorics, namely to coding theory and combinatorial number theory. Additionally, our proofs are constructive, and thus yield the best known algorithms for several natural computational problems.

Our original motivation comes from a problem in coding theory. Low density parity check codes were introduced by Gallager in the 1960s, and have since found numerous theoretical and practical applications in engineering and computer science (see [22, 29, 30] for an account of this theory. We also refer to [11] for a nice introduction to the geometry of codes). Given a linear code $C \subseteq \mathbb{F}_2^m$ of dimension $\ell$ and minimum Hamming weight $t$, an $(m - \ell) \times m$ matrix $H$ is called a *parity check matrix* of the code $C$ if $C = \{v \in \mathbb{F}_2^m : Hv = 0\}$. We shall say that $H$ is $r$-sparse if every column of $H$ has at most $r$ non-zero entries. The *Syndrome Decoding Algorithm* for such codes works as follows: given a corrupted signal $z$ one computes the vector $x$ of minimal weight

satisfying $Hz = Hx$, and decodes $z$ to $z - x$ (this algorithm corrects at most $t/2$ errors). As such computations are faster if the sparseness of $H$ is exploited, it is desirable to obtain codes with sparse parity check matrices. Indeed, sparse parity check matrices occur in many of the known constructions of codes, e.g. codes based on bounded degree graphs such as expander codes [35, 36], and we also refer to [28] for theoretical and experimental coding theory applications of very sparse matrices (we stress here that the present paper deals with a different range of parameters – our bounds will be for codes in which the minimal weight is not proportional to the dimension. Such codes occur in several contexts, e.g. certain BCH and Reed-Solomon codes [29], Turbo and Turbo-like codes [8, 25, 14, 7]). Additionally, the above discussion makes sense for parity check matrices over arbitrary finite fields, which are also used in coding theory (see [29, 30] for basic information on this topic, and [16] for empirical results on such codes). Finally, sparse parity check matrices are the key ingredient in the construction of small probability spaces and deterministic simulations of $k$-wise independent random variables, which are a key tool in derandomization [1, 2, 4, 10].

Somewhat surprisingly, in spite of their importance, there was a large gap between the known upper and lower bounds for the maximal number of columns of sparse parity check matrices. For a finite field $\mathbb{F}$ let $N_{\mathbb{F}}(n, k, r)$ be the maximal number of vectors in $\mathbb{F}^n$ with at most $r$ non-zero coordinates such that no $k$ of them are linearly dependent (observe that the linear independence condition corresponds to the fact that the kernel of the matrix whose rows are the given vectors is a code with minimal distance at least $k + 1$). When $\mathbb{F} = \mathbb{F}_2$ we use the notation $N_{\mathbb{F}_2}(n, k, r) = N(n, k, r)$. The problem dealt with in this paper, namely that of estimating $N_{\mathbb{F}}(n, k, r)$, differs from the classical Gilbert-Varshamov bounds (see e.g. [29]), since the classical bounds on sizes of codes are geometric packing bounds which depend only on the minimum distance of the code. Here we introduce an additional algebraic restriction on the code (the existence of a sparse parity check matrix) which is motivated by computational issues. Thus, we are dealing with a mixture of a geometric and algebraic problem. In this paper we are primarily interested in the case that $k$ and $r$ are fixed and $n \to \infty$, although the results are valid for arbitrary $k$ and $r$.

Throughout this paper we use the following notation: given two non-negative sequences $\{a_n\}_{n=1}^{\infty}$, $\{b_n\}_{n=1}^{\infty}$, we write $a_n \ll b_n$ if there exists a constant $C > 0$ such that for all $n$, $a_n \leq C \cdot b_n$.

## 1.1 Bounds on $N_{\mathbb{F}}(n, k, r)$

A probabilistic construction [27] (using the first moment method) shows that

$$N(n, k, r) \gg n^{\frac{r}{2} + \frac{r}{2k-2}},$$

and this was generalized to arbitrary finite fields in [26] to $N_{\mathbb{F}}(n, k, r) \gg n^{\frac{r}{2} + \frac{r}{2k-2}}$ for even $k$ and $N_{\mathbb{F}}(n, k, r) \gg n^{\frac{r}{2} + \frac{r}{2k-4}}$ for odd $k$. When $k \geq 4$ is even, and $\gcd(k-1, r) = 1$, the probabilistic lower bound above was improved in [9] to

$$N(n, k, r) \gg n^{\frac{r}{2} + \frac{r}{2k-2}} \cdot (\log n)^{\frac{1}{k-1}}.$$

This lower bound was generalized to arbitrary finite fields in [26] (in which case the constant also depends on the size of the field).

In [27] it was shown that when $k$ is a power of 2, $N(n, k, r) \ll n^{\frac{1}{2}\lceil r + \frac{r}{k-1} \rceil}$, which coincides with the probabilistic lower bound (up to factors independent of $n$) when $k - 1$ divides $r$. This upper bound was generalized to arbitrary finite fields in [26]. Observe that in the important case $k > r$,

i.e. when the number of correctible errors is greater than the weight, this upper bound becomes: for $k$ a power of 2, $N_{\mathbb{F}}(n, k, r) \ll n^{\frac{r}{2} + \frac{1}{2}}$. Thus the gap between the exponent of $n$ in this bound and the probabilistic lower bounds deteriorates as $k$ grows. Here we resolve this problem by proving the following theorem in Section 3:

**Theorem 1.1.** *For every integer $k \geq 8$ and every finite field $\mathbb{F}$*

$$N_{\mathbb{F}}(n, k, r) \ll n^{\frac{r}{2} + \frac{\lceil r/3 \rceil}{2 \lfloor k/8 \rfloor}}$$

*where the implied constant depends only on $k, r$ and $|\mathbb{F}|$.*

The exponent in the displayed inequality behaves roughly like $\frac{r}{2} + \frac{4r}{3k}$ when $k$ is large. This should be compared with the exponent $\frac{r}{2} + \frac{r}{2k-2}$ in the probabilistic lower bound on $N_{\mathbb{F}}(n, k, r)$. In particular it follows from Theorem 1.1 that for any positive integer $r$,

$$\lim_{k \to \infty} \left[ \liminf_{n \to \infty} \frac{\log N_{\mathbb{F}}(n, k, r)}{\log n} \right] = \lim_{k \to \infty} \left[ \limsup_{n \to \infty} \frac{\log N_{\mathbb{F}}(n, k, r)}{\log n} \right] = \frac{r}{2}.$$

It is worthwhile to note here that the proof of Theorem 1.1 when $r$ is even differs markedly from the proof in the case of odd $r$. In fact, it turns out that the case of odd $r$ involves a substantially more subtle argument. The difference between these cases will be explained in section 2. It is likely to be difficult to determine the exact value of the bracketed terms above for all $k$ and $r$ (the answer probably depends on arithmetic properties of $k$ and $r$).

The proof of Theorem 1.1 is based on a novel reduction of the problem to the following Turán type problem: What is the maximum number of edges in an $n$-vertex graph which doesn't contain an even cycle of length $2k$? We then employ recent results on this problem [37, 3, 24, 31] to deduce bounds on $N_{\mathbb{F}}(n, k, r)$. Some of the previous results on $N_{\mathbb{F}}(n, k, r)$ reduced the problem to the study of certain Turán type questions on hypergraphs (see [9, 7]). Since very little is known on hypergraph Turán problems, our main contribution is the method of reducing such questions to a problem on graphs. We believe that this approach is of independent interest. Indeed, as an example we apply our result to a problem in combinatorial number theory, improving a theorem of Erdős, Sárközy and Sós [21] (this application relies heavily on the more difficult part of Theorem 1.1, namely the case of odd $r$). We also apply our result to the even cover problem for hypergraphs.

## 1.2 Even covers in hypergraphs

The methods used to prove Theorem 1.1 in the case $\mathbb{F} = \mathbb{F}_2$ extends the Even Cycle Theorem of Erdős [13] to hypergraphs. If we rephrase the theorem in the terminology of hypergraphs, then we define an *even cover* to be a non-empty collection of sets each point of which is in a even number of sets. For example, a graph contains an even cover of size at most $k$ if and only if it has girth at most $k$. Applying Theorem 1.1 with $\mathbb{F} = \mathbb{F}_2$ to the incidence vectors of edges in a hypergraph, we obtain the following theorem:

**Theorem 1.2.** *Let $\mathcal{S}$ be a hypergraph on $n$ points whose edges have size at most $r$ each, and which does not contain an even cover of size $k$. Then*

$$|\mathcal{S}| \ll n^{\frac{r}{2} + \frac{\lceil r/3 \rceil}{2 \lfloor k/8 \rfloor}}.$$

3

As an example of an application to extremal hypergraph theory, it is notoriously difficult to determine which configurations of triples are guaranteed to appear in every large enough Steiner triple system (for example, many of the questions discussed in [19] remain open). In the present context, it is known that there are infinitely many Steiner triple systems which do not contain an even cover of four triples (one is constructed in [27]). Using Theorem 1.2, we can at least guarantee small even covers in Steiner triple systems. For if $\mathcal{S}$ is a Steiner triple system on $n$ points, then $\mathcal{S}$ has $\frac{1}{3}\binom{n}{2}$ edges [15], which is larger than the bound in Theorem 1.2 for $r = 3$ and $k = 16$, so every large enough Steiner triple system contains an even cover of size at most sixteen. We leave open the problem of determining if smaller even covers exist in every Steiner triple system.

## 1.3   Product representations of squares

Denote by $\mathrm{Rep}_k(n)$ the largest $N$ such that there exists $A \subseteq \{1, \ldots n\}$ with $|A| = N$ such that for every $1 \le \ell \le k$ there are no distinct $a_1, \ldots, a_\ell \in A$ and $x \in \mathbb{Z}$ satisfying $a_1 \cdot a_2 \cdots a_\ell = x^2$. The behavior of the sequence $\mathrm{Rep}_k(n)$ has been studied by several authors [21, 34, 23]. One of the motivations for studying this sequence is that the problem of finding product representations of squares is a key step in certain sub-exponential factoring algorithms (see the surveys [38, 32] for an account of this fascinating field, and [17] for the first rigorous analysis of a sub-exponential randomized factoring algorithm). In these algorithms the goal is to find efficiently a subset of a certain set of integers whose product is a square – the sets that are analyzed are carefully constructed so that such a product representation is guaranteed to exist, but it is of interest to ask how large such a set should be in order to ensure the existence of the required product representation. Moreover, once the cardinality of a set is above this threshold, one would like to efficiently find such a product representation. The best known results on this question follow from the work of Erdős [18] and Erdős, Sárközy and Sós [21], who showed that for every $k \ge 6$,

$$\left(\frac{n}{\log n}\right)^{\frac{1}{2}+\frac{1}{6k}} \ll \mathrm{Rep}_k(n) - \pi(n) \ll \left(\frac{n}{(\log n)^2}\right)^{\frac{2}{3}},$$

where the implied constants are absolute and $\pi(n)$ is the number of primes less than or equal to $n$. Here we show that for large $k$, the order of magnitude of $\mathrm{Rep}_k(n) - \pi(n)$ is roughly $N^{\frac{1}{2}}$, namely in Section 4 we prove the following theorem:

**Theorem 1.3.** *For every $k \ge 1$ and every integer $n$,*

$$\left(\frac{n}{\log n}\right)^{\frac{1}{2}+\frac{1}{48k}} \ll \mathrm{Rep}_{8k}(n) - \pi(n) \ll kn^{\frac{1}{2}+\frac{1}{2k}}(\log n)^2.$$

## 1.4   An algorithm for linear dependences

Our proof of Theorem 1.1 yields an algorithm which, given a set $X \subseteq \mathbb{F}^n$ of vectors of weight at most $r$ with

$$|X| \gg n^{\frac{r}{2}+\frac{\lceil r/3 \rceil}{2\lfloor k/8 \rfloor}},$$

finds $k$ linearly dependent vectors in $X$ in time quadratic in $|X|$ – that is in time $O(|X|^2) = O(n^{2r})$. Observe that an exhaustive check of all possible $\binom{|X|}{k}$ linear dependencies requires time $\Omega(n^{\frac{kr}{2}})$. In

4

the case that $r$ is even, we will give a short proof of Theorem 1.1 which itself yields an algorithm whose running time is linear in $|X|$. These algorithms will follow directly from our proof of Theorem 1.1 and the Alon-Yuster-Zwick algorithm [5, 6] for finding cycles in graphs, or alternatively the proof of the main theorem in [37], which gives better constants. Additionally, our proof yields an algorithm such that given a set $A \subseteq \{1, \ldots, n\}$ with $|A| \gg kn^{\frac{1}{2}+\frac{1}{2k}}(\log n)^2$, finds in quadratic time, distinct $a_1, \ldots, a_j \in A$ with $1 \leq j \leq 8k$, such that $a_1 a_2 \cdots a_j$ is a square.

# 2 Forbidden cycles and sparse parity check matrices

Throughout this section, $\mathbb{F}$ is a finite field and $\mathbb{F}^* = \mathbb{F} \setminus \{0\}$ is the set of non-zero elements of $\mathbb{F}$. A set $X \subseteq \mathbb{F}^n$ is $k$-wise independent if no $k$ vectors in $X$ are linearly dependent. A vector $v \in \mathbb{F}^n$ is said to have weight $r$ if it has exactly $r$ non-zero coordinates. Then $N_{\mathbb{F}}(n, k, r)$ is the maximum size of a set of $k$-wise independent vectors of weight at most $r$ in $\mathbb{F}^n$. The following result on forbidden cycles will be used in the proof of Theorem 1.1:

**Theorem 2.1.** *Let $k$ be a positive integer and let $G = (V, E)$ be an $N$-vertex graph. If $G$ contains no cycle of length exactly $2k$ then*

$$|E| < 2kN^{1+\frac{1}{k}}. \tag{1}$$

*If $G$ is an $M$ by $N$ bipartite graph containing no cycle of length $2k$, then*

$$|E| < 2k \cdot [M^{\frac{1}{2}} N^{\frac{1}{2}+\frac{1}{k}} + M + N]. \tag{2}$$

*If $G$ has girth at least $2k + 1$, then the factor $2k$ may be removed in each of the upper bounds.*

Theorem 2.1 is proved in [31], the last statement is proved in [3], the bipartite case of which is proved in [24]. The first assertion (1), with the same dependence on $n$ but a worse constant, is Erdős' Even Cycle Theorem – see [13].

## 2.1 Vectors of even weight

We will reduce the problem of estimating $N_{\mathbb{F}}(n, k, r)$ to the bounds in Theorem 2.1. The reduction we give is simple in the case that $r$ is even, but substantially more involved when $r$ is odd. The first theorem we prove is as follows:

**Theorem 2.2.** *Let $n, k, r$ be positive integers, and define*

$$M = \sum_{\ell=0}^{\lfloor r/2 \rfloor} (|\mathbb{F}| - 1)^\ell \binom{n}{\ell} \qquad N = \sum_{\ell=0}^{\lceil r/2 \rceil} (|\mathbb{F}| - 1)^\ell \binom{n}{\ell}.$$

*Let $X \subseteq \mathbb{F}^n$ be a set of vectors of weight at most $r$ such that*

$$|X| \;>\; 2k \cdot [M^{\frac{1}{2}} N^{\frac{1}{2}+\frac{1}{k}} + M + N].$$

*Then there are disjoint sets $A, B \subseteq X$, each of size $k$, such that $\sum_{a \in A} a \;=\; \sum_{b \in B} b$. In particular, there is a set of $2k$ linearly dependent vectors in $X$.*

5

The bound in Theorem 2.2 implies in particular that $N_{\mathbb{F}}(n, 2k, r) \ll n^{\frac{r}{2} + \frac{1}{2}}$ for all $k$. This generalizes the same bound which was proved for $k$ a power of two in [27] to all values of $k$. The proof of Theorem 2.2 is short, and we present it here:

**Proof of Theorem 2.2.** For each vector $v \in X$ of weight $\omega \leq r$, fix a pair $e(v) = \{x, y\}$ of vectors in $\mathbb{F}^n$ of weights at most $\lfloor \omega/2 \rfloor$ and $\lceil \omega/2 \rceil$, respectively, satisfying $x + y = v$. If $r$ is even, let $G$ be the graph whose vertex set consists of all vectors in $\mathbb{F}^n$ of weight at most $r/2$ and whose edge set is $E = \{e(v) : v \in X\}$. Then $|E| = |X|$ and $G$ has $M$ vertices. By the first assertion (1) in Theorem 2.1, $G$ contains a cycle of length $2k$. By the definition of $G$, there exist distinct vectors $v_1, v_2, \ldots, v_{2k} \in X$ such that the edge set of this cycle consist of the pairs $e(v_i) = \{x_i, x_{i+1}\}$ for $i = 1, 2, \ldots, 2k$, where $x_{2k+1} = x_1$. Then

$$\sum_{\ell=1}^{2k} (-1)^{\ell+1} v_i = (x_1 + x_2) - (x_2 + x_3) + \cdots - (x_{2k} + x_1) = 0,$$

and the disjoint sets $A = \{v_1, v_3, \ldots, v_{2k-1}\}$ and $B = \{v_2, v_4, \ldots, v_{2k}\}$ have the same sum. If $r$ is odd, then $G$ is a bipartite graph whose parts have sizes $M$ and $N$. By the second assertion (2) of Theorem 2.1, $C_{2k} \subseteq G$, and we conclude by the same argument as that presented above. $\qquad \square$

Theorem 2.2 implies Theorem 1.1 in the case of even $r$: in fact in this case we obtain the stronger bound $N_{\mathbb{F}}(n, k, r) \ll n^{\frac{r}{2} + \frac{r}{k}}$ for $k$ even. We also remark that the proof above gives a linear time algorithm – that is time $O(|X|)$ – for finding $2k$ linearly dependent vectors in a set $X \subseteq \mathbb{F}^n$ satisfying the requirements of the theorem – this follows from the linear time algorithm in [5, 6] for finding a cycle of length $2k$ in a graph with appropriately many edges.

## 2.2 Vectors of odd weight

In the case of odd $r$, Theorem 2.2 is insufficient to deduce Theorem 1.1, since the rounding up of $\frac{r}{2}$ only yields an upper bound of order $n^{\frac{r}{2} + \frac{1}{2}}$. This difference between even and odd values of $r$ leads to a more involved argument in the case of odd $r$, which nevertheless yields the bounds of Theorem 1.1 in this case as well. We conclude the proof of Theorem 1.1 by giving the following refinement of Theorem 2.2 when $k > r$ and $r$ is odd. For convenience we define $[\frac{r}{3}] = r - \lceil \frac{r}{3} \rceil - \lfloor \frac{r}{3} \rfloor$ and $q = |\mathbb{F}^*|$:

**Theorem 2.3.** *Let $n, k$ and $r$ be positive integers, and define*

$$L = \sum_{\ell=0}^{\lfloor r/3 \rfloor} q^\ell \binom{n}{\ell} \qquad M = \sum_{\ell=0}^{[r/3]} q^\ell \binom{n}{\ell} \qquad N = \sum_{\ell=0}^{\lceil r/3 \rceil} q^\ell \binom{n}{\ell}.$$

*Let $X \subseteq \mathbb{F}^n$ be an $8k$-wise independent set of vectors. Then*

$$|X| \quad < \quad 12kr \cdot [(LN)^{\frac{1}{2}} M^{\frac{1}{2} + \frac{1}{2k}} + N^{1 + \frac{1}{2k}}].$$

In particular, it follows that $N_{\mathbb{F}}(n, 8k, r) \ll n^{\frac{r}{2} + \frac{[r/3]}{2k}}$, which gives Theorem 1.1 when we replace $k$ with $\lfloor k/8 \rfloor$. The remainder of the paper is devoted to the proof of Theorem 2.3.

# 3 Proof of Theorem 1.1

We already proved Theorem 1.1 in the case that $r$ is even. In this section we prove Theorem 2.3 which implies Theorem 1.1 when $r$ is odd. Since the proof is quite involved, we begin with an outline of the proof. For simplicity we omit all the multiplicative constants. Starting with a set $X \subseteq \mathbb{F}^n$ where $|X| \gg (LN)^{\frac{1}{2}} M^{\frac{1}{2} + \frac{1}{2k}} + N^{1 + \frac{1}{2k}}$, we wish to find a linearly dependent set of $8k$ vectors in $X$. Suppose that there is no such subset of $X$. Via an averaging argument (see §3.6), we will show that it is sufficient to find such a linear dependence in a set $Z \subseteq \mathbb{F}^L \times \mathbb{F}^M \times \mathbb{F}^N$, where $|Z| \gg |X|$ and where the projection of $Z$ onto any of the subspaces $\mathbb{F}^L, \mathbb{F}^M$ and $\mathbb{F}^N$ consists of vectors of weight one. We call such sets *balanced*.

To find linear dependences in balanced sets $Z$ where $|Z| \gg |X|$, we first prove in §3.1 that there is a set $Y \subseteq Z$ where $|Y| > |Z| - L^{1 + \frac{1}{2k}} + M^{1 + \frac{1}{2k}} + N^{1 + \frac{1}{2k}}$ and the projection of $Y$ onto any one of the subspaces $\mathbb{F}^L \times \mathbb{F}^M$ and $\mathbb{F}^L \times \mathbb{F}^N$ and $\mathbb{F}^M \times \mathbb{F}^N$ uniquely determines $Y$. In this case we say that $Y$ is *determined by projection*. This allows us in §3.2 to count special four-element subsets of $Y$ called *partially dependent quadruples*: these are sets of four vectors in $Y$ whose projection onto $\mathbb{F}^L \times \mathbb{F}^N$ is a linearly dependent set of four vectors. The key point in the proof is to prove that there are enough of these quadruples in $Y$ to ensure that the projection of $2k$ of these quadruples onto $\mathbb{F}^M$ results in a linearly dependent set. Then this gives $8k$ linearly dependent vectors in $Y$, and the required contradiction.

To obtain a feeling for where the bound in Theorem 2.2 comes from, we give some of the bounds involved in the proof. Let $\mathcal{Q}$ denote the set of partially dependent quadruples in $Y$. We will show in §3.2 that

$$|\mathcal{Q}| > \frac{|Y|^4}{4L^2 N^2}.$$

Now $|Y| \gg (LN)^{\frac{1}{2}} M^{\frac{1}{2} + \frac{1}{2k}}$, so this inequality gives $|\mathcal{Q}| \gg M^{2 + 2/k}$. Since the projection of $Y$ onto $\mathbb{F}^M$ consists of vectors of weight one, the projection of each quadruple in $\mathcal{Q}$ onto $\mathbb{F}^M$ consists of four vectors of weight one. If we treat these projections as vectors of weight four in $\mathbb{F}^M$, then there are $|\mathcal{Q}| \gg M^{2 + 2/k}$ of these vectors. By Theorem 2.2 with $r = 4$, this shows that there are $2k$ quadruples in $\mathcal{Q}$ whose projections onto $\mathbb{F}^M$ form a linearly dependent set of $2k$ vectors. In principle, we then consider all the vectors in the corresponding $2k$ quadruples in $\mathcal{Q}$ to obtain a linearly dependent set of at most $8k$ vectors in $Y$, and this contradiction completes the proof of Theorem 2.2.

However, there is a subtlety, which is that the vectors in the $2k$ quadruples in $\mathcal{Q}$ might form a trivial dependence in $Y$. One reason might be that each vector appears with coefficient zero modulo $p$ in the linear dependence, where $p$ is the characteristic of $\mathbb{F}$, in which case the linear dependence is trivial. In the proof we consider only special types of linear dependences of projections of quadruples onto $\mathbb{F}^M$, which guarantee that when we lift back to the quadruples themselves, the linear dependence is not trivial. This approach is covered in §3.3 and §3.4.

## 3.1 Sets determined by projection

We represent the elements of $\mathbb{F}^\alpha \times \mathbb{F}^\beta \times \mathbb{F}^\gamma$ as triples $(v_\alpha, v_\beta, v_\gamma)$. For convenience we write $\lambda = |\mathbb{F}^*| \alpha$, $\mu = |\mathbb{F}^*| \beta$ and $\nu = |\mathbb{F}^*| \gamma$. A set $Z \subseteq \mathbb{F}^\alpha \times \mathbb{F}^\beta \times \mathbb{F}^\gamma$ is *balanced* if the components of each $(v_\alpha, v_\beta, v_\gamma) \in Z$ have weight one. The *projection* of a set $Z$ onto $\mathbb{F}^\alpha$ is

$$Z_\alpha = \{v \in \mathbb{F}^\alpha : \exists \, (v_\beta, v_\gamma) \in \mathbb{F}^\beta \times \mathbb{F}^\gamma, \text{ such that } (v, v_\beta, v_\gamma) \in Z\}.$$

Similarly, $Z_{\alpha\beta}$ denotes the projection of $Z$ onto $\mathbb{F}^\alpha \times \mathbb{F}^\beta$. For $v \in \mathbb{F}^\alpha \times \mathbb{F}^\beta$ let $Z^v = \{z \in Z : z_{\alpha\beta} = v\}$. A set $Y \subseteq \mathbb{F}^\alpha \times \mathbb{F}^\beta \times \mathbb{F}^\gamma$ is *determined by projection* if any one of the sets $Y_{\alpha\beta}$, $Y_{\beta\gamma}$, $Y_{\alpha\gamma}$ uniquely determines $Y$. The following lemma says that $8k$-wise independent sets contain large subsets which are determined by projection:

**Lemma 3.1.** *Let $Z \subseteq \mathbb{F}^\alpha \times \mathbb{F}^\beta \times \mathbb{F}^\gamma$ be a balanced $8k$-wise independent set. Then there exists a set $Y \subseteq Z$ such that $Y$ is determined by projection and*

$$|Y| \; > \; |Z| - \lambda^{1+\frac{1}{2k}} - \mu^{1+\frac{1}{2k}} - \nu^{1+\frac{1}{2k}}. \tag{3}$$

*Proof.* For $v \in Z_{\alpha\beta}$, let $T(v)$ be a spanning tree of the complete graph on $Z^v$. So $|E(T(v))| = |Z^v| - 1$. We claim that the trees $\{T(v)\}_{v \in Z_{\alpha\beta}}$ can be chosen so that the *multigraph* $G_{\alpha\beta}$ consisting of all edges in all the trees $\{T(v)\}_{v \in Z_{\alpha\beta}}$ has girth greater than $4k$. This is done by choosing the trees $\{T(v)\}_{v \in Z_{\alpha\beta}}$ so that the girth of $G_{\alpha\beta}$ is a minimum. This choice implies that if $C$ is a shortest cycle in $G$, then $|C \cap T(v)| \leq 1$ for all $v \in Z_{\alpha\beta}$. We aim to show that $|C| > 4k$. Suppose the edges of $C$ are $\{\{w_1, w_2\}, \{w_2, w_3\}, \ldots, \{w_\ell, w_1\}\}$. Then there are distinct $v_j = (v_\alpha^j, v_\beta^j) \in Z_{\alpha\beta}$ such that $\{w_j, w_{j+1}\} \in T(v_j)$ for all $j \leq \ell$. Let $x_j = (v_j, w_j) \in Z$ and $y_j = (v_j, w_{j+1}) \in Z$ for $j \leq \ell$. Then

$$\sum_{j=1}^{\ell} x_j - \sum_{j=1}^{\ell} y_j = 0.$$

Now $x_i$ and $y_j$ are distinct for all $i, j \leq \ell$. This means that $\{x_1, x_2, \ldots, x_\ell, y_1, y_2, \ldots, y_\ell\} \subseteq Z$ is linearly dependent. Since $Z$ is $8k$-wise independent, it follows that $\ell > 4k$, as required, so $G$ has girth greater than $4k$. The number of edges in $G_{\alpha\beta}$ is

$$|E(G)| \;=\; \sum_{v \in Z_{\alpha\beta}} (|Z^v| - 1). \tag{4}$$

On the other hand, since $G_{\alpha\beta}$ has at most $c$ vertices (since $Z$ is balanced), we conclude from Theorem 2.1 that $|E(G_{\alpha\beta})| < \nu^{1+\frac{1}{2k}}$. We may define $G_{\beta\gamma}$ and $G_{\alpha\gamma}$ similarly, and by symmetry $|E(G_{\beta\gamma})| < \lambda^{1+\frac{1}{2k}}$. $|E(G_{\alpha\gamma})| < \mu^{1+\frac{1}{2k}}$. Using (4), these inequalities translate to

$$\sum_{v \in Z_{\alpha\beta}} (|Z^v| - 1) \;<\; \nu^{1+\frac{1}{2k}}.$$

$$\sum_{v \in Z_{\alpha\gamma}} (|Z^v| - 1) \;<\; \mu^{1+\frac{1}{2k}}.$$

$$\sum_{v \in Z_{\beta\gamma}} (|Z^v| - 1) \;<\; \lambda^{1+\frac{1}{2k}}.$$

Finally, the number of vectors in $Z$ which are *not* determined by projection is exactly the sum of the three terms on the left in the above inequalities. So the number of vectors in $Z$ which are determined by projection is greater than

$$|Z| - \lambda^{1+\frac{1}{2k}} - \mu^{1+\frac{1}{2k}} - \nu^{1+\frac{1}{2k}}.$$

Let $Y$ be the set of all these vectors; then $Y$ satisfies the requirements of the lemma. $\qquad\square$

## 3.2 Partially dependent quadruples

For the remainder of the proof of Theorem 2.3, we restrict our attention to a set $Y \subseteq Z$ which is $8k$-wise independent and determined by projection, and satisfies (3). A *partially dependent quadruple* in $Y$ is a quadruple $\{w, x, y, z\} \subseteq Y$ such that $(w_\alpha, x_\gamma, y_\alpha, z_\gamma) = (x_\alpha, y_\gamma, z_\alpha, w_\gamma)$. When it is convenient, we represent this quadruple as an ordered 4-tuple $(w, x, y, z)$ with the understanding that the non-zero co-ordinate of $w_\alpha$ precedes the non-zero co-ordinate of $y_\alpha$ and the non-zero co-ordinate of $w_\gamma$ precedes the non-zero co-ordinate of $y_\gamma$. A partially dependent quadruple is illustrated in Figure 1.
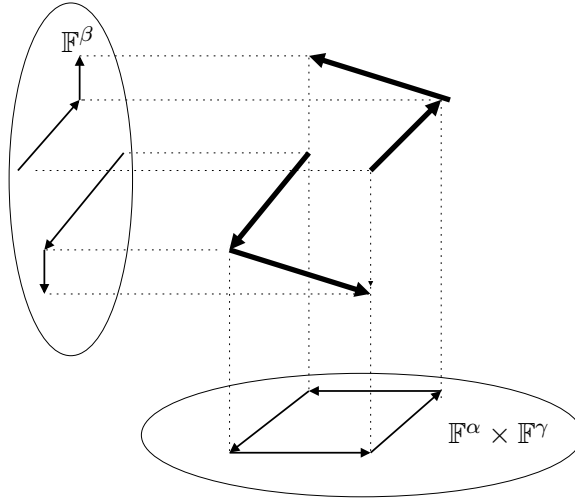


Figure 1 : A partially dependent quadruple

The key point is that the projection of a partially dependent quadruple onto $\mathbb{F}^\alpha \times \mathbb{F}^\gamma$ is a linearly dependent set. Therefore our aim is to try to find a linearly dependent set of $2k$ projections of quadruples in $Y$ onto $\mathbb{F}^\beta$. Since each quadruple consists of four vectors in $Y$, altogether this gives $8k$ linearly dependent vectors in $Y$. Let $\mathcal{Q}$ denote the set of partially dependent quadruples in $Y$.

**Lemma 3.2.** *Suppose that $|Y| > \nu + 2\lambda\nu^{\frac{1}{2}}$. Then*

$$|\mathcal{Q}| \;>\; \frac{|Y|^4}{4\lambda^2\nu^2}. \tag{5}$$

*Proof.* For $v \in \mathbb{F}^\alpha$, recall that $Y^v = \{(v_\alpha, v_\beta, v_\gamma) \in Y : v_\alpha = v\}$. We define $Y^v$ similarly when $v \in \mathbb{F}^\beta$ or $v \in \mathbb{F}^\gamma$. To prove (5), we use the following identity, which follows from the fact that $Y$ is determined by projection:

$$\sum_{\{v,x\} \subseteq \mathbb{F}^\alpha} |Y^v \cap Y^x| \;=\; \sum_{w \in \mathbb{F}^\gamma} \binom{|Y^w|}{2}. \quad (*)$$

9

It follows that

$$
\begin{aligned}
|\mathcal{Q}| &= \sum_{\{v,x\}\subseteq \mathbb{F}^\alpha} \binom{|Y^v \cap Y^x|}{2} \\
&\geq \binom{\lambda}{2}\left(\frac{\frac{1}{\binom{\lambda}{2}}\sum_{\{v,x\}\subseteq \mathbb{F}^\alpha}|Y^v \cap Y^x|}{2}\right) \\
&\overset{(*)}{=} \binom{\lambda}{2}\left(\frac{\frac{1}{\binom{\lambda}{2}}\sum_{w\in \mathbb{F}^\gamma}\binom{|Y^w|}{2}}{2}\right) \\
&\geq \binom{\lambda}{2}\left(\frac{\frac{\nu}{\binom{\lambda}{2}}\binom{\frac{1}{\nu}\sum_{w\in \mathbb{F}^\gamma}|Y^w|}{2}}{2}\right) \\
&= \binom{\lambda}{2}\left(\frac{\frac{\nu}{\binom{\lambda}{2}}\binom{\nu^{-1}|Y|}{2}}{2}\right) \\
&> \frac{|Y|^4}{4\lambda^2\nu^2}.
\end{aligned}
$$

This is exactly (5). In each of the inequalities we used the convexity of the function $a \mapsto \binom{a}{2}$. In the last inequality we used the lower bound on $|Y|$ assumed in the lemma. $\qquad\square$

## 3.3 Constructing linear dependences

For each partially dependent quadruple $Q = (w, x, y, z)$ in $\mathcal{Q}$, define $\pi(Q) = \{(w_\beta, y_\beta), (x_\beta, z_\beta)\}$. It is convenient to define the multigraph $G$ consisting of all pairs $\pi(Q)$ where $Q \in \mathcal{Q}$ (the vertex set of $G$ is a subset of $\mathbb{F}^\beta \times \mathbb{F}^\beta$). A *chain of length $k$* is a sequence $Q = (Q_1, Q_2, \ldots, Q_k)$ where $Q_i \in \mathcal{Q}$ such that $(\pi(Q_1), \pi(Q_2), \ldots, \pi(Q_k))$ is a non-returning walk (a walk in which consecutive edges are distinct) of length $k$ in $G$ Two chains are *concyclic* if the walks corresponding to them in $G$ have the same endpoints. If $Q = (Q_1, Q_2, \ldots, Q_k)$ and $R = (R_1, R_2, \ldots, R_k)$ are concyclic chains, where $Q_i = (s_i, t_i, u_i, v_i)$ and $R_i = (w_i, x_i, y_i, z_i)$, then $\{(s_1)_\beta, (u_1)_\beta\} = \{(w_1)_\beta, (y_1)_\beta\}$ and $\{(t_k)_\beta, (v_k)_\beta\} = \{(x_k)_\beta, (z_k)_\beta\}$. It follows that

$$
\sum_{i=1}^{k}(s_i + u_i - t_i - v_i) - \sum_{i=1}^{k}(w_i + y_i - x_i - z_i) = 0. \tag{6}
$$

We want to find concyclic chains $Q$ and $R$ such that the equation above is a non-trivial linear dependence of at most $8k$ vectors in $Y$ (it is possible to construct many examples where this linear dependence is trivial). To find a non-trivial dependence, it is sufficient to show that some vector appears in the equation (6) with a coefficient which is not zero modulo $p$, where $p$ is the characteristic of $\mathbb{F}$. This will hold for certain special chains which we call *nondegenerate* chains.

Let $Q = (Q_1, Q_2, \ldots, Q_k)$ be a chain of length $k$ where $Q_i = (s_i, t_i, u_i, v_i)$ for $i \in \{1, 2, \ldots, k\}$. Then the *reduction of $Q$* is the set of vectors defined by

$$
\widehat{Q} := Q_k \,\triangle_p\, Q_{k-1} \,\triangle_p\, \cdots \,\triangle_p\, Q_1. \tag{7}
$$

The expression above is read from right to left, and the symmetric difference operator $\triangle_p$ is defined as follows: we delete any vector once it appears with coefficient zero mod $p$ in the sum:

$$\sum_{i=1}^{k}(s_i + u_i - t_i - v_i)$$

We say that $Q$ is *nondegenerate* if $|Q_i \cap (Q_1 \cup Q_2 \cup \cdots \cup Q_{i-1})| \leq 1$, for $2 \leq i \leq k$, and $Q$ is *degenerate* otherwise.

**Lemma 3.3.** *Suppose $Q, R$ are concyclic nondegenerate chains of length $k$ in $\mathcal{Q}$. Then $\widehat{Q} = \widehat{R}$.*

*Proof.* Suppose $Q = (Q_1, Q_2, \ldots, Q_k)$ and $R = (R_1, R_2, \ldots, R_k)$ where $Q_i = (s_i, t_i, u_i, v_i)$ and $R_i = (w_i, x_i, y_i, z_i)$. Since $Q$ and $R$ are concyclic, equation (6) holds. If we restrict the sum on the left of (6) to vectors in $\widehat{Q}$, the equality still holds, since $p$ is the characteristic of $\mathbb{F}$. Similarly, we may delete all terms on the right of (6) which are not in $\widehat{R}$. Since $Q$ and $R$ are nondegenerate, $\widehat{Q}$ and $\widehat{R}$ are each non-empty. If $\widehat{Q} \triangle \widehat{R} \neq \emptyset$, then $\widehat{Q} \triangle \widehat{R}$ is a non-empty linearly dependent set of at most $8k$ vectors in $Y$, a contradiction. So $\widehat{Q} = \widehat{R}$, as required. $\qquad\square$

## 3.4 Counting nondegenerate chains

**Lemma 3.4.** *If $Q = (Q_1, Q_2, \ldots, Q_\ell)$ is a nondegenerate chain in $\mathcal{Q}$, then at most $8\ell^2$ degenerate chains of length $\ell + 1$ contain $Q$.*

*Proof.* Let $(w_\beta, y_\beta)$ be the endvertex of the walk $(\pi(Q_1), \pi(Q_2), \ldots, \pi(Q_\ell))$ in $G$ and let $f(Q) = \bigcup_{i=1}^{\ell} Q_i$. Then the number of degenerate chains of length $\ell+1$ containing $Q$ is equal to the number of partially dependent quadruples $R \in \mathcal{Q}$ such that $|R \cap f(Q)| \geq 2$. Let us assume

$$R = (w, x, y, z) = ((w_\alpha, w_\beta, w_\gamma), (y_\alpha, x_\beta, w_\gamma), (y_\alpha, y_\beta, y_\gamma), (w_\alpha, z_\beta, y_\gamma)). \qquad (8)$$

We claim that exactly two elements of $S = \{w, x, y, z\} \cap f(Q)$ uniquely determine $R$. Once that is proved, it follows that the number of $R$ such that $(Q_1, Q_2, \ldots, Q_\ell, R)$ is degenerate is at most

$$\binom{|f(Q)|}{2} \;\leq\; \binom{|Q_1| + |Q_2| + \cdots + |Q_\ell|}{2} \;=\; \binom{4\ell}{2} \;<\; 8\ell^2$$

as required. We now prove the claim. Since $Y$ is determined by projection, $R$ is uniquely determined upon specifying the projection of $R$ onto $\mathbb{F}^\alpha$ and onto $\mathbb{F}^\gamma$. So the claim is proved if $S_\alpha = \{w_\alpha, y_\alpha\}$ and $S_\gamma = \{w_\gamma, y_\gamma\}$ – since in that case two co-ordinates of each vector in the expression (8) defining $R$ are specified. If this is not the case, then one checks that $S$ is one of the pairs $\{w, x\}$, $\{x, y\}$, $\{y, z\}$ or $\{z, w\}$. These cases are all dealt with in the same way, so we check only the case $S = \{w, x\}$. Since $y_\gamma$ is uniquely determined by $y_\alpha$ and $y_\beta$, $y$ is uniquely determined by $S$. Hence $\{w, x, y\}$ are specified, and therefore $\{w_\alpha, y_\alpha\}$ and $\{w_\gamma, y_\gamma\}$ are uniquely determined. This uniquely determines $R$, and proves the claim. $\qquad\square$

**Lemma 3.5.** *Let $\mathcal{P}_k$ denote the set of nondegenerate chains of length $k$ in $\mathcal{Q}$. Then*

$$|\mathcal{P}_k| \;>\; 4^{-k}m(d-32k^2)^k, \tag{9}$$

*where $d > 64k^2$ is the average degree in $G$ and $m$ is the number of vertices in $G$.*

*Proof.* We claim that for all $\ell \le k$ and $d > 64\ell^2$,

$$|\mathcal{P}_\ell| \;>\; 4^{-\ell}m(d-32\ell^2)^\ell. \tag{10}$$

The proof is by induction on $m+\ell$. If $G$ contains a vertex of degree less than $d/4$, then we remove such a vertex to obtain a graph of average degree greater than $\tilde{d} = d + d/(2m-2)$. By induction, the number of walks in $\mathcal{P}_\ell$ for this new graph is at least

$$4^{-\ell}(m-1)(\tilde{d}-32\ell^2)^\ell \;>\; 4^{-\ell}m(d-32\ell^2)^\ell.$$

In particular, the number of non-returning walks of $G$ which are in $\mathcal{P}_\ell$ is at least $4^{-\ell}m(d-32\ell^4)^\ell$, as required. Suppose every vertex of $G$ has degree at least $d/4$, and

$$|\mathcal{P}_{\ell-1}| \;>\; 4^{-\ell+1}m(d-32(\ell-1)^2)^{\ell-1}.$$

Since every walk in $\mathcal{P}_{\ell-1}$ has at most $8(\ell-1)^2$ extensions to a degenerate walk of length $\ell$, by Lemma 3.4, there are at least $d/4 - 8(\ell-1)^2 - \ell > d/4 - 8\ell^2$ extensions of each walk in $\mathcal{P}_{\ell-1}$ to a walk in $\mathcal{P}_{\ell+1}$. This proves (10). Now the lemma follows from Lemma 3.4 with $\ell = k$. $\qquad\square$

**Lemma 3.6.** *Let $m$ be the number of vertices in $G$. Then*

$$|\mathcal{P}_k| \;<\; 4^k k^{4k}\binom{m}{2}. \tag{11}$$

*Proof.* By Lemma 3.3, if $Q = (Q_1, Q_2, \ldots, Q_k)$ is a chain in $\mathcal{P}_k$, then the number of chains $R = (R_1, R_2, \ldots, R_k)$ in $\mathcal{P}_k$ such that $Q$ and $R$ are concyclic is at most the number of choices of $R$ such that $\widehat{R} = \widehat{Q}$. Let $f(Q) = \bigcup_{i=1}^k Q_i$ and $f(R) = \bigcup_{i=1}^k R_i$. The important point is that since $Q$ and $R$ are nondegenerate, their projections onto $\mathbb{F}^\alpha$ and $\mathbb{F}^\gamma$ are invariant under reduction:

$$\begin{aligned}
\widehat{Q}_\alpha = f(Q)_\alpha \qquad & \widehat{Q}_\gamma = f(Q)_\gamma \\
\widehat{R}_\alpha = f(R)_\alpha \qquad & \widehat{R}_\gamma = f(R)_\gamma.
\end{aligned}$$

Therefore the number choices of $R \in \mathcal{P}_k$ concyclic with $Q \in \mathcal{P}_k$ is at most the number of choices of $R$ such that $f(R)_\alpha = f(Q)_\alpha$ and $f(R)_\gamma = f(Q)_\gamma$. Since $Y$ is determined by projections, any quadruple $R_i \in \mathcal{Q}$ is specified by its projection onto $\mathbb{F}^\alpha \times \mathbb{F}^\gamma$. The number of ways of choosing quadrilaterals $R_1, R_2, \ldots, R_k$ so that $R = (R_1, R_2, \ldots, R_k)$ is at most

$$\binom{|f(Q)_\alpha|}{2}^k \binom{|f(Q)_\gamma|}{2}^k \;\le\; \binom{2k}{2}^k \binom{2k}{2}^k \;<\; 4^k k^{4k}$$

since $|f(Q)_\alpha| \le 2k$ and $|f(Q)_\beta| \le 2k$. This gives the required upper bound on $|\mathcal{P}_k|$. $\qquad\square$

## 3.5 Linear Dependences

In this section we prove Theorem 2.3, using the lemmas we have developed in the last few sections. The following theorem combines all of these lemmas, and will also be used to prove Theorem 1.3.

**Theorem 3.7.** *Let $Z \subseteq \mathbb{F}^\alpha \times \mathbb{F}^\beta \times \mathbb{F}^\gamma$ be a balanced $8k$-wise independent set of vectors. Then*

$$|Z| \; < \; 4k(\lambda\mu\nu)^{\frac{1}{2}}\mu^{\frac{1}{2k}} + \lambda^{1+\frac{1}{2k}} + \mu^{1+\frac{1}{2k}} + \nu^{1+\frac{1}{2k}} + 2\lambda\nu^{\frac{1}{2}}. \tag{12}$$

*Proof.* Since $Z$ is $8k$-wise independent, we may apply Lemma 3.1: there exists a set $Y \subseteq Z$ such that $Y$ is determined by projection and

$$|Y| \; > \; \left(16k^2\mu^{1+\frac{1}{k}}\lambda\nu\right)^{\frac{1}{2}} + \nu + 2\lambda\nu^{\frac{1}{2}}. \tag{13}$$

Let $\mathcal{P}_k$ denote the set of nondegenerate chains of length $k$ in $\mathcal{Q}$, let $d$ and $m$ be the average degree and number of vertices in $G$, respectively. Combining (11) and (9) gives $d - 32k^2 < 4k^4m^{\frac{1}{k}}$ and therefore $d < 64k^4m^{\frac{1}{k}}$. It follows that since $|\mathcal{Q}| = |E(G)| = \frac{1}{2}dm$,

$$|\mathcal{Q}| \; < \; 64k^4m^{1+\frac{1}{k}}. \tag{14}$$

For a contradiction, suppose that $|Z|$ is at least the expression claimed in (12). Now from (13),

$$|\mathcal{Q}| \; > \; \frac{|Y|^4}{4\lambda^2\nu^2}. \tag{15}$$

Using (14) and $m = \mu^2$ this gives

$$|Y|^4 \; < \; 4\lambda^2\nu^2 \cdot 64k^4m^{1+\frac{1}{k}} \; < \; 256k^4\lambda^2\nu^2\mu^{2+\frac{2}{k}} \; = \; \left(16k^2\mu^{1+\frac{1}{k}}\lambda\nu\right)^2$$

This contradicts (13), and proves the theorem. $\qquad\square$

**Remark.** Theorem 3.7 can be used to derive a more precise version of Theorem 1.2: if $\mathcal{S}$ is an $r$-partite hypergraph with parts of sizes $N_1, N_2, \ldots, N_r$, then the above theorem can be used to prove that if $|\mathcal{S}| \gg (N_1N_2\ldots N_r)^{\frac{1}{2}+\frac{1}{2k}} + (N_1 + N_2 + \cdots + N_r)^{\lceil\frac{r}{3}\rceil(1+\frac{1}{2k})}$ then $\mathcal{S}$ contains an even cover of size at most $8k$. This result may be viewed as an extension of Theorem 2.1 from bipartite graphs to $r$-partite hypergraphs.

*Proof of Theorem 1.1.* Let $X \subseteq \mathbb{F}^n$ be a set of vectors of weight at most $r$. Let $\chi : \{1, 2, \ldots, n\} \to \{1, 2, 3\}$ be a random three-coloring of the co-ordinates, where distinct co-ordinates are colored independently and each color is equiprobable. For a vector $x \in X$ of weight $\omega(x) = \omega$, the probability that $x$ has exactly $\lfloor\frac{\omega}{3}\rfloor$ non-zero co-ordinates of color 1, exactly $\lceil\frac{\omega}{3}\rceil$ non-zero co-ordinates of color 2, and exactly $\lceil\frac{\omega}{3}\rceil$ non-zero co-ordinates of color three is exactly

$$\frac{1}{3^\omega}\binom{\omega}{\lfloor\frac{\omega}{3}\rfloor}\binom{\omega - \lfloor\frac{\omega}{3}\rfloor}{\lceil\frac{\omega}{3}\rceil} \; > \; \frac{1}{3^\omega},$$

where we used the numerical Lemma 3.8 (see below) to obtain the inequality. In this case we say that $x$ is equipartitioned by $\chi$. Therefore the expected number of vectors $x \in X$ which are equipartitioned is greater than

$$\sum_{x \in X} \frac{1}{3\omega(x)} \; > \; \frac{|X|}{3r}.$$

This implies that there is a subset $Z$ of $X$ of size greater than $\frac{|X|}{3r}$ and a three-coloring $\chi$ such that every vector $z \in Z$ is equipartitioned by $\chi$. Then $Z$ may be regarded as a balanced subset of $\mathbb{F}^L \times \mathbb{F}^M \times \mathbb{F}^N$ where $L, M$ and $N$ are defined in Theorem 2.3. Applying Theorem 3.7 to $Z$ with $\lambda = L$, $\mu = M$ and $\nu = N$, we obtain

$$\frac{|X|}{3r} \; < \; 4k(LN)^{\frac{1}{2}}M^{\frac{1}{2}+\frac{1}{2k}} + L^{1+\frac{1}{2k}} + M^{1+\frac{1}{2k}} + N^{1+\frac{1}{2k}} + 2LN^{\frac{1}{2}},$$

which implies that

$$|X| \; < \; 12kr \cdot [(LN)^{\frac{1}{2}}M^{\frac{1}{2}+\frac{1}{2k}} + N^{1+\frac{1}{2k}}].$$

This is precisely the statement of Theorem 2.3. $\qquad\square$

**Lemma 3.8.** *Let $\omega$ be a positive integer. Then*

$$\binom{\omega}{\lfloor\frac{\omega}{3}\rfloor}\binom{\omega - \lfloor\frac{\omega}{3}\rfloor}{\lceil\frac{\omega}{3}\rceil} > \frac{3^{\omega-1}}{\omega}.$$

*Proof.* Let $f(\omega)$ denote the expression on the left in the inequality above. It is not hard to verify that the result is true for $\omega \in \{1, 2, 3\}$. Suppose $\omega > 3$. Using the inequalities,

$$n^n e^{-n}(2\pi n)^{\frac{1}{2}} < n! < n^n e^{-n}(2\pi n)^{\frac{1}{2}} e^{\frac{1}{12n}},$$

which are valid for all positive integers $n$, we get that for all integers $s \geq 1$,

$$f(3s) = \binom{3s}{s}\binom{2s}{s} = \frac{(3s)!}{(s!)^3} > \frac{(3s)^{3s}e^{-3s}(6\pi s)^{1/2}}{s^{3s}e^{-3s}(2\pi s)^{3/2}e^{1/4s}} > \frac{3^{3s+\frac{1}{2}}}{2\pi s e^{1/4s}} > \frac{3^{3s}}{4s}.$$

This implies the required inequality when $\omega$ is a multiple of 3. We pass to general $\omega$ by noting that

$$f(3s+1) = \frac{3s+1}{s+1}f(3s) \quad \text{and} \quad f(3s+2) = \frac{(3s+2)(3s+1)}{(s+1)^2}f(3s).$$

In particular, for $s \geq 1$, $f(3s+1) \geq 2f(3s)$ and $f(3s+2) \geq 5f(3s)$, which implies the required inequality. $\qquad\square$

# 4 Product representations of squares

In this section we prove Theorem 1.3. Before doing so, we require the following simple lemma:

**Lemma 4.1.** *Let $n > 1$ be a positive integer. Then either $n$ has a prime factor larger than $N^{\frac{1}{2}}$, or $n = xyz$ where $x, y, z \leq N^{\frac{1}{2}}$.*

*Proof.* Let $n = p_1 p_2 \ldots p_r$ denote the prime factorization of $n$ into (not necessarily distinct) primes $p_i$ where $p_1 \geq p_2 \geq \cdots \geq p_r$. Suppose $p_1 \leq N^{\frac{1}{2}}$. Then we can find a set $X$ of $p_i$s whose product $x$ is at most $N^{\frac{1}{2}}$ but as close to $N^{\frac{1}{2}}$ as possible. Let $y$ be a prime factor of $n$ which isn't in $X$. Then $xy \geq N^{\frac{1}{2}}$, and we may take $z = n/(xy)$. $\qquad\square$

In what follows we denote by $\Pi(n)$ the set of all primes in $\{1, \ldots, n\}$, and let $\pi(n) = |\Pi(n)|$ be the usual prime counting function.

*Proof of Theorem 1.3.* Let $A \subseteq \{1, 2, \ldots, n\}$ be a set such that no product of at most $8k$ distinct elements of $A$ is a square. Denote by $B \subseteq A$ the set of integers in $A$ which have a prime factor larger than $N^{\frac{1}{2}}$, and write $C = A \setminus B$. By Lemma 4.1 we have that $C = \{a \in A : a = xyz, \ x, y, z \leq N^{\frac{1}{2}}\}$. Denote for $0 \leq i \leq \frac{1}{2} \log_2 n$,

$$P_i = \left\{p \in \Pi(n) : \ \frac{n}{2^{i+1}} < p \leq \frac{n}{2^i}\right\}.$$

Form a bipartite graph $G_i$ with parts $P_i$ and $\{1, \ldots 2^{i+1}\}$ such that $p \in P_i$ is joined to $q \in \{1, \ldots, 2^{i+1}\}$ if $pq \in B$. Then $G_i$ does not contain a cycle of length at most $8k$ since if for some $2 \leq \ell \leq 8k$, $p_1 q_1, q_1 p_2, p_2 q_2, \ldots, q_{\ell-1} p_1$ is such a cycle then $p_j q_j, p_{j+1} q_j \in A$ are distinct and their product is a square. It was proved in [24] that an $M$ by $N$ bipartite graph of girth at least $2k + 2$ has at most $(MN)^{\frac{1}{2}+\frac{1}{2k}} + M + N$ edges. Since $G_i$ has girth at least $8k + 2$, we deduce that

$$
\begin{aligned}
|E(G_i)| &\leq \left(2^{i+1}|P_i|\right)^{\frac{1}{2}+\frac{1}{8k}} + 2^{i+1} + |P_i| \\
&\leq \left(2^{i+1}\left[\pi\left(\frac{n}{2^i}\right) - \pi\left(\frac{n}{2^{i+1}}\right)\right]\right)^{\frac{1}{2}+\frac{1}{8k}} + 2^{i+1} + |P_i|.
\end{aligned}
$$

Adding these inequalities for $i = 0, \ldots, \lfloor \frac{1}{2} \log_2 n \rfloor$ gives

$$|B| = \sum_{i=0}^{\lfloor \frac{1}{2} \log_2 n \rfloor} |E(G_i)| \leq \pi(n) + O(n^{\frac{1}{2}+\frac{1}{8k}}).$$

We now estimate $|C|$. For each $t \in C$ fix a factorization $t = x_t y_t z_t$ with $x_t y_t, z_t \leq N^{\frac{1}{2}}$. We assume $x_t \geq y_t \geq z_t$, so in particular $z_t \leq n^{\frac{1}{3}}$. Denote

$$S = \left\{(i, j) : \ 0 \leq i \leq j : \ i + j + 2 \leq \frac{1}{3} \log_2 n\right\}$$

For $(i, j) \in S$ let $C_{ij}$ denote the set of all $t \in C$ such that

$$\frac{N^{\frac{1}{2}}}{2^{i+1}} < x_t \leq \frac{N^{\frac{1}{2}}}{2^i} \qquad \frac{N^{\frac{1}{2}}}{2^{j+1}} < y_t \leq \frac{N^{\frac{1}{2}}}{2^j} \qquad 1 \leq z_t \leq 2^{i+j+2}.$$

We now apply Theorem 3.7 with $\mathbb{F} = \mathbb{F}_2$ and $\lambda = 2^{i+j+2}$, $\mu = N^{\frac{1}{2}}/2^{i+1}$ and $\nu = N^{\frac{1}{2}}/2^{j+1}$. Each $t \in C_{ij}$ may be considered as a vector of weight three in $\mathbb{F}^\lambda \times \mathbb{F}^\mu \times \mathbb{F}^\nu$: if $t = x_t y_t z_t$ is the prescribed factorization of $t$ then the vector associated with $t$ is the vector of weight three with a one in positions $x_t$, $\mu + y_t$ and $\mu + \nu + z_t$, and zeros elsewhere. Clearly no $8k$ of these vectors are linearly dependent, otherwise the product of the corresponding elements of $C_{ij}$ is a square. By Theorem 3.7, we have that

$$
\begin{aligned}
|C_{ij}| \;&<\; 12kr(\lambda\mu\nu)^{\frac{1}{2}}\mu^{\frac{1}{2k}} + \lambda^{1+\frac{1}{2k}} + \mu^{1+\frac{1}{2k}} + \nu^{1+\frac{1}{2k}} + 2\lambda\nu^{\frac{1}{2}} \\
&\ll\; k\left(\frac{N^{\frac{1}{2}}}{2^{i+1}} \cdot \frac{N^{\frac{1}{2}}}{2^{j+1}} \cdot 2^{i+j+2}\right)^{\frac{1}{2}}\left(\frac{N^{\frac{1}{2}}}{2^{i+1}}\right)^{\frac{1}{2k}} + \left(\frac{N^{\frac{1}{2}}}{2^{j+1}}\right)^{\frac{1}{2}+\frac{1}{2k}} + 2^{i+j}\left(\frac{N^{\frac{1}{2}}}{2^{j+1}}\right)^{\frac{1}{2}} \\
&\ll\; kn^{\frac{1}{2}+\frac{1}{2k}} + 2^{i+\frac{j}{2}}n^{\frac{1}{4}}.
\end{aligned}
$$

Finally, we sum this inequality over all $(i, j) \in S$. We chose $\lambda, \mu, \nu$ carefully to ensure that the sum of the last term over $(i, j) \in S$ is $O(n^{\frac{1}{2}})$. Therefore we have

$$
|C| \ll kn^{\frac{1}{2}+\frac{1}{2k}} \cdot |S| + O(n^{\frac{1}{2}}) \ll kn^{\frac{1}{2}+\frac{1}{2k}}(\log n)^2.
$$

This concludes the proof of Theorem 1.3. $\qquad\square$

# 5  Acknowledgements

# References

[1] N. Alon, L. Babai, and A. Itai. A fast and simple randomized parallel algorithm for the maximal independent set problem. *J. Algorithms*, 7(4):567–583, 1986.

[2] N. Alon, O. Goldreich, J. Håstad, and R. Peralta. Simple constructions of almost $k$-wise independent random variables. *Random Structures Algorithms*, 3(3):289–304, 1992.

[3] N. Alon, S. Hoory, and N. Linial. The Moore bound for irregular graphs. *Graphs Combin.*, 18(1):53–57, 2002.

[4] N. Alon and J. H. Spencer. *The probabilistic method*. Wiley-Interscience Series in Discrete Mathematics and Optimization. Wiley-Interscience [John Wiley & Sons], New York, second edition, 2000. With an appendix on the life and work of Paul Erdős.

[5] N. Alon, R. Yuster, and U. Zwick. Color-coding. *J. Assoc. Comput. Mach.*, 42(4):844–856, 1995.

[6] N. Alon, R. Yuster, and U. Zwick. Finding and counting given length cycles. *Algorithmica*, 17(3):209–223, 1997.

[7] L. Bazzi, M. Mahdian, and D. A. Spielman. The minimum distance of Turbo-like codes. Preprint, 2003.

[8] C. Berrou, A. Glavieux, and P. Thitimajshima. Near Shannon Limit Error Correcting Codes and Decoding: Turbo Codes. In *Proceedings of IEEE International Communications Conference*, pages 1064–1070. 1993.

[9] C. Bertram-Kretzberg, T. Hofmeister, and H. Lefmann. Sparse 0-1 matrices and forbidden hypergraphs. *Combin. Probab. Comput.*, 8(5):417–427, 1999.

[10] C. Bertram-Kretzberg and H. Lefmann. $MOD_p$-tests, almost independence and small probability spaces. *Random Structures Algorithms*, 16(4):293–313, 2000.

[11] A. Beutelspacher and U. Rosenbaum. *Projective geometry: from foundations to applications.* Cambridge University Press, Cambridge, 1998.

[12] B. Bollobás. *Extremal graph theory.* Dover Publications Inc., Mineola, NY, 2004. Reprint of the 1978 original.

[13] J. Bondy and M. Simonovits. Cycles of even length in graphs. *J. Combinatorial Theory B*, 16:97–105, 1974.

[14] M. Breiling. A logarithmic upper bound on the minimum distance of Turbo codes. Preprint, 2001.

[15] A. E. Brouwer. Block designs. In *Handbook of combinatorics, Vol. 1, 2*, pages 693–745. Elsevier, Amsterdam, 1995.

[16] M. C. Davey and D. J. C. MacKay. Low-density parity check codes over $GF(q)$. *IEEE Communications Letters*, 2(6):165–167, 1998.

[17] J. D. Dixon. Asymptotically fast factorization of integers. *Math. Comp.*, 36(153):255–260, 1981.

[18] Erdős, P. On some applications of graph theory to number theoretic problems. Publ. Ramanujan Inst. 1, 131–136, 1969.

[19] Erdős, Brown, W. G. and Sós, V. T. Some extremal problems on $r$-graphs. New directions in the theory of graphs. Proc 3rd Ann Arbor Conference on Graph Theory, Academic Press, New York, 55–63, 1973.

[20] P. Erdős and D. J. Kleitman. On coloring graphs to maximize the proportion of multicolored $k$-edges. *J. Combinatorial Theory*, 5:164–169, 1968.

[21] P. Erdős, A. Sárközy, and V. T. Sós. On product representations of powers. I. *European J. Combin.*, 16(6):567–588, 1995.

[22] R. G. Gallager. *Low Density Parity Check Codes.* MIT Press, Cambridge MA, 1963. Research Monograph Series, no. 21.

[23] E. Györi. $C_6$-free bipartite graphs and product representation of squares. *Discrete Math.*, 165/166:371–375, 1997. Graphs and combinatorics (Marseille, 1995).

[24] S. Hoory. The size of bipartite graphs with a given girth. *J. Combin. Theory Ser. B*, 86(2):215–220, 2002.

[25] N. Kahale and R. Urbanke. On the minimum distance of parallel and serially concatenated codes. *IEEE Trans. Inform. Theory*. To appear.

[26] H. Lefmann. Sparse parity-check matrices over finite fields (extended abstract). In *Computing and combinatorics*, volume 2697 of *Lecture Notes in Comput. Sci.*, pages 112–121. Springer, Berlin, 2003.

[27] H. Lefmann, P. Pudlák, and P. Savický. On sparse parity check matrices. *Des. Codes Cryptogr.*, 12(2):107–130, 1997.

[28] D. J. C. MacKay. Good error-correcting codes based on very sparse matrices. *IEEE Trans. Inform. Theory*, 45(2):399–431, 1999.

[29] F. J. MacWilliams and N. J. A. Sloane. *The theory of error-correcting codes. I.* North-Holland Publishing Co., Amsterdam, 1977. North-Holland Mathematical Library, Vol. 16.

[30] F. J. MacWilliams and N. J. A. Sloane. *The theory of error-correcting codes. II.* North-Holland Publishing Co., Amsterdam, 1977. North-Holland Mathematical Library, Vol. 16.

[31] A. Naor and J. Verstraëte. A note on bipartite graphs without a $2k$-cycle. Preprint, 2003.

[32] C. Pomerance. A tale of two sieves. *Notices Amer. Math. Soc.*, 43(12):1473–1485, 1996.

[33] C. Pomerance and A. Sárközy. Combinatorial number theory. In *Handbook of combinatorics, Vol. 1, 2*, pages 967–1018. Elsevier, Amsterdam, 1995.

[34] G. N. Sárközy. Cycles in bipartite graphs and an application in number theory. *J. Graph Theory*, 19(3):323–331, 1995.

[35] M. Sipser and D. A. Spielman. Expander codes. *IEEE Trans. Inform. Theory*, 42(6, part 1):1710–1722, 1996. Codes and complexity.

[36] D. A. Spielman. Linear-time encodable and decodable error-correcting codes. *IEEE Trans. Inform. Theory*, 42(6, part 1):1723–1731, 1996. Codes and complexity.

[37] J. Verstraëte. On arithmetic progressions of cycle lengths in graphs. *Combin. Probab. Comput.*, 9(4):369–373, 2000.

[38] H. C. Williams and J. O. Shallit. Factoring integers before computers. In *Mathematics of Computation 1943–1993: a half-century of computational mathematics (Vancouver, BC, 1993)*, volume 48 of *Proc. Sympos. Appl. Math.*, pages 481–531. Amer. Math. Soc., Providence, RI, 1994.