

Separating hash families: A Johnson-type bound and new constructions

Chong Shangguan^a, Gennian Ge^{b,c,*}

^a School of Mathematical Sciences, Zhejiang University, Hangzhou 310027, Zhejiang, China

^b School of Mathematical Sciences, Capital Normal University, Beijing 100048, China

^c Beijing Center for Mathematics and Information Interdisciplinary Sciences, Beijing 100048, China

Abstract

Separating hash families are useful combinatorial structures which are generalizations of many well-studied objects in combinatorics, cryptography and coding theory. In this paper, using tools from graph theory and additive number theory, we solve several open problems and conjectures concerning bounds and constructions for separating hash families.

Firstly, we discover that the cardinality of a separating hash family satisfies a Johnson-type inequality. As a result, we obtain a new upper bound, which is superior to all previous ones.

Secondly, we present a construction for an infinite class of perfect hash families. It is based on the Hamming graphs in coding theory and generalizes many constructions that appeared before. It provides an affirmative answer to both Bazrafshan-Trung's open problem on separating hash families and Alon-Stav's conjecture on parent-identifying codes.

Thirdly, let $p_t(N, q)$ denote the maximal cardinality of a t -perfect hash family of length N over an alphabet of size q . Walker and Colbourn conjectured that $p_3(3, q) = o(q^2)$. We verify this conjecture by proving $q^{2-o(1)} < p_3(3, q) = o(q^2)$. Our proof can be viewed as an application of Ruzsa-Szemerédi's (6,3)-theorem. We also prove $q^{2-o(1)} < p_4(4, q) = o(q^2)$. Two new notions in graph theory and additive number theory, namely rainbow cycles and R -sum-free sets, are introduced to prove this result. These two bounds support a question of Blackburn, Etzion, Stinson and Zaverucha.

Finally, we establish a bridge between perfect hash families and hypergraph Turán problems. This connection has not been noticed before. As a consequence, many new results and problems arise.

Keywords: separating hash family, perfect hash family, Johnson-type bound, rainbow cycle, R -sum-free set.

Mathematics subject classifications: 05B30, 94A60, 68R05, 94B60

1 Introduction

Separating hash families are useful combinatorial structures introduced by Stinson, Wei and Chen [38]. They are generalizations of many combinatorial objects, for example, perfect hash families, frameproof codes and codes with the identifiable parent property.

Let us begin with some definitions.

Definition 1.1. Let X and Y be sets of cardinalities n and q , respectively. We call a set \mathcal{F} of N functions $f : X \rightarrow Y$ an $(N; n, q)$ -hash family.

Definition 1.2. Let $f : X \rightarrow Y$ be a function, and let pairwise disjoint subsets $C_1, C_2, \dots, C_t \subseteq X$. We say that f separates C_1, C_2, \dots, C_t if $f(C_1), \dots, f(C_t)$ are pairwise disjoint. In particular, we say that f separates a subset $C \subseteq X$ if $f(C) \subseteq Y$ has $|C|$ distinct values.

Definition 1.3. Let X and Y be sets of cardinalities n and q , respectively, and let \mathcal{F} be an $(N; n, q)$ -hash family of functions from X to Y . We say that \mathcal{F} is an $(N; n, q, \{w_1, \dots, w_t\})$ -separating hash family (which we will also denote as an $SHF(N; n, q, \{w_1, \dots, w_t\})$) if it satisfies

*Corresponding author. Email address: gnge@zju.edu.cn. Research supported by the National Natural Science Foundation of China under Grant Nos. 11431003 and 61571310.

the following property: for all pairwise disjoint subsets $C_1, C_2, \dots, C_t \subseteq X$ with $|C_i| = w_i$ for $1 \leq i \leq t$, there exists at least one function $f \in \mathcal{F}$ that separates C_1, C_2, \dots, C_t . We call the multiset $\{w_1, \dots, w_t\}$ the type of this separating hash family.

For a positive integer q , we denote $[q]$ for the set $\{1, \dots, q\}$. Without loss of generality, we may fix the alphabet set Y to be the set of first q positive integers. And for the sake of simplicity, we set $u = \sum_{i=1}^t w_i$ throughout this paper. To avoid trivial cases, we assume that $n > q$, $q \geq t \geq 2$ and $u \leq n$.

The concept of separating hash families was first introduced in the special case $t = 2$ by Stinson, Trung and Wei [36] and then generalized by Stinson, Wei and Chen [38]. This notion has relations with many well-studied objects in combinatorics, cryptography and coding theory, see [18, 38] for a detailed introduction. We will summarise some objects in which we are interested.

- If $w_1 = w_2 = \dots = w_t = 1$, an $SHF(N; n, q, \{1, \dots, 1\})$ is known as a t -perfect hash family, which will be denoted as $PHF(N; n, q, t)$. Perfect hash families are basic combinatorial structures and have important applications in cryptography [14, 17, 35, 36], database management [30], circuit design [31] and the design of deterministic analogues of probabilistic algorithms [4].
- If $t = 2$ with $w_1 = 1$ and $w_2 = w$, an $SHF(N; n, q, \{1, w\})$ is known as a w -frameproof code. The frameproof code is a kind of fingerprinting codes and has applications in the protection of copyrighted materials. See [15, 19, 35, 37] for results on frameproof codes.
- Codes with the *identifiable parent property* (or 2 -*IPP* codes) are separating hash families which are simultaneously of type $\{1, 1, 1\}$ and $\{2, 2\}$, see [1, 3, 6, 16, 27].

Bounds and constructions for separating hash families are central problems in this research area. Given positive integers N , q and w_1, \dots, w_t , it is of interest how large the cardinality n of the preimage set X can be. We use $C(N, q, \{w_1, \dots, w_t\})$ to denote this maximal cardinality.

By a method known as grouping coordinates, the problem of bounding $C(N, q, \{w_1, \dots, w_t\})$ can be reduced to bounding $C(u - 1, q, \{w_1, \dots, w_t\})$, since it has been observed in [7, 18, 38] that $C(N, q, \{w_1, \dots, w_t\}) \leq C(u - 1, q^{\lceil N/(u-1) \rceil}, \{w_1, \dots, w_t\})$.

In the literature, researchers are seeking for the minimal positive real number γ such that $C(u - 1, q, \{w_1, \dots, w_t\}) \leq \gamma q$ holds for arbitrary q . The reader is referred to [7, 18, 35, 36, 38] for the attempts that have been made. In 2008, Stinson, Wei and Chen [38] proved $C(3, q, \{1, 1, 2\}) \leq 3q + 2 - 2\sqrt{3m + 1}$ and $C(3, q, \{2, 2\}) \leq 4q - 3$ for two special cases. In the same year, Blackburn, Etizon, Stinson and Zaverucha [18] proved $C(u - 1, q, \{w_1, \dots, w_t\}) \leq (w_1 w_2 + u - w_1 - w_2)q$, where $w_1, w_2 \leq w_i$ for $3 \leq i \leq t$. In 2011, Bazrafshan and Trung [7] proved the following theorem:

Theorem 1.4. ([7]) $C(u - 1, q, \{w_1, \dots, w_t\}) \leq (u - 1)q$.

Moreover, they conjectured that (see Question 1.6) $\gamma = u - 1$ is the minimal real number such that the above bound holds for arbitrary q .

We improve Theorem 1.4 in various aspects, including some tighter bounds and asymptotically optimal constructions. The novelty of our work is that we develop two new approaches to study bounds and constructions for codes and hash families with the separating property. We will explain them in detail in the conclusion section of this paper.

We state our main results as follows.

1.1 Separating hash families

Following the steps of previous papers [7, 18, 38], we discover an important property for separating hash families that the growth of $C(N, q, \{w_1, \dots, w_t\})$ satisfies a Johnson-type inequality. Roughly speaking, $C(N, q, \{w_1, \dots, w_t\}) \leq q^l + \max\{u - 1, C(N - l, q, \{w_1 - 1, \dots, w_t\})\}$ holds for every positive integer l (see Lemma 3.1 below). As a result, we obtain the following new upper bound for separating hash families which is the best known one.

Theorem 1.5. *Suppose there exists an $SHF(N; n, q, \{w_1, \dots, w_t\})$. Let $u = \sum_{i=1}^t w_i$ and let $1 \leq r \leq u - 1$ be the positive integer such that $N \equiv r \pmod{u - 1}$. If $C(\lfloor N/(u - 1) \rfloor, q, \{w_1, \dots, w_t\}) \geq u$, then it holds that $n \leq rq^{\lceil N/(u-1) \rceil} + (u - 1 - r)q^{\lfloor N/(u-1) \rfloor}$.*

A novelty of our proof is that we avoid the use of the grouping coordinates method, which has appeared in all previous proofs. The constraint $C(\lfloor N/(u-1) \rfloor, q, \{w_1, \dots, w_t\}) \geq u$ can be omitted when $N \geq u-1$ and $q \geq u$.

For the coefficient γ defined in Theorem 1.4, the authors of [7] posed the following question:

Question 1.6. ([7]) *Is there any type $\{w_1, \dots, w_t\}$ for which the constant $(u-1)$ in Theorem 1.4 can be replaced by another constant strictly smaller than $(u-1)$?*

We give a negative answer to their question by presenting the following construction:

Theorem 1.7. *There exists a PHF($N; Nq^{N-1}, q^{N-1} + (N-1)q^{N-2}, N+1$) for any integer $q \geq 2$ and $N \geq 2$. As a consequence, $\gamma = u-1$ is the minimal real number such that $C(u-1, q, \{w_1, \dots, w_t\}) \leq \gamma q$ holds for arbitrary q .*

To see that our construction is actually a negative answer to Question 1.6, one just needs to notice that a u -perfect hash family is also $\{w_1, \dots, w_t\}$ -separating for arbitrary $\sum_{i=1}^t w_i = u$. If we set $N = u-1$ then our construction implies the existence of an SHF($u-1; n, q, \{w_1, \dots, w_t\}$) such that $\lim_{q \rightarrow \infty} \frac{n}{q} = u-1$ holds for arbitrary $\sum_{i=1}^t w_i = u$. Therefore, the constant γ can never be less than $u-1$.

1.2 Codes with the identifiable parent property

We have mentioned 2-IPP codes before and the notion was generalized to codes with the t -identifiable parent property (t -IPP codes) in [35]. We postpone the definition to Section 2 for the sake of saving space.

Let $i_t(N, q)$ denote the maximal cardinality of a t -IPP code of length N over an alphabet of size q . Let $v = \lfloor (t/2 + 1)^2 \rfloor$. One can verify that $i_t(N, q) \leq i_t(v-1, q^{\lceil N/(v-1) \rceil})$ (just as the case for separating hash families). Thus the problem of bounding $i_t(N, q)$ can be reduced to bounding $i_t(v-1, q)$. Alon and Stav [6] proved that $i_t(v-1, q) \leq (v-1)q$, and they conjectured:

Conjecture 1.8. ([6]) *There are constructions showing that $(v-1)$ is the best constant in the inequality $i_t(v-1, q) \leq (v-1)q$.*

Our Theorem 1.7 not only answers Question 1.6 but also verifies this conjecture, since it was observed in [6, 35] that a v -perfect hash family also satisfies the t -identifiable parent property.

1.3 Perfect hash families

As claimed in [18], the exponent $\lceil N/(u-1) \rceil$ in the bound of Theorem 1.5 is realistic. We can understand this in two aspects. On the one hand, a probabilistic construction of Blackburn [13] showed that for any fixed u and any positive real number δ such that $\delta < N/(u-1)$, there exists a PHF($N; \lfloor q^\delta \rfloor, q, u$) whenever q is sufficiently large. On the other hand, let $p_t(N, q)$ denote the maximal cardinality of a PHF($N; n, q, t$), it was respectively observed in [6, 28, 32] that $p_u(N, q) \geq (c_u q)^{N/(u-1)}$ holds for some constant c_u . So we can conclude that the exponent $\lceil N/(u-1) \rceil$ is tight when $(u-1) \mid N$.

But the problem becomes much more difficult when $(u-1) \nmid N$. It is not known that whether the exponent is tight. Even for the smallest case, $u=3$ and $N=3$, Walker and Colbourn [39] posed the following conjecture:

Conjecture 1.9. ([39]) $p_3(3, q) = o(q^2)$.

Note that Theorem 1.5 shows $p_3(3, q) = O(q^2)$. A recent paper [24] showed that $p_3(3, q) = \Omega(q^{5/3})$. Results from finite geometry were used to construct such families. There is still a huge gap between the upper and lower bounds. For general types of separating hash families, Blackburn et al. [18] asked a similar question:

Question 1.10. ([18]) *Let N and w_i be fixed integers. If $(u-1) \nmid N$, then for sufficiently large q and arbitrary small $\epsilon > 0$, does there exist an SHF($N; n, q, \{w_1, \dots, w_t\}$) such that $n \geq q^{\lceil N/(u-1) \rceil - \epsilon}$?*

We prove Conjecture 1.9 in Section 5 (see Theorems 5.4 and 5.7 below). We find that perfect hash families are closely related to a hypergraph Turán problem. With some transformations, Walker-Colbourn’s conjecture can be proved by a direct application of the famous (6,3)-theorem of Ruzsa and Szemerédi [34]. In fact, we show

$$q^{2-\epsilon} < p_3(3, q) = o(q^2)$$

holds for sufficiently large q and arbitrary $\epsilon > 0$. We also prove

$$q^{2-\epsilon} < p_4(4, q) = o(q^2)$$

(see Theorem 6.5 below). Two new notions in graph theory and additive number theory, namely rainbow cycles and R -sum-free sets, are introduced to prove this result. One can see that these two bounds suggest that there may be a positive answer to Question 1.10.

1.4 Organization

The rest of this paper is organised as follows. Section 2 is for some preparations. Theorem 1.5 is proved in Section 3 and Theorem 1.7 is proved in Section 4. The subsequent sections will focus on perfect hash families. We prove $q^{2-o(1)} < p_3(3, q) = o(q^2)$ in Section 5. As an application of the Johnson-type bound, this result will be extended to $p_t(t, q)$ and related separating hash families. We prove $q^{2-o(1)} < p_4(4, q) = o(q^2)$ in Section 6. In Section 7 we will build the connection between perfect hash families and a class of hypergraph Turán problems. Section 8 consists of some concluding remarks and open problems.

2 Preliminaries

In this section, we will introduce some notations and terminology. We will also introduce some simple lemmas that will be used in the subsequent sections.

2.1 Separating hash families and IPP codes

The matrix representation of a separating hash family is very useful when discussing its properties. An $(N; n, q)$ -hash family can be described as an $N \times n$ matrix on q symbols, which will be usually denoted as M . The rows of M correspond to the functions in the hash family and the columns of M correspond to the elements of X . The entry of M in row $f \in \mathcal{F}$ and column $x \in X$ is just $f(x) \in Y$. We denote the entry of M as $M(f, x)$ for $f \in \mathcal{F}$, $x \in X$ or $M(i, j)$ for $1 \leq i \leq N$, $1 \leq j \leq n$.

The matrix representation of an $SHF(N; n, q, \{w_1, \dots, w_t\})$ satisfies the following property: given disjoint sets of columns C_1, \dots, C_t , where $|C_i| = w_i$ for $1 \leq i \leq t$, there exists a row r of M such that

$$\{M(r, x) : x \in C_i\} \cap \{M(r, x) : x \in C_j\} = \emptyset$$

for all $i \neq j$. We say row r separates a subset of columns $C \subseteq X$ if $\{M(r, x) : x \in C\}$ has exactly $|C|$ distinct values in Y . The column x of M will be written as a q -ary vector of length N , $x = (x(1), x(2), \dots, x(N))$, where $x(i) \in [q]$ for $i \in [N]$. For a subset L of the rows of M , the coordinates of x restricted to L give a word of length $|L|$, which is denoted as $x|_L = (x(i_1), x(i_2), \dots, x(i_{|L|}))$, where i_j , $1 \leq j \leq |L|$ are the row indices. We say a column $x \in X$ of M has a unique coordinate i if for any other column $y \in X$, $y \neq x$, it holds that $y(i) \neq x(i)$. If there is no confusion, we will not distinguish between a hash family and its representation matrix.

Next we will introduce the definition of IPP codes.

Let $\mathcal{C} \subseteq Y^N$ be a code of length N and let $D \subseteq \mathcal{C}$ be a set of codewords. The set of descendants of D , denoted as $desc(D)$, is defined by

$$desc(D) = \{d \in Y^N : \text{for all } i \in \{1, 2, \dots, N\}, d(i) = x(i) \text{ for some } x \in D\}.$$

A set $D \subseteq \mathcal{C}$ is said to be a parent set of a word $d \in Y^N$ if $d \in desc(D)$. For $d \in Y^N$, let $\mathcal{P}_t(d)$ denote the collection of parent sets of d such that $|D| \leq t$ and $D \subseteq \mathcal{C}$. Then we call $\mathcal{C} \subseteq Y^N$ a t -IPP code if for all $d \in Y^N$, either $\mathcal{P}_t(d) = \emptyset$ or

$$\bigcap_{D \in \mathcal{P}_t(d)} D \neq \emptyset.$$

2.2 Graph theory

We will use the notion of Hamming graphs when constructing perfect hash families in Section 4. Let k and q be positive integers, the Hamming graph (see [26] for details) $H(k, q)$ has the set of all k -tuples from an alphabet of q symbols as its vertex set, and two k -tuples are adjacent if and only if they differ in exactly one coordinate position. This graph is also known as the q -ary hypercube of dimension k . Here we will fix this q -symbol alphabet set to be $[q]$.

When speaking about a hypergraph we mean a pair $\mathcal{G} = (V(\mathcal{G}), E(\mathcal{G}))$, where the vertex set $V(\mathcal{G})$ is identified as the set of first integers $[n]$ and the edge set $E(\mathcal{G})$ is identified as a collection of subsets of $[n]$. \mathcal{G} is said to be linear if for all distinct $A, B \in E(\mathcal{G})$ it holds that $|A \cap B| \leq 1$. We say \mathcal{G} is r -uniform if $|A| = r$ for all $A \in E(\mathcal{G})$.

An r -uniform hypergraph \mathcal{G} is r -partite if its vertex set $V(\mathcal{G})$ can be colored in r colors in such a way that no edge of \mathcal{G} contains two vertices of the same color. In such a coloring, the color classes of $V(\mathcal{G})$, the sets of all vertices of the same color, are called parts of \mathcal{G} . In this paper we mainly concern r -uniform r -partite hypergraphs with equal part size q . We will see later that the edge set of such hypergraph is equivalent to an $r \times |E(\mathcal{G})|$ matrix over a q -symbol alphabet.

Given a set \mathcal{H} of r -uniform hypergraphs, an \mathcal{H} -free r -uniform hypergraph is a graph containing none of the members of \mathcal{H} . The Turán number $ex_r(n, \mathcal{H})$ denotes the maximum number of edges in an \mathcal{H} -free r -uniform hypergraph on n vertices. In this paper, we will talk about several hypergraph Turán problems.

Brown, Erdős and Sós [20, 21] introduced the function $f_r(n, v, e)$ to denote the maximum number of edges in an r -uniform hypergraph on n vertices which does not contain e edges spanned by v vertices. In other words, in such hypergraphs the size of the union of arbitrary e edges is at least $v + 1$. These hypergraphs are called $G(v, e)$ -free (more precisely, $G_r(v, e)$ -free). The famous (6,3)-theorem of Ruzsa and Szemerédi [34] pointed out that

$$n^{2-o(1)} < f_3(n, 6, 3) = o(n^2). \quad (1)$$

This was extended by Alon and Shapira [5] to

$$n^{k-o(1)} < f_r(n, 3(r-k) + k + 1, 3) = o(n^k). \quad (2)$$

These bounds will be used when considering problems about perfect hash families in the sequel. For more results on Turán problems of this type, see [25] and the references therein.

The definitions of $f_r(n, v, e)$ can be restricted to the case for r -uniform r -partite hypergraphs with equal part size q . We use $f_r^*(q, v, e)$ to denote the corresponding formula. Note that $f_r^*(q, v, e) \leq f_r(rq, v, e)$.

In the literature, there are several definitions of hypergraph cycles. The one we use in this paper was introduced by Berge [10, 11]. For $k \geq 2$, a cycle in a hypergraph \mathcal{G} is an alternating sequence of vertices and edges of the form $v_1, E_1, v_2, E_2, \dots, v_k, E_k, v_1$ such that

- (a) v_1, v_2, \dots, v_k are distinct vertices of \mathcal{G} ,
- (b) E_1, E_2, \dots, E_k are distinct edges of \mathcal{G} ,
- (c) $v_i, v_{i+1} \in E_i$ for $1 \leq i \leq k-1$ and $v_k, v_1 \in E_k$.

Next we will introduce the definition of rainbow cycles. Note that in the literature “rainbow cycles” always stand for edge-colorings, but in this paper we consider rainbow cycles due to vertex colorings. Given a hypergraph \mathcal{G} and a vertex-coloring of \mathcal{G} , a subgraph $\mathcal{H} \subseteq \mathcal{G}$ is called a rainbow subgraph of \mathcal{G} if all joint vertices in \mathcal{H} have different colors. In other words, for arbitrary distinct vertices $x, y \in \{A \cap B : A, B \in E(\mathcal{H})\}$, x and y are colored by different colors. This definition is most meaningful when discussing linear hypergraphs. Let \mathcal{G} be an r -uniform r -partite linear hypergraph, a k -cycle $v_1, E_1, v_2, E_2, \dots, v_k, E_k, v_1$ is said to be a rainbow cycle of \mathcal{G} if v_1, \dots, v_k locate in different parts of $V(\mathcal{G})$. For r -partite graphs, a rainbow k -cycle exists only if $k \leq r$.

Let \mathcal{G} be an r -uniform r -partite linear hypergraph with equal part size q . Assume that \mathcal{G} can not have rainbow cycles, then we use $g_r^*(q)$ to denote the maximal number of edges that can be contained in \mathcal{G} . Lemma 6.1 shows that hypergraphs with large $g_r^*(q)$ can be used to construct good perfect hash families.

2.3 Additive number theory

It has been shown in [3, 25] that tools from additive number theory can be used to construct codes with some specified properties. We will introduce some notions from additive number theory.

Assume $m_1, m_2, m_3 \in M \subseteq [q]$ and c_1, c_2 are positive integers such that $c_1 + c_2 \leq r$, we call the set M r -sum-free if the equation

$$c_1 m_1 + c_2 m_2 = (c_1 + c_2) m_3$$

has no solution except the one with $m_1 = m_2 = m_3$. A result proved by Erdős, Frankl and Rödl [22] and Ruzsa [33] will be needed.

Lemma 2.1. ([22, 33]) *For arbitrary positive integer r there exists a $\gamma_r > 0$ such that for any integer q , one can find an r -sum-free subset $M \subseteq [q]$ with $|M| > qe^{-\gamma_r \sqrt{\log q}}$.*

Note that the case $r_1 = r_2 = 1$ was originally proved by Behrend [9].

A linear equation with integer coefficients

$$\sum_{i=1}^k a_i x_i = 0$$

in the unknowns x_i is homogeneous if $\sum_{i=1}^k a_i = 0$. We say that $M \subseteq [q]$ has no nontrivial solution to above equation, if whenever $m_i \in M$ and $\sum_{i=1}^k a_i m_i = 0$, it follows that all m_i 's are equal. Note that if M has no nontrivial solution to above function, then the same holds for any shift $(M + x) \cap [q]$ with $x \in \mathbb{Z}$, where $M + x := \{m + x : m \in M\}$. This property suggests that one can use probabilistic method to construct sets with no nontrivial solution to a system of homogeneous linear equations. Note that this definition of the nontrivial solution is a simplification of the original one of Ruzsa [33].

Now we will generalize the definition of the r -sum-free set. Given a set $R = \{b_1, \dots, b_r\}$ of r distinct nonnegative integers. A set M is said to be R -sum-free if for any $3 \leq k \leq r$ and any k -element subset $S = \{b_{j_1}, b_{j_2}, \dots, b_{j_k}\} \subseteq R$, the equation

$$(b_{j_2} - b_{j_1})m_1 + (b_{j_3} - b_{j_2})m_2 + \dots + (b_{j_k} - b_{j_{k-1}})m_{k-1} + (b_{j_1} - b_{j_k})m_k = 0$$

has no solution in M except the trivial one $m_1 = m_2 = \dots = m_k$. The rank of R is defined to be the maximal difference between the elements of R :

$$r(R) = \max_{1 \leq i < j \leq r} |b_i - b_j|.$$

We are interested in R -sum-free sets $M \subseteq [q]$ with relatively small rank, namely, $r(R) = o(q^\epsilon)$ for arbitrary $\epsilon > 0$. Lemma 6.2 shows that R -sum-free sets can be used to construct hypergraphs with large $g_r^*(q)$.

2.4 Some lemmas

The following lemma is a variant of a result of Erdős and Kleitman [23].

Lemma 2.2. *Every r -uniform hypergraph \mathcal{G} contains an r -uniform r -partite hypergraph \mathcal{H} with equal part size q or $q + 1$ such that*

$$\frac{|E(\mathcal{H})|}{|E(\mathcal{G})|} \geq \frac{r!}{r^r}.$$

Proof. Let $|V(\mathcal{G})| = n$, take q to be the integer such that $rq \leq n < r(q+1)$. We only prove the lemma for $n = rq$, otherwise we can set the part size of the desired subgraph to be $q+1$. It suffices to find a partition π of $V(\mathcal{G})$ with $\pi = \{B_1, \dots, B_r\}$ and $|B_i| = q$ for $1 \leq i \leq r$, such that $\mathcal{F}_\pi = \{A \in E(\mathcal{G}) : |A \cap B_i| = 1 \text{ for all } 1 \leq i \leq r\}$ contains the desired number of edges. Let $P(\mathcal{G})$ denote the collection of all appropriate partitions of $V(\mathcal{G})$. Let us count the number of the pairs $N := |\{(A, \pi) : A \in E(\mathcal{G}), \pi \in P(\mathcal{G}), |A \cap B_i| = 1 \text{ for every } B_i \in \pi\}|$. One can compute that any $A \in E(\mathcal{G})$ is contained in $\frac{|P(\mathcal{G})| \cdot q^r}{\binom{rq}{r}}$ members of $P(\mathcal{G})$ satisfying the desired property. Therefore, by double counting, there exists a $\pi \in P(\mathcal{G})$ such that \mathcal{F}_π contains at least

$$\frac{|E(\mathcal{G})| \cdot |P(\mathcal{G})| \cdot q^r / \binom{rq}{r}}{|P(\mathcal{G})|} = \frac{|E(\mathcal{G})| \cdot q^r}{\binom{rq}{r}}$$

members of $E(\mathcal{G})$. Then this specified π will induce an r -uniform r -partite hypergraph \mathcal{H} containing the desired number of edges. \square

This lemma implies that for any r -uniform hypergraph \mathcal{G} with sufficiently large $|V(\mathcal{G})|$, there exists an r -partite subgraph $\mathcal{H} \subseteq \mathcal{G}$ such that $|E(\mathcal{H})|$ and $|E(\mathcal{G})|$ are of the same order of magnitude. In other words, one can infer $f_r(rq, v, e) = \Theta(f_r^*(q, v, e))$ by Lemma 2.2.

Another simple lemma will be used.

Lemma 2.3. *Suppose G is a finite graph with n vertices. If G has no cycles, then G can have at most $n - 1$ edges.*

Proof. G must have a vertex with degree one since every path in G is finite and must have an end point. Choose a vertex in G with degree one, then the statement follows trivially by applying induction on $|V|$. \square

With some reformulations, one can combine Lemma 3.2 and Corollary 3.3 of Alon, Fischer and Szegedy [3] to prove the following result:

Lemma 2.4. ([3]) *There exists a set $M \subseteq \{0, 1, \dots, \lfloor (q-1)/(\mu+5) \rfloor\}$ satisfying*

$$|M| \geq qe^{-\gamma(\log q)^{3/4}}$$

such that M has no non-trivial solution to all the following equations

$$\begin{cases} 2m_1 + 3m_2 + \mu m_3 - (\mu + 5)m_4 & = 0 \\ 5m_1 + (\mu + 3)m_2 - 3m_3 - (\mu + 5)m_4 & = 0 \\ 5m_1 + \mu m_2 - 2m_3 - (\mu + 3)m_4 & = 0 \\ 2m_1 + 3m_2 - 5m_3 & = 0 \\ 5m_1 + \mu m_2 - (\mu + 5)m_3 & = 0 \\ 2m_1 + (\mu + 3)m_2 - (\mu + 5)m_3 & = 0 \\ 3m_1 + \mu m_2 - (\mu + 3)m_3 & = 0 \end{cases} \quad (3)$$

where γ is a constant and $\mu = \lceil 2\sqrt{\log q} \rceil$.

Sketch of the proof. Using the technique introduced in the proof of Lemma 3.2 of [3], for $1 \leq i \leq 7$, one can prove that there exists a set $M_i \subseteq \{0, 1, \dots, \lfloor (q-1)/(\mu+5) \rfloor\}$ and a constant γ_i satisfying

$$|M_i| \geq qe^{-\gamma_i(\log q)^{3/4}}$$

such that M_i has no nontrivial solution to the i -th equation in the above system. In order to prove the existence of the set M which has no nontrivial solution to all equations, we can apply a probabilistic method. Take six integers x_i such that $-\lfloor (q-1)/(\mu+5) \rfloor \leq x_i \leq \lfloor (q-1)/(\mu+5) \rfloor$, $2 \leq i \leq 7$, randomly, uniformly and independently. $M = M_1 \cap (M_2 + x_2) \cap \dots \cap (M_7 + x_7)$ has no nontrivial solution to any of the above equations. Since $M_i + x_i \in [-\lfloor (q-1)/(\mu+5) \rfloor, \lfloor (q-1)/(\mu+5) \rfloor]$ for each $2 \leq i \leq 7$, then one can compute that every $m \in M_1$ has probability at least $e^{-\sum_{i=2}^7 \gamma_i(\log q)^{3/4}}$ to lie in the intersection. Therefore, the result follows from the linearity of the expectation, where $|M| \geq qe^{-\gamma(\log q)^{3/4}}$ with $\gamma \leq \sum_{i=1}^7 \gamma_i$. \square

3 A Johnson-type upper bound

The aim of this section is to establish a Johnson-type bound for separating hash families and we will use it to prove Theorem 1.5. To establish this bound, the idea is to delete some rows and corresponding carefully chosen columns from the representation matrix of the separating hash family. Our goal is to show the remaining submatrix satisfies some weaker separating property. We call this recursive bound a “Johnson-type bound” due to its similarity with the traditional recursive Johnson bound in coding theory. Note that we always use M to denote the representation matrix of a separating hash family.

Lemma 3.1. *Let $1 \leq l \leq N$ be a positive integer, then it holds that $C(N, q, \{w_1, \dots, w_t\}) \leq q^l + \max\{u - 1, C(N - l, q, \{w_1 - 1, \dots, w_t\})\}$. In fact, in the right hand side of the inequality we can choose the minus of 1 to be after an arbitrary w_i , $1 \leq i \leq t$.*

Proof. Choose arbitrary l rows of M and let L denote the collection of these chosen rows. Denote $\mathcal{A} \subseteq Y^l$ the maximal collection of columns whose restrictions to L are all distinct (we just choose one column if there are several columns with the same restrictions to L). It is easy to see $|\mathcal{A}| \leq q^l$ since there are at most q^l distinct words of length l . Delete these l rows and the columns contained in \mathcal{A} from M . Let M' denote the remaining submatrix. Then M' is a q -ary $(N - l) \times (n - |\mathcal{A}|)$ matrix. If $n - |\mathcal{A}| \leq u - 1$, we are done. Otherwise it suffices to show M' is a representation matrix of a separating hash family of type $\{w_1, \dots, w_i - 1, \dots, w_t\}$ for arbitrary $1 \leq i \leq t$.

Assume the contrary, M' is not $\{w_1, \dots, w_i - 1, \dots, w_t\}$ -separating for some $1 \leq i \leq t$. Without loss of generality, we set $i = 1$. Then there exist t subsets C_1, \dots, C_t of the columns of M' with $|C_1| = w_1 - 1$ and $|C_i| = w_i$ for $2 \leq i \leq t$, such that no row of M' can separate C_1, \dots, C_t . Let c be an arbitrary column of C_2 and let c' be a column in \mathcal{A} such that $c'|_L = c|_L$. Such $c' \in \mathcal{A}$ must exist by our definition of \mathcal{A} . Consequently, no row can separate $C_1 \cup \{c'\}, C_2, \dots, C_t$ in the original matrix M , which contradicts the fact that M is $\{w_1, \dots, w_t\}$ -separating. Thus M' satisfies the desired separating property and the lemma follows from $n - |\mathcal{A}| \leq C(N - l, q, \{w_1 - 1, \dots, w_t\})$ and $|\mathcal{A}| \leq q^l$. \square

Remark 3.2. *This lemma is obviously an extension of Lemma 2 of [7]. We think this Johnson-type bound is very interesting and important since it points out the information hidden in the structure of separating hash families.*

As the first application of Lemma 3.1, we will use it to prove Theorem 1.5. Note that we can omit the constraint $C(\lfloor N/(u - 1) \rfloor, q, \{w_1, \dots, w_t\}) \geq u$ in the theorem by introducing a maximum term in the expression of the upper bound (just as the case in Lemma 3.1). And $C(\lfloor N/(u - 1) \rfloor, q, \{w_1, \dots, w_t\}) \geq u$ always holds for $N \geq u - 1$ and sufficiently large q , for example, $q \geq u$.

Proof of Theorem 1.5. One can verify that $N = r \lceil N/(u - 1) \rceil + (u - 1 - r) \lfloor N/(u - 1) \rfloor$. We apply Lemma 3.1 repeatedly for $u - 1$ times, in which l is chosen to be $\lceil N/(u - 1) \rceil$ (r times) and $\lfloor N/(u - 1) \rfloor$ ($u - 1 - r$ times), respectively. The theorem follows from a simple fact that $C(0, q, \{1\}) = 0$. \square

Remark 3.3. *It is not hard to see our bound is an improvement of Theorem 1.4, and hence an improvement of [7, 18]. One can see $\lceil N/(u - 1) \rceil$ is the best exponential term that can be obtained by our method, since to reduce the exponential term, one should reduce the maximum value of l involved in the deletions. In other words, we should find a finer partition of $[N]$ and hence more deletion rounds are needed. However, at most $(u - 1)$ deletion rounds can be used, because $C(N, q, \{w_1, \dots, w_t\})$ can be arbitrary large if $t = 1$ and $N > 0$.*

Remark 3.4. *Since the frameproof code is a special class of separating hash families, it is not surprising to see our bound contains Theorem 1 of [15] as a special case. By Constructions 2 and 3 in [15], one can find that Theorem 1.5 is asymptotically optimal when $q \geq N$, $\{w_1, \dots, w_t\} = \{1, w\}$, $N \equiv 1 \pmod{w}$, or $q = \Omega(N^2)$, $\{w_1, \dots, w_t\} = \{1, 2\}$. The following section presents a construction which shows Theorem 1.5 is also asymptotically optimal when $N = u - 1$.*

4 A construction for t -perfect hash families with $t - 1$ rows

The aim of this section is to present a construction for $PHF(N; Nq^{N-1}, q^{N-1} + (N - 1)q^{N-2}, N + 1)$ for arbitrary positive integer $q \geq 2$ and $N \geq 2$.

One nice feature of our construction is that it is a generalization of many previous ones. When $N = 2$, the construction of $PHF(2; 2q, q + 1, 3)$ has appeared in [29, 39]. And when $N = 3$, the construction of $PHF(3; 3q^2, q^2 + 2q, 4)$ has appeared in numerous papers, for example, Hollmann et al. [27], Blackburn [12], Stinson et al. [38] and Bazrafshan et al. [7].

Let us begin with $N = 3$ as a simple example to illustrate our idea.

Example 4.1. ([7, 12, 27, 38]) *There exists a $PHF(3; 3q^2, q^2 + 2q, 4)$ for any integer $q \geq 2$.*

Proof. We first construct a $3 \times q^2$ submatrix, in which the alphabet set is the $(q^2 + 2q)$ -element set defined as $\{(x, y), (x, 0), (0, y) : 1 \leq x, y \leq q, x, y \in \mathbb{Z}\}$,

$$\begin{pmatrix} (1, 1) & (1, 2) & \cdots & (1, q) & (2, 1) & \cdots & (2, q) & \cdots & (q, 1) & \cdots & (q, q) \\ (0, 1) & (0, 2) & \cdots & (0, q) & (0, 1) & \cdots & (0, q) & \cdots & (0, 1) & \cdots & (0, q) \\ (1, 0) & (1, 0) & \cdots & (1, 0) & (2, 0) & \cdots & (2, 0) & \cdots & (q, 0) & \cdots & (q, 0) \end{pmatrix}.$$

We denote the three rows of this submatrix as A_0, A_1, A_2 , respectively. Then the representation matrix of the desired perfect hash family can be presented as follows:

$$\begin{pmatrix} A_0 & A_2 & A_1 \\ A_1 & A_0 & A_2 \\ A_2 & A_1 & A_0 \end{pmatrix}.$$

We can easily see it is a $3 \times 3q^2$ matrix over an alphabet of size $q^2 + 2q$. One can verify (or see the proof of Theorem 1.7 below) that it is indeed a representation matrix of a 4-perfect hash family. \square

In the above matrix A_0 acts like an identity map that preserves each element in $\{(x, y) : 1 \leq x, y \leq q, x, y \in \mathbb{Z}\}$, while $A_i, i = 1, 2$, acts like a projection that projects the i -th entry of (x, y) to zero. Actually, the idea behind this simple construction can be generalized.

Recall the definition of the Hamming graphs in Section 2. Take a q -ary hypercube \mathcal{A} of dimension k , then $|V(\mathcal{A})| = q^k$. For $1 \leq i \leq k$ and arbitrary $\alpha = (\alpha(1), \dots, \alpha(k)) \in V(\mathcal{A})$, define π_i to be the map that sets $\alpha(i)$ into zero but preserves all other coordinates of α . We say π_i separates a set $S \subseteq V(\mathcal{A})$ if $\pi_i(\alpha) \neq \pi_i(\beta)$ for arbitrary distinct $\alpha, \beta \in S$. Proposition 1 of [12] establishes an important property of these maps. We present the proof here for the sake of reader's convenience.

Lemma 4.2. ([12]) *Let $S \subseteq V(\mathcal{A})$ be an arbitrary t -element subset with $t \leq k$, then S is separated by at least $k - t + 1$ of the functions π_1, \dots, π_k .*

Proof. Assume the contrary. Without loss of generality, let $S = \{\alpha_1, \dots, \alpha_t\}$ and let π_1, \dots, π_t be the t functions which can not separate S . Define a colored graph $G = (V, E)$ by $V = S$ and connect $\alpha, \beta \in V$ by an edge of color i if $\pi_i(\alpha) = \pi_i(\beta)$. Note that graph G is a subgraph of a Hamming graph. Since π_1, \dots, π_t can not separate S , then for every $i \in [t]$, there exist $1 \leq j < l \leq t$ such that $\pi_i(\alpha_j) = \pi_i(\alpha_l)$. So G contains a subgraph $G' = (V, E')$ with t vertices and t edges of distinct colors. By Lemma 2.3 we can deduce that G' contains a cycle which can be denoted as $(\alpha_1, \alpha_2, \dots, \alpha_c)$, where c is an integer such that $1 \leq c \leq t$.

We are done if we can show such cycle must not exist. Assume that the edge between α_1 and α_2 is colored by the i -th color. Then α_1 and α_2 must differ in their i -th coordinate. Since every edge in this cycle is of distinct color and every pair of connected vertices differ in exactly one coordinate, then for every $j \in \{2, 3, \dots, c\}$, α_j and α_{j+1} must agree in their i -th coordinate. In particular, α_{c+1} is recognised as α_1 , which implies that $\alpha_2(i) = \alpha_3(i) = \dots = \alpha_c(i) = \alpha_1(i)$. Thus the desired contradiction follows. \square

The following lemma is an easy consequence of above lemma.

Lemma 4.3. *Let π_i ($1 \leq i \leq k$) be the functions defined as above and let π_0 denote the identity map which satisfies $\pi_0(\alpha) = \alpha$ for every $\alpha \in V(\mathcal{A})$. Suppose $S \subseteq V(\mathcal{A})$ is a t -element subset with $t \leq k + 1$, then at most $t - 1$ of the functions $\pi_0, \pi_1, \dots, \pi_k$ can not separate S .*

Proof. Apply Lemma 4.2 and remember the fact that π_0 separates every subset of $V(\mathcal{A})$. \square

Now we can prove Theorem 1.7.

Proof of Theorem 1.7. Take a q -ary hypercube \mathcal{A} of dimension $N - 1$. Obviously $|V(\mathcal{A})| = q^{N-1}$. Let $\pi_0, \pi_1, \dots, \pi_{N-1}$ be the maps defined as above. Then our desired perfect hash family can be represented as the following matrix

$$\begin{pmatrix} \pi_0(\mathcal{A}) & \pi_{N-1}(\mathcal{A}) & \cdots & \cdots & \pi_1(\mathcal{A}) \\ \pi_1(\mathcal{A}) & \pi_0(\mathcal{A}) & \cdots & \cdots & \pi_2(\mathcal{A}) \\ \vdots & \vdots & \ddots & & \vdots \\ \vdots & \vdots & & \ddots & \vdots \\ \pi_{N-1}(\mathcal{A}) & \pi_{N-2}(\mathcal{A}) & \cdots & \cdots & \pi_0(\mathcal{A}) \end{pmatrix},$$

where for every $0 \leq i \leq N - 1$, $\pi_i(\mathcal{A}) := (\pi_i(\alpha))_{\alpha \in V(\mathcal{A})}$ is a $1 \times |V(\mathcal{A})|$ submatrix. Denote this representation matrix as M , then M is an $N \times Nq^{N-1}$ matrix. Let $Y = \cup_{i=0}^{N-1} \pi_i(\mathcal{A})$ denote the alphabet set. It is not hard to see $|\{\pi_0(\alpha) : \alpha \in V(\mathcal{A})\}| = q^{N-1}$ and $|\{\pi_i(\alpha) : \alpha \in V(\mathcal{A})\}| = q^{N-2}$ for every $1 \leq i \leq N - 1$. Then one can verify that $|Y| = q^{N-1} + (N - 1)q^{N-2}$. Thus we can conclude that M is the representation matrix of an $(N; Nq^{N-1}, q^{N-1} + (N - 1)q^{N-2})$ -hash family.

Now it remains to verify that this hash family is indeed an $(N + 1)$ -perfect hash family. Consider M as the concatenation of N column patterns denoted as $(C_1|C_2|\cdots|C_N)$ with $|C_1| = |C_2| = \cdots = |C_N| = q^{N-1}$. Take an arbitrary $(N + 1)$ -subset S of the columns of M . We are going to show that there must exist a row of M that separates S . If $S \subseteq C_i$ for some $1 \leq i \leq N$, then the i -th row of C_i , which corresponds to π_0 , can separate S , since $\pi_0(\alpha) \neq \pi_0(\beta)$ for arbitrary distinct $\alpha, \beta \in V(\mathcal{A})$. Otherwise, let C_{i_1}, \dots, C_{i_j} be the column patterns which have non-empty intersection with S , where $j \geq 2$ is a positive integer. For $1 \leq l \leq j$, denote $C_{i_l} \cap S = S_l$. Then $\sum_{l=1}^j |S_l| = N + 1$ and $|S_l| \leq N$ for every l . By Lemma 4.3, at most $|S_l| - 1$ rows of C_{i_l} can not separate S_l . Since $\sum_{l=1}^j (|S_l| - 1) = N + 1 - j \leq N + 1 - 2 = N - 1 < N$, then there must exist a row of $(C_1|C_2|\cdots|C_N)$ that separates $\cup_{l=1}^j S_l = S$. \square

Remark 4.4. Our construction has an important property satisfying

$$\lim_{q \rightarrow \infty} \frac{Nq^{N-1}}{q^{N-1} + (N - 1)q^{N-2}} = N$$

and hence it is asymptotically optimal since we have $p_{N+1}(N, q) \leq Nq$ by Theorem 1.5. Note that a u -perfect hash family is $\{w_1, \dots, w_t\}$ -separating for arbitrary w_i such that $w_i \geq 1$ and $\sum_{i=1}^t w_i = u$. Theorems 1.5 and 1.7 can be combined to show

$$\lim_{q \rightarrow \infty} \frac{C(u - 1, q, \{w_1, \dots, w_t\})}{q} = u - 1,$$

which gives a negative answer to Question 1.6. Furthermore, taking into account the fact that any $(\lfloor (t/2 + 1)^2 \rfloor)$ -perfect hash family is also a t -IPP code, one can see that our construction also confirms the validity of Conjecture 1.8.

Remark 4.5. It is worth mentioning that Proposition 2 of [12] (an unpublished paper) also noticed that $\lim_{q \rightarrow \infty} \frac{p_u(u-1, q)}{q} = u - 1$. The author used an optimization method and no explicit construction was given in that paper.

The proof of Lemma 4.2 also leads to a conclusion on Hamming graphs, which we think may be of independent interest.

Corollary 4.6. Color the edges of $H(k, q)$ with k colors such that the edge (α, β) is colored by color i if α and β differ in their i -th coordinate. Then $H(k, q)$ contains no cycles with pairwise distinct colors.

5 Perfect hash families of strength three with three rows

Constructions for perfect hash families can induce constructions for corresponding separating hash families. And with the aid of Lemma 3.1, upper bounds for perfect hash families can also induce upper bounds for related separating hash families. Therefore, from this section we will focus on perfect hash families.

We have mentioned in Section 1 that if $(u-1) \nmid N$, it is very difficult to determine whether the exponent $\lceil N/(u-1) \rceil$ in Theorem 1.5 is tight. In the following two sections we will handle two small cases in such problems, namely, $N = u = 3$ and $N = u = 4$. When $N = 3$ and $u = 3$, the corresponding separating hash families only have two alternative types, namely, $\{1, 2\}$ -separating and 3-perfect hashing. Bazrafshan and Trung [8] proved that $C(3, q, \{1, 2\}) \leq q^2$ and an $SHF(3; q^2, q, \{1, 2\})$ does exist for $q \geq 2$. Walker and Colbourn [39] conjectured that $p_3(3, q) = o(q^2)$. In this section, we will verify this conjecture by proving $q^{2-o(1)} < p_3(3, q) = o(q^2)$. Furthermore, the upper bound is extended to $p_t(t, q)$ and $C(u, q, \{w_1, \dots, w_t\})$ with $\sum_{i=1}^t w_i = u$.

Let us begin with a simple lemma. Note that we will not distinguish between a perfect hash family and its representation matrix. We say a word x of the hash family (resp. a column of the representation matrix) has a unique coordinate i if for any other word (resp. column) y , $y \neq x$, it holds that $y(i) \neq x(i)$.

Lemma 5.1. *Let X denote the column set (words) of a $PHF(N; n, q, t)$. Then by deleting at most Nq words from X , we can get a subset $X^* \subseteq X$ such that no word in X^* has a unique coordinate in X^* .*

Proof. We use a greedy algorithm to construct X^* . Delete x_1 from X if x_1 has a unique coordinate in X . Denote $X_1 = X - \{x_1\}$. In general, if $x_{i+1} \in X_i$ has a unique coordinate in X_i , we delete x_{i+1} from X_i and then denote $X_{i+1} = X_i - \{x_{i+1}\}$. Continue this procedure until we get an X^* with no words containing a unique coordinate in it. At most Nq words will be deleted from X since we can delete any symbol $y \in [q]$ at most one time for any coordinate $i \in [N]$. \square

Since all perfect hash families being considered in the following are of size at least $q^{1+\epsilon}$ for some positive constant ϵ , then the deletion of at most Nq words from X can be neglected. Let $PHF^*(N; n, q, t)$ denote the perfect hash family (obtained from $PHF(N; n, q, t)$) such that no word in it contains a unique coordinate. We use $p_t^*(N, q)$ to denote the corresponding maximal cardinality.

Lemma 5.2. *In a $PHF^*(t; n, q, t)$, any two words can agree with at most one coordinate.*

Proof. Assume the contrary, then the following submatrix is contained in the representation matrix of such $PHF^*(t; n, q, t)$

$$\begin{pmatrix} \alpha_1(1) & \alpha_2(1) & * & * & * & * \\ \alpha_1(2) & \alpha_2(2) & * & * & * & * \\ \alpha_1(3) & * & \alpha_3(3) & * & * & * \\ \vdots & * & * & \ddots & * & * \\ \vdots & * & * & * & \ddots & * \\ \alpha_1(t) & * & * & * & * & \alpha_t(t) \end{pmatrix},$$

where in each row, the two bold coordinates are equal. α_1, α_2 are two words such that $\alpha_1(i) = \alpha_2(i)$ for $i = 1, 2$ and since α_1 has no unique coordinates, there exist $\alpha_3, \dots, \alpha_t$ such that $\alpha_j(j) = \alpha_1(j)$ for each $3 \leq j \leq t$. Therefore, no row of the submatrix can separate $\{\alpha_1, \dots, \alpha_t\}$, violating the t -perfect hashing property. \square

The following two observations are very useful.

Observation 1. On one hand, any $N \times n$ q -ary matrix M can be viewed as an N -uniform N -partite hypergraph $\mathcal{G} = (V(\mathcal{G}), E(\mathcal{G}))$ with equal part size q , where the vertex set is defined as $V(\mathcal{G}) = \cup_{i=1}^N V_i$, $V_i = \{(i, j) : 1 \leq j \leq q\}$ for $1 \leq i \leq N$, and the edge set is defined as $E(\mathcal{G}) = \{\{(i, x(i))\}_{i=1}^N : x = \{x(i)\}_{i=1}^N \text{ is a column of } M\}$.

Observation 2. On the other hand, given an N -uniform N -partite hypergraph $\mathcal{G} = (V(\mathcal{G}), E(\mathcal{G}))$ with equal part size q . We can regard $E(\mathcal{G})$ as some $N \times |E(\mathcal{G})|$ q -ary matrix M . Note that $V(\mathcal{G})$ can be partitioned into N pairwise disjoint sets with size q . We can set $V_i = \{(i, j) : 1 \leq j \leq q\}$ for $1 \leq i \leq N$, where the first coordinate i corresponds to the i -th part V_i and the second coordinate j corresponds to the j -th vertex in V_i . Then the matrix M is formed by setting its column set as $\{x = \{x(i)\}_{i=1}^N : \{(i, x(i))\}_{i=1}^N \in E(\mathcal{G})\}$. Such M is said to be the representation matrix of $E(\mathcal{G})$.

These two observations establish a bridge between q -ary matrices and multipartite hypergraphs. Recall the definition of $f_r^*(n, v, e)$ in Section 2.

Lemma 5.3. $p_3^*(3, q) \leq f_3^*(q, 6, 3) \leq p_3(3, q)$.

Proof. It is not hard to see that a $PHF^*(3; n, q, 3)$ exists if and only if the following configuration is not contained in its representation matrix

$$\begin{pmatrix} a & * & a \\ b & b & * \\ * & c & c \end{pmatrix},$$

where none of the stars belong to $\{a, b, c\}$.

We call this configuration a triangle since these three columns have no identical coordinates and every pair of columns have exactly one common coordinate. On one hand, it holds that $f_3^*(q, 6, 3) \leq p_3(3, q)$, since for arbitrary three columns of a hash family, if no row can separate them then for each row there exists some coordinate equal to another one. Therefore, these columns (or corresponding edges) must be spanned by at most six points, which violates the (6,3)-free property. On the other hand, if some three columns of a $PHF^*(3; n, q, 3)$ contain at most six points, then either there exists a pair of two columns having two coordinates in common or these three columns form a triangle. Both cases are forbidden in a $PHF^*(3; n, q, 3)$. Therefore, it holds that $p_3^*(3, q) \leq f_3^*(q, 6, 3)$ and hence our lemma follows. \square

Theorem 5.4. $p_3(3, q) = f_3(3q, 6, 3) + O(q)$ and hence for arbitrary $\epsilon > 0$, $q^{2-\epsilon} < p_3(3, q) = o(q^2)$ holds for sufficiently large q .

Proof. Apply Lemmas 2.2, 5.1, 5.3 and inequality (1). \square

As the second application of Lemma 3.1, the upper bound of $p_3(3, q)$ can be extended to $p_t(t, q)$ and $C(u, q, \{w_1, \dots, w_t\})$.

Corollary 5.5. $C(u, q, \{w_1, \dots, w_t\}) = o(q^2)$ for any $t \geq 3$ and $\sum_{i=1}^t w_i = u$. In particular, $p_t(t, q) = o(q^2)$ for any $t \geq 3$.

Proof. Apply Lemma 3.1 and Theorem 5.4. \square

Remark 5.6. One can also prove $p_t(t, q) = o(q^2)$ by applying the graph removal lemma [2], see [3, 6] for examples of applications of graph removal lemma in such problems. Here our proof applying Lemma 3.1 is much simpler. When $1+w \leq q$, it was shown in [8] that $C(1+w, q, \{1, w\}) \leq q^2$. And for any prime power q , there exists an $SHF(w+1; q^2, q, \{1, w\})$. Therefore, for $C(u, q, \{w_1, \dots, w_t\})$ with $t = 2$, we can not determine whether $C(w_1 + w_2, q, \{w_1, w_2\}) = \Omega(q^2)$ or $C(w_1 + w_2, q, \{w_1, w_2\}) = o(q^2)$. It is an interesting problem to determine the right order of the magnitude of $C(w_1 + w_2, q, \{w_1, w_2\})$.

Although we can get the lower bound $q^{2-\epsilon}$ by a direct application of the (6,3)-theorem and Lemma 2.2, we prefer a construction which provides the explicit cardinality. A method introduced in Section 3 of [25] can be used to construct such q -ary codes of length N . Our method is similar to that one except some transformations which will be mentioned later.

Given integers $q \geq N \geq 2$, $M \subseteq \{0, 1, \dots, q-1\}$, we define an N -uniform N -partite hypergraph \mathcal{G}_M (whose edge set can be viewed as the representation matrix of our desired code) as follows. The vertex set $V(\mathcal{G}_M)$ is defined to be

$$V(\mathcal{G}_M) := \{(j, y) : j \in [N], y \in \mathbb{Z}_q\}.$$

It is easy to see $|V(\mathcal{G}_M)| = Nq$. For each $j \in [N]$, we use $V_j = \{(j, y) : y \in \mathbb{Z}_q\}$ to denote the vertex set of the j -th part of $V(\mathcal{G})$. For integers $0 \leq y, m \leq q$, the hyperedge of \mathcal{G} is defined to be the N -element set

$$A(y, m) = \{(1, y + b_1 m), (2, y + b_2 m), \dots, (N, y + b_N m)\},$$

where $\mathcal{B} := \{b_1, \dots, b_N\} \subseteq \{0, 1, \dots, q-1\}$ is an undetermined N -element set and the second coordinates $y + b_i m$ are taken modulo q . We call \mathcal{B} the tangent set of $A(y, m)$. $A(y, m)$ can also be viewed as a q -ary word of length N . If q is a prime, one can verify that

$$|A(y, m) \cap A(y', m')| \leq 1 \tag{4}$$

holds for $(y, m) \neq (y', m')$ by solving a system of two congruence equations.

From now on, we fix the size of the alphabet set q to be a prime or the prime nearest to it. For a subset $M \subseteq \{0, 1, \dots, q-1\}$, we set

$$E(\mathcal{G}_M) := \{A(y, m) : y \in \mathbb{Z}_q, m \in M\}$$

to be the edge set of our desired hypergraph, where the set M is determined by the subgraphs that needed to be forbidden (these subgraphs can also be viewed as the configurations that needed to be forbidden in the desired code). Obviously $|E(\mathcal{G}_M)| = q|M|$ and by (4) we can verify that \mathcal{G}_M is also linear.

Now, we are going to choose appropriate \mathcal{B} and M according to the properties of our desired codes. For example, to construct a 3-perfect hash family with three rows, we first set $N = 3$ and then choose $\mathcal{B} = \{0, 1, 2\}$, where $b_i = i - 1$ for $1 \leq i \leq 3$. Therefore, to show this specified \mathcal{G}_M can indeed induce a $PHF(3, n, q, 3)$, by Lemma 5.3 we only need to guarantee that $E(\mathcal{G})$ is triangle-free, since it is already linear (we have set q to be a prime). We claim that it suffices to choose $M \subseteq \{0, 1, \dots, \lfloor (q-1)/2 \rfloor\}$ to be a 2-sum-free set such that the equation $m_1 + m_2 = 2m_3$ has no solution except $m_1 = m_2 = m_3$.

Theorem 5.7. *There exists a constant γ such that $p_3(3, q) > q^2 e^{-\gamma \sqrt{\log q}}$.*

Proof. It suffices to show \mathcal{G}_M contains no triangles for arbitrary 2-sum-free set $M \subseteq \{0, 1, \dots, \lfloor (q-1)/2 \rfloor\}$. If otherwise, assume that $\{A(y_i, m_i) \in \mathcal{G}_M : 1 \leq i \leq 3\}$ forms a triangle. One can verify that the vertices of this triangle must locate on different parts of V_1, V_2, V_3 . Thus we can assume that

$$\begin{cases} A(y_1, m_1) \cap A(y_2, m_2) = \{(j_2, a_2)\} \\ A(y_2, m_2) \cap A(y_3, m_3) = \{(j_3, a_3)\} \\ A(y_3, m_3) \cap A(y_1, m_1) = \{(j_1, a_1)\} \end{cases}$$

where $\{j_1, j_2, j_3\} = \{1, 2, 3\}$ and a_1, a_2, a_3 are some positive integers. Then the following three equations hold simultaneously

$$\begin{cases} y_1 + (j_2 - 1)m_1 \equiv y_2 + (j_2 - 1)m_2 \pmod{q} \\ y_2 + (j_3 - 1)m_2 \equiv y_3 + (j_3 - 1)m_3 \pmod{q} \\ y_3 + (j_1 - 1)m_3 \equiv y_1 + (j_1 - 1)m_1 \pmod{q}. \end{cases}$$

Because of the symmetry of a triangle, we can always assume that $j_1 < j_2 < j_3$. By a simple elimination we can infer

$$(j_2 - j_1)m_1 + (j_3 - j_2)m_2 \equiv (j_3 - j_1)m_3 \pmod{q},$$

or simply

$$m_1 + m_2 \equiv 2m_3 \pmod{q}.$$

This implies $m_1 + m_2 = 2m_3$ since $m_i \leq \lfloor (q-1)/2 \rfloor$ for all $1 \leq i \leq 3$, which contradicts the fact that M is 2-sum-free. By Lemma 2.1 there exists a 2-sum-free set M with $|M| > q e^{-\gamma \sqrt{\log q}}$ for some constant γ . Therefore, it follows that $|E(\mathcal{G}_M)| = q|M| > |M| > q^2 e^{-\gamma \sqrt{\log q}}$. \square

6 Perfect hash families of strength four with four rows

It is much more complicated to construct 4-perfect hash families such that $p_4(4, q) > q^{2-o(1)}$. We will use the notion of rainbow cycles and R -sum-free sets defined in Section 2. In fact, we are going to prove the following result:

Lemma 6.1. $p_t^*(t, q) \leq g_t^*(q) \leq p_t(t, q)$.

Proof. First we are going to show that any $PHF^*(t; n, q, t)$ can induce a t -uniform t -partite linear hypergraph \mathcal{G} containing no rainbow cycles. Let M denote the representation matrix of the hash family, then M can also be viewed as the representation matrix of $E(\mathcal{G})$ by Observation 1. M (resp. $E(\mathcal{G})$) is already linear by Lemma 5.2. It suffices to show M (resp. $E(\mathcal{G})$) contains no rainbow cycles. Assume otherwise, the columns (resp. hyperedges) of M (resp. $E(\mathcal{G})$) indexed by $\alpha_1, \dots, \alpha_k$ form a rainbow k -cycle $v_1, \alpha_1, v_2, \alpha_2, \dots, v_k, \alpha_k, v_1$ with $k \leq t$. By Observation 1, the i -th part of $V(\mathcal{G})$ can be defined as $V_i = \{(i, j) : j \in [q]\}$, where the first coordinate corresponds to the i -th row of M and the second coordinate corresponds to the j -th element in $[q]$. Without loss of generality, we can assume that v_i is from the i -th part of the vertex set. Then it holds that $\alpha_i(i) = \alpha_{i+1}(i)$ for $1 \leq i \leq k-1$ and $\alpha_k(k) = \alpha_1(k)$. The following submatrix induced by such k -cycle is contained in M :

$$\begin{pmatrix} \alpha_1(\mathbf{1}) & \alpha_2(\mathbf{1}) & \alpha_3(\mathbf{1}) & \alpha_4(\mathbf{1}) & & \alpha_{k-1}(\mathbf{1}) & \alpha_k(\mathbf{1}) \\ \alpha_1(\mathbf{2}) & \alpha_2(\mathbf{2}) & \alpha_3(\mathbf{2}) & \alpha_4(\mathbf{2}) & & \alpha_{k-1}(\mathbf{2}) & \alpha_k(\mathbf{2}) \\ \alpha_1(\mathbf{3}) & & \alpha_3(\mathbf{3}) & \alpha_4(\mathbf{3}) & & \alpha_{k-1}(\mathbf{3}) & \alpha_k(\mathbf{3}) \\ \vdots & & & \ddots & & \vdots & \vdots \\ \vdots & & & & \ddots & \alpha_{k-1}(\mathbf{k}-1) & \alpha_k(\mathbf{k}-1) \\ \alpha_1(\mathbf{k}) & & & & & & \alpha_k(\mathbf{k}) \end{pmatrix},$$

where in each row, the two bold coordinates are equal. Note that in this matrix, the columns represent the hyperedges and the coordinates in each column represent the vertices contained in the corresponding hyperedge. It is easy to see none of the first k rows of M can separate $\{\alpha_1, \dots, \alpha_k\}$. Note that no column of M has unique coordinates, then there exist $\alpha_{k+1}, \dots, \alpha_t$ such that $\alpha_j(j) = \alpha_1(j)$ for $k+1 \leq j \leq t$, which can also be depicted by

$$\begin{pmatrix} \alpha_1(\mathbf{k}+1) & \alpha_{k+1}(\mathbf{k}+1) & * & * & * & * \\ \alpha_1(\mathbf{k}+2) & * & \alpha_{k+2}(\mathbf{k}+2) & * & * & * \\ \vdots & * & * & \ddots & * & * \\ \vdots & * & * & * & \ddots & * \\ \alpha_1(\mathbf{t}) & * & * & * & * & \alpha_t(\mathbf{t}) \end{pmatrix}.$$

Therefore, the left $t-k$ rows of M can not separate $\{\alpha_1, \alpha_{k+1}, \dots, \alpha_t\}$. So we can conclude that no row of M can separate $\{\alpha_1, \dots, \alpha_t\}$, violating the t -perfect hashing property.

It remains to show that any t -uniform t -partite linear hypergraph (with equal part size q) \mathcal{G} without rainbow cycles can induce a $PHF(t; n, q, t)$ such that $n = |E(\mathcal{G})|$. We also use M to denote the representation matrix of $E(\mathcal{G})$. We claim that if there exists a $t \times t$ submatrix T of M such that no row can separate it, then the hypergraph induced by T will contain a rainbow k -cycle with $k \leq t$.

We will argue by induction on t . When $t = 2$, a 2×2 submatrix can always be separated by one of its two rows provided that the two columns of this submatrix are distinct. When $t = 3$, if a 3×3 submatrix of a 3-uniform 3-partite linear hypergraph can not be separated by one of its three rows, then this submatrix actually forms a triangle defined in Lemma 5.3. One can verify that this triangle can be represented as a rainbow 3-cycle $\{a, E_1, b, E_2, c, E_3\}$ for some edges E_1, E_2, E_3 . Now assume the statement is true for $t-1$. Take a $t \times t$ matrix T with columns indexed by $C = \{\alpha_1, \dots, \alpha_t\}$ and rows indexed by $R = \{r_1, \dots, r_t\}$ such that no row can separate C . We denote $C_i = C - \{\alpha_i\}$ and $R_i = R - \{r_i\}$ for each $1 \leq i \leq t$. Furthermore, we use T_{ij} to denote the $(t-1) \times (t-1)$ submatrix formed by R_i and C_j . Then for any submatrix T_{ij} , there must exist a row that separates all columns of T_{ij} since otherwise T_{ij} contains a rainbow k -cycle with $k \leq t-1$ by the induction hypothesis.

Without loss of generality, assume r_1 separates C_t . Note that this row can not separate C , so we can assume further that $\alpha_t(1) = \alpha_1(1)$. Then consider T_{11} , there exists a row in $R - \{r_1\}$ that separates $C - \{\alpha_1\}$. We can set this row to be r_2 . Similarly, there exists $2 \leq j \leq t$ such that $\alpha_1(2) = \alpha_j(2)$ since r_2 can not separate C . Then $j \neq t$ since α_1 and α_t have already agreed on one coordinate, say, $\alpha_t(1) = \alpha_1(1)$. Assume that $\alpha_1(2) = \alpha_2(2)$. Now consider T_{22} , then there exists a row in $R - \{r_2\}$ that separates $C - \{\alpha_2\}$. Note that this row can not be r_1 since α_1 and α_t agree on their first coordinate. We can set this row to be r_3 . For the same reason, there exists $j \in [t], j \neq 2$ such that $\alpha_2(3) = \alpha_j(3)$. Then $j \neq 1$ since it already holds $\alpha_1(2) = \alpha_2(2)$. If $j = t$, we are done since $\{\alpha_1, \alpha_2, \alpha_t\}$ forms a rainbow 3-cycle. So we can set $j = 3$.

The above discussion can be depicted by the following matrix:

$$\left(\begin{array}{cccccccc} \alpha_1(\mathbf{1}) & \alpha_2(1) & \alpha_3(1) & \alpha_4(1) & \cdots & \cdots & \alpha_{t-1}(1) & \alpha_t(\mathbf{1}) \\ \alpha_1(\mathbf{2}) & \alpha_2(\mathbf{2}) & \alpha_3(2) & \alpha_4(2) & \cdots & \cdots & \alpha_{t-1}(2) & \alpha_t(2) \\ \alpha_1(3) & \alpha_2(\mathbf{3}) & \alpha_3(\mathbf{3}) & \alpha_4(3) & \cdots & \cdots & \alpha_{t-1}(3) & \alpha_t(3) \\ \alpha_1(4) & \alpha_2(4) & \alpha_3(\mathbf{4}) & \alpha_4(\mathbf{4}) & \cdots & \cdots & \alpha_{t-1}(4) & \alpha_t(4) \\ & & & \ddots & & & & \\ & & & & \ddots & & & \\ \alpha_1(i-1) & \cdots & \alpha_{i-2}(i-1) & \alpha_{i-1}(i-1) & & \cdots & \cdots & \alpha_t(i+1) \\ \alpha_1(i) & \cdots & & \alpha_{i-1}(i) & \alpha_i(i) & & \cdots & \alpha_t(i+1) \\ \alpha_1(i+1) & \cdots & & & \alpha_i(i+1) & \alpha_{i+1}(i+1) & \cdots & \alpha_t(i+1) \\ & & & \ddots & & & & \\ & & & & \ddots & & & \end{array} \right),$$

where in each row, the two bold coordinates are equal. We continue this procedure for $T_{i,i}$ with $i \geq 3$. By our choice, for all $1 \leq j \leq i$, in row r_j it holds that $\alpha_{j-1}(j) = \alpha_j(j)$ (α_0 is recognised as α_t). Thus no row in $\{r_1, \dots, r_i\}$ can separate $T_{i,i}$. We can always assume that $r_{i+1} \in R - \{r_i\}$ is the row that separates $C - \{\alpha_i\}$. Then there exists a $j \in [t], j \neq i$ such that $\alpha_i(i+1) = \alpha_j(i+1)$ since r_{i+1} can not separate the whole C . Obviously, $j \neq i-1$. If $j \in \{1, \dots, i-2\}$ or $j = t$, then such choice of j will induce a rainbow $(i-j+1)$ -cycle formed by $\{\alpha_j, \dots, \alpha_i\}$

$$\left(\begin{array}{cccccc} \alpha_j(j+1) & \alpha_{j+1}(j+1) & * & * & * \\ * & \alpha_{j+1}(j+2) & \alpha_{j+2}(j+2) & * & * \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ * & * & * & \alpha_{i-1}(i) & \alpha_i(i) \\ \alpha_j(i+1) & * & * & * & \alpha_i(i+1) \end{array} \right)$$

or a rainbow $(i+1)$ -cycle formed by $\{\alpha_1, \dots, \alpha_i, \alpha_t\}$

$$\left(\begin{array}{cccccc} \alpha_1(\mathbf{1}) & * & * & & \cdots & \alpha_t(\mathbf{1}) \\ \alpha_1(\mathbf{2}) & \alpha_2(\mathbf{2}) & * & & \cdots & * \\ * & \alpha_2(\mathbf{3}) & \alpha_3(\mathbf{3}) & & \cdots & * \\ \vdots & \vdots & & \ddots & \cdots & \vdots \\ \vdots & \vdots & \vdots & \ddots & \cdots & \vdots \\ * & * & * & \alpha_{i-1}(i) & \alpha_i(i) & * \\ * & * & * & \alpha_i(i+1) & \cdots & \alpha_t(i+1) \end{array} \right).$$

If neither one of the above cases holds, we can always assume that $j = i+1$ and continue this procedure.

This procedure will end when it comes to $T_{t-1,t-1}$ with $\alpha_{t-1}(t) = \alpha_t(t)$. Then $\{\alpha_1, \dots, \alpha_t\}$ will form a rainbow t -cycle and our desired contradiction follows. \square

We can use a similar method as that of the previous section to construct 4-perfect hash family with four rows. However, we can not simply take $\mathcal{B} = \{0, 1, 2, 3\}$ since such choice will lead to an equation

$$2m_1 + 2m_2 - 3m_3 - m_4 = 0,$$

whose solution is not easy to determine as suggested by Ruzsa [33]. In order to show $p_4(4, q) > q^{2-o(1)}$, we should choose \mathcal{B} more carefully. Recall that we have set q to be a prime.

Lemma 6.2. *Let $R = \{b_1, \dots, b_r\} \subseteq \{0, \dots, q-1\}$ be an r -element subset with rank $r(R)$. If $M \subseteq \{0, 1, \dots, \lfloor (q-1)/r(R) \rfloor\}$ is an R -sum-free set, then the hypergraph defined by*

$$E(\mathcal{G}_M) = \{A(y, m) : y \in \mathbb{Z}_q, m \in M\},$$

where $A(y, m) = \{(i, y + b_i m) : b_i \in R\}$, is an r -uniform r -partite linear hypergraph containing no rainbow cycles.

Proof. First it is easy to see \mathcal{G}_M is r -uniform and r -partite with $V(\mathcal{G}_M) = \cup_{j=1}^r V_j$, where $V_j = \{(j, y) : y \in \mathbb{Z}_q\}$, $1 \leq j \leq r$. To see \mathcal{G}_M is also linear, one just needs to notice that if $|A(y, m) \cap A(y', m')| \geq 2$, then there are $b_1, b_2 \in R$, $b_1 \neq b_2$ such that

$$\begin{cases} y + b_1 m \equiv y' + b_1 m' \pmod{q} \\ y + b_2 m \equiv y' + b_2 m' \pmod{q}. \end{cases}$$

Then we can infer $(b_1 - b_2)(m - m') \equiv 0 \pmod{q}$, which is a contradiction with q prime.

Now it remains to show that \mathcal{G}_M indeed contains no rainbow cycles. Assume the contrary, it contains a rainbow k -cycle with $k \leq r$, denoted by $v_1, A(y_1, m_1), v_2, A(y_2, m_2), \dots, v_k, A(y_k, m_k), v_1$, where $v_i \in V_{j_i}$ and $j_{i_1} \neq j_{i_2}$ for $i_1 \neq i_2$ by the definition of a rainbow cycle. The following k equations hold simultaneously:

$$\begin{cases} y_1 + b_{j_2} m_1 \equiv y_2 + b_{j_2} m_2 \pmod{q} \\ y_2 + b_{j_3} m_2 \equiv y_3 + b_{j_3} m_3 \pmod{q} \\ \vdots \\ y_{k-1} + b_{j_k} m_{k-1} \equiv y_k + b_{j_k} m_k \pmod{q} \\ y_k + b_{j_1} m_k \equiv y_1 + b_{j_1} m_1 \pmod{q}. \end{cases}$$

By a simple elimination, one can infer

$$(b_{j_2} - b_{j_1})m_1 + (b_{j_3} - b_{j_2})m_2 + \dots + (b_{j_k} - b_{j_{k-1}})m_{k-1} + (b_{j_1} - b_{j_k})m_k \equiv 0 \pmod{q},$$

or

$$(b_{j_2} - b_{j_1})m_1 + (b_{j_3} - b_{j_2})m_2 + \dots + (b_{j_k} - b_{j_{k-1}})m_{k-1} + (b_{j_1} - b_{j_k})m_k = 0,$$

since $m_i \leq \lfloor (q-1)/r(R) \rfloor$ for each $1 \leq i \leq k$, which implies $m_1 = \dots = m_k$ taking into account the fact that M is R -sum-free. Thus $y_1 = \dots = y_k$, which is a contradiction. Therefore, we can conclude that \mathcal{G}_M contains no rainbow cycles. \square

Lemmas 6.1 and 6.2 suggest that we can use tools from additive number theory to construct good perfect hash families. As discussed before Theorem 5.7, we use \mathcal{B} to denote the set of tangents of $A(y, m)$. To construct $PHF(4; n, q, 4)$, we take $\mathcal{B} = \{0, 2, 5, \mu + 5\}$, where $b_0 = 0, b_1 = 2, b_2 = 5$ and $b_3 = \mu + 5$ with $\mu = \lceil 2^{\sqrt{\log q}} \rceil$. Note that $\mu = o(q^\epsilon)$ for arbitrary small constant $\epsilon > 0$. By previous lemmas, our goal is to construct a \mathcal{B} -sum-free subset M of \mathbb{Z}_q with sufficiently large cardinality. The desired hyperedge $A(y, m)$ is defined to be

$$A(y, m) = \{(1, y), (2, y + 2m), (3, y + 5m), (4, y + (\mu + 5)m)\}. \quad (5)$$

The following lemma (together with Lemma 6.2) shows that if we choose M as the set defined in Lemma 2.4, then the corresponding $E(\mathcal{G}_M)$ contains no rainbow cycles.

Lemma 6.3. *Choose M as the set defined in Lemma 2.4 and let \mathcal{B} be the 4-element set defined above, then M is \mathcal{B} -sum-free.*

Proof. Note that M has no nontrivial solution to all equations in (3), then one can verify this lemma directly by definition. \square

Lemma 6.4. *The hypergraph defined by*

$$\mathcal{G}_M = \{A(y, m) : y \in \mathbb{Z}_q, m \in M\},$$

has no rainbow cycles, where M is the set defined in Lemma 2.4 and $A(y, m)$ is defined in (5).

Proof. Apply Lemmas 6.2 and 6.3 and note that $r(\mathcal{B}) = \mu + 5$. \square

Theorem 6.5. *There exists a constant γ such that $p_4(4, q) > q^2 e^{-\gamma(\log q)^{3/4}}$.*

Proof. Apply Lemmas 2.4, 6.1 and 6.4. Then the theorem follows from

$$p_4(4, q) \geq g_4^*(q) \geq |E(\mathcal{G}_M)| = q|M| > q^2 e^{-\gamma(\log q)^{3/4}}.$$

\square

Remark 6.6. *In the above construction of $PHF^*(4, n, q, 4)$, we choose the tangent set \mathcal{B} of the hyperedge $A(y, m)$ to be $\mathcal{B} = \{0, 2, 5, \mu + 5\}$ with $\mu = \lceil 2\sqrt{\log q} \rceil$. This choice of \mathcal{B} has appeared in [3], where the authors used such \mathcal{B} to construct 2-IPP codes. In this paper we choose the same \mathcal{B} as they did since in this way we can save the space for proving Lemma 2.4. Actually, when $|R| = 4$ there are many choices of \mathcal{B} satisfying the following conditions*

- (a) $M \subseteq \{0, 1, \dots, \lfloor (q-1)/r(R) \rfloor\}$ is R -sum-free,
- (b) $|M| > q^{1-o(1)}$,
- (c) $r(R) = o(q^\epsilon)$ for arbitrary small $\epsilon > 0$.

However, for $|R| \geq 5$, we do not know whether such \mathcal{B} exists.

7 Connections to hypergraph Turán problems

In this section we will study perfect hash families in view of hypergraph Turán problems.

Theorem 7.1. *For arbitrary positive integers t, N, q , it holds that $f_N^*(q, tN - N, t) \leq p_t(N, q)$. Furthermore, $\frac{N!}{N^N} f_N(Nq, tN - N, t) \leq p_t(N, q)$.*

Proof. By Lemma 2.2, it suffices to prove the first statement of the theorem. Recall that if a hypergraph \mathcal{G} is N -uniform N -partite with equal part size q , then $E(\mathcal{G})$ can be represented by an $N \times |E(\mathcal{G})|$ q -ary matrix M . If \mathcal{G} is $G(tN - N, t)$ -free, then given any collection of t edges $S \subseteq E(\mathcal{G})$, it is not hard to verify that in its representation matrix there must exist a row that separates S , since otherwise S can contain at most $tN - N$ vertices, violating the fact that \mathcal{G} is $G(tN - N, t)$ -free. Therefore, M can be viewed as the representation matrix of the desired perfect hash family. \square

A direct application of Theorem 7.1 gives the following result.

Corollary 7.2. *If $2 \nmid N$, then for arbitrary $\epsilon > 0$, it holds that $p_3(N, q) > q^{\lceil N/2 \rceil - \epsilon}$.*

Proof. This corollary follows from the inequality (2), $n^{k-o(1)} < f_r(n, 3(r-k) + k + 1, 3) = o(n^k)$. Set $N = 2k - 1$ and $t = 3$, by Theorem 7.1 one can infer

$$p_3(N, q) \geq p_3^*(N, q) \geq f_N^*(q, 3N - N, 3) \geq \frac{N!}{N^N} f_N(Nq, 3N - N, 3) > \frac{N!}{N^N} (Nq)^{\lceil N/2 \rceil - o(1)}.$$

\square

8 Concluding remarks

In this paper we mainly study codes and hash families with the separating property. Several open problems and conjectures concerning the upper or lower bounds are solved. Our two essential methods to study these objects can be summarized as follows.

The first method is to discover the structural information hidden in the separating property. As an example, our Johnson-type bound (Lemma 3.1) is used to establish Theorem 1.5 and Corollary 5.5.

The second one is that we establish a bridge between perfect hash families, graph theory and additive number theory. For example, we solve Conjecture 1.9 by considering a related hypergraph Turán problem. We also showed that tools from additive number theory can be used to construct good perfect hash families. As a result, Theorems 5.7, 6.5 and Corollary 7.2 suggest that there may exist a positive answer to Question 1.10.

Besides these two new methods, we believe that the construction in Section 4 is of interest since it generalizes many previous ones. Further generalizations of our method are expected.

As a conclusion, we would like to mention several open problems which we think are interesting.

Open Problem 1. If $2 \nmid N$, Corollary 7.2 shows that $p_3(N, q) > q^{\lceil N/2 \rceil - o(1)}$. Determine whether $p_3(N, q) = o(q^{\lceil N/2 \rceil})$ or $p_3(N, q) = \Theta(q^{\lceil N/2 \rceil})$.

Open Problem 2. For r -uniform r -partite linear hypergraph without rainbow cycles, we have proved that $g_r^*(q) = o(q^2)$ and $g_i^*(q) > q^{2-o(1)}$ for $i = 3, 4$. Then does it hold that $g_r^*(q) > q^{2-o(1)}$ for all $r \geq 3$?

Open Problem 3. For arbitrary $r \geq 3$, does there exist an r -element set R and $M \subseteq [q]$ such that the conditions in Remark 6.6 are satisfied? Note that the question is true when $r = 3, 4$.

Open Problem 4. It has been shown in Theorem 7.1 that $p_t(N, q) \geq f_N^*(q, tN - N, t)$. Then does there exist an upper bound for $p_t(N, q)$ only using $f_N^*(q, v, t)$?

References

- [1] N. Alon, G. Cohen, M. Krivelevich, and S. Litsyn. Generalized hashing and parent-identifying codes. *J. Combin. Theory Ser. A*, 104(1):207–215, 2003.
- [2] N. Alon, R.A. Duke, H. Lefmann, V. Rödl, and R. Yuster. The algorithmic aspects of the regularity lemma. *J. Algorithms*, 16(1):80–109, 1994.
- [3] N. Alon, E. Fischer, and M. Szegedy. Parent-identifying codes. *J. Combin. Theory Ser. A*, 95(2):349–359, 2001.
- [4] N. Alon and M. Naor. Derandomization, witnesses for Boolean matrix multiplication and construction of perfect hash functions. *Algorithmica*, 16(4-5):434–449, 1996.
- [5] N. Alon and A. Shapira. On an extremal hypergraph problem of Brown, Erdős and Sós. *Combinatorica*, 26(6):627–645, 2006.
- [6] N. Alon and U. Stav. New bounds on parent-identifying codes: the case of multiple parents. *Combin. Probab. Comput.*, 13(6):795–807, 2004.
- [7] M. Bazrafshan and T. Trung. Bounds for separating hash families. *J. Combin. Theory Ser. A*, 118(3):1129–1135, 2011.
- [8] M. Bazrafshan and T. Trung. Improved bounds for separating hash families. *Des. Codes Cryptogr.*, 69(3):369–382, 2013.
- [9] F.A. Behrend. On sets of integers which contain no three terms in arithmetical progression. *Proc. Nat. Acad. Sci. U. S. A.*, 32:331–332, 1946.

- [10] C. Berge. Hypergraphs. In *Selected topics in graph theory*, 3, pages 189–206. Academic Press, San Diego, CA, 1988.
- [11] C. Berge. *Hypergraphs*, volume 45 of *North-Holland Mathematical Library*. North-Holland Publishing Co., Amsterdam, 1989. Combinatorics of finite sets, Translated from the French.
- [12] S.R. Blackburn. Perfect hash families with few functions. *Unpublished manuscript, 2000; available online as IACR research report 2003/17; see <http://eprint.iacr.org/2003/017>*.
- [13] S.R. Blackburn. Perfect hash families: probabilistic methods and explicit constructions. *J. Combin. Theory Ser. A*, 92(1):54–60, 2000.
- [14] S.R. Blackburn. Combinatorial schemes for protecting digital content. In *Surveys in combinatorics, 2003 (Bangor)*, volume 307 of *London Math. Soc. Lecture Note Ser.*, pages 43–78. Cambridge Univ. Press, Cambridge, 2003.
- [15] S.R. Blackburn. Frameproof codes. *SIAM J. Discrete Math.*, 16(3):499–510 (electronic), 2003.
- [16] S.R. Blackburn. An upper bound on the size of a code with the k -identifiable parent property. *J. Combin. Theory Ser. A*, 102(1):179–185, 2003.
- [17] S.R. Blackburn, M. Burmester, Y. Desmedt, and P.R. Wild. Efficient multiplicative sharing schemes. In *Advances in cryptology—EUROCRYPT '96 (Saragossa, 1996)*, volume 1070 of *Lecture Notes in Comput. Sci.*, pages 107–118. Springer, Berlin, 1996.
- [18] S.R. Blackburn, T. Etzion, D.R. Stinson, and G.M. Zaverucha. A bound on the size of separating hash families. *J. Combin. Theory Ser. A*, 115(7):1246–1256, 2008.
- [19] D. Boneh and J. Shaw. Collusion-secure fingerprinting for digital data. *IEEE Trans. Inform. Theory*, 44(5):1897–1905, 1998.
- [20] W.G. Brown, P. Erdős, and V.T. Sós. On the existence of triangulated spheres in 3-graphs, and related problems. *Period. Math. Hungar.*, 3(3-4):221–228, 1973.
- [21] W.G. Brown, P. Erdős, and V.T. Sós. Some extremal problems on r -graphs. In *New directions in the theory of graphs (Proc. Third Ann Arbor Conf., Univ. Michigan, Ann Arbor, Mich, 1971)*, pages 53–63. Academic Press, New York, 1973.
- [22] P. Erdős, P. Frankl, and V. Rödl. The asymptotic number of graphs not containing a fixed subgraph and a problem for hypergraphs having no exponent. *Graphs Combin.*, 2(2):113–121, 1986.
- [23] P. Erdős and D.J. Kleitman. On coloring graphs to maximize the proportion of multicolored k -edges. *J. Combinatorial Theory*, 5:164–169, 1968.
- [24] R. Fuji-Hara. Perfect hash families of strength three with three rows from varieties on finite projective geometries. *Des., Codes Cryptogr.*, to appear. DOI: 10.1007/s10623-015-0052-z.
- [25] Z. Füredi and M. Ruszinkó. Uniform hypergraphs containing no grids. *Adv. Math.*, 240:302–324, 2013.
- [26] C.D. Godsil. *Algebraic combinatorics*. Chapman and Hall Mathematics Series. Chapman & Hall, New York, 1993.
- [27] D.L. Hollmann, J.H. van Lint, J.-P. Linnartz, and L.M.G.M. Tolhuizen. On codes with the identifiable parent property. *J. Combin. Theory Ser. A*, 82(2):121–133, 1998.
- [28] J. Körner and K. Marton. New bounds for perfect hashing via information theory. *European J. Combin.*, 9(6):523–530, 1988.
- [29] S. Martirosyan and T. Trung. Explicit constructions for perfect hash families. *Des. Codes Cryptogr.*, 46(1):97–112, 2008.
- [30] K. Mehlhorn. *Data structures and algorithms. 1*. EATCS Monographs on Theoretical Computer Science. Springer-Verlag, Berlin, 1984. Sorting and searching.

- [31] I. Newman and A. Wigderson. Lower bounds on formula size of Boolean functions using hypergraph entropy. *SIAM J. Discrete Math.*, 8(4):536–542, 1995.
- [32] A. Nilli. Perfect hashing and probability. *Combin. Probab. Comput.*, 3(3):407–409, 1994.
- [33] I.Z. Ruzsa. Solving a linear equation in a set of integers. I. *Acta Arith.*, 65(3):259–282, 1993.
- [34] I.Z. Ruzsa and E. Szemerédi. Triple systems with no six points carrying three triangles. In *Combinatorics (Proc. Fifth Hungarian Colloq., Keszthely, 1976), Vol. II*, volume 18 of *Colloq. Math. Soc. János Bolyai*, pages 939–945. North-Holland, Amsterdam-New York, 1978.
- [35] J.N. Staddon, D.R. Stinson, and R. Wei. Combinatorial properties of frameproof and traceability codes. *IEEE Trans. Inform. Theory*, 47(3):1042–1049, 2001.
- [36] D.R. Stinson, T. Trung, and R. Wei. Secure frameproof codes, key distribution patterns, group testing algorithms and related structures. *J. Statist. Plann. Inference*, 86(2):595–617, 2000. Special issue in honor of Professor Ralph Stanton.
- [37] D.R. Stinson and R. Wei. Combinatorial properties and constructions of traceability schemes and frameproof codes. *SIAM J. Discrete Math.*, 11(1):41–53 (electronic), 1998.
- [38] D.R. Stinson, R. Wei, and K. Chen. On generalized separating hash families. *J. Combin. Theory Ser. A*, 115(1):105–120, 2008.
- [39] R.A. Walker II and C.J. Colbourn. Perfect Hash families: constructions and existence. *J. Math. Cryptol.*, 1(2):125–150, 2007.