# Subexponential and Linear Subpacketization Coded Caching via Projective Geometry

Hari Hara Suthan Chittoor, Prasad Krishnan, K V Sushena Sree, and Bhavana MVN

## Abstract

Large gains in the rate of cache-aided broadcast communication are obtained using coded caching, but to obtain this most existing centralized coded caching schemes require that the files at the server be divisible into a large number of parts (this number is called subpacketization). In fact, most schemes require the subpacketization to be growing asymptotically as exponential in $\sqrt[r]{K}$ for some positive integer $r$ and $K$ being the number of users. On the other extreme, few schemes having subpacketization linear in $K$ are known; however, they require large number of users to exist, or they offer only little gain in the rate. In this work, we propose two new centralized coded caching schemes with low subpacketization and moderate rate gains utilizing projective geometries over finite fields. Both the schemes achieve the same asymptotic subpacketization, which is exponential in $O((\log K)^2)$ (thus improving on the $\sqrt[r]{K}$ exponent). The first scheme has a larger cache requirement but has at most a constant rate (with increasing $K$), while the second has small cache requirement but has a larger rate. As a special case of our second scheme, we get a new linear subpacketization scheme, which has a more flexible range of parameters than the existing linear subpacketization schemes. Extending our techniques, we also obtain low subpacketization schemes for other multi-receiver settings such as distributed computing and the cache-aided interference channel. We validate the performance of all our schemes via extensive numerical comparisons. For a special class of symmetric caching schemes with a given subpacketization level, we propose two new information theoretic lower bounds on the optimal rate of coded caching.

## Index Terms

Coded caching for broadcast channel, projective geometry, distributed computing, interference networks, subpacketization for coded caching.

## I. INTRODUCTION

Present and future wireless communication systems (4G,5G and beyond) are becoming more content-centric. The majority of this content is video, which is generated well ahead of transmission. Further, it is predicted in [5] that by 2022, four-fifths of all the internet traffic will be video. Therefore, we require new strategies to manage the data-heavy communication systems while ensuring quality of service.

Caching has been in vogue to lay off the traffic during the peak times in the network by storing part of the information demanded by the users (clients) in local storage known as *caches*. In this way, during the peak hours, the server can transmit only the non-cached information, thus reducing the traffic. *Coded caching* was proposed in a landmark paper by Ali-Niesen [6] to exploit this aspect of cached content to reduce the network congestion in the peak traffic time by prefetching part of the content in cost-effective cache available at the users during the off-peak time, while using coded transmissions during the peak times.

In [6], the authors considered an error-free broadcast channel with a server containing $N$ files of the same size and $K$ users (clients) each having a cache capable of storing $M$ files, where $M \leq N$. According to the scheme presented in [6] the system operates in two phases. During the *caching phase* (happens in the off-peak time) each file in the server is divided into $F$ equal-sized subfiles ($F$ is known as the *subpacketization* parameter), and placed in the caches of the clients. The caching phase occurs well ahead of the appearance of receiver demands, and thus the caching phase has to be designed in a demand-oblivious manner. During the *delivery phase* (happens in the peak time) each user demands a file from the server. Based on

the demands and the cache contents of the users, the server makes multiple coded transmissions. The goal is to design the caching and delivery phase so that the demands of all the users are satisfied. Since the caching phase is also designed centrally, this framework for coded caching is known as the *centralized coded caching*. Decentralized coded caching was introduced in [7], in which a coded delivery scheme is shown to achieve large gains in the rate, under a random or decentralized caching phase. Other important variations of this setting include coded caching in a popularity-based caching setting [8], online coded caching [9], and hierarchical coded caching [10]. We assume the basic centralized coded caching framework in the present work.

The delivery scheme in [6] serves $\gamma = 1 + \frac{MK}{N}$ users per transmission. The parameter $\gamma$ is known as the *global caching gain* and the *rate* of the scheme is given as $R = \frac{K\left(1 - \frac{M}{N}\right)}{\gamma}$, which has a $\Theta(K)$ gain over the uncoded delivery rate $\left(\text{for constant } \frac{M}{N}\right)$ which is $K\left(1 - \frac{M}{N}\right)$. The rate achieved by Ali-Niesen scheme [6] was shown to be optimal for a given cache size $M$ in [11], under the assumption of uncoded cache placement and $N \geq K$. Further, for large $K$, this scheme surprisingly achieves (approximately) a rate that is independent of $K$. However, it suffers from the problem of large subpacketization, which we now describe.

The achievability of the scheme shown in [6] is ensured by splitting each file into $F = \binom{K}{MK/N}$ equal-sized subfiles, where $F$ is known as the *subpacketization level* or simply, subpacketization. It was noticed in [12] that for this scheme, the subpacketization required grows exponentially in $K$ for constant $\frac{M}{N}$ as $K$ grows large (as $\binom{K}{Kp} \approx 2^{KH(p)}$ for constant $0 < p < 1$, where $H(p)$ is the binary entropy). For instance, with $K = 25$ users and with the capability to store one-fifth of the file library in the cache of each user $\left(\frac{M}{N} = \frac{1}{5}\right)$, the subpacketization becomes $\binom{25}{5}$, which is 53130. A high subpacketization requirement poses multiple issues in the implementation of coded caching. A straightforward issue is that of the size of the file itself; the file-size has to be at least as large as the product of subpacketization $F$ and the size of any accessible file-segment. Assuming a file-segment size of about 512 KB, the file-size has to be at least 27 GB for a subpacketization level of 53130, which is prohibitive in practice. Higher subpacketization levels also mean higher indexing overheads to identify the subfiles. Also, a large subpacketization level implies smaller chunks of the file, which in turn means higher normalized read times from the storage media, along with search-and-read overheads. Having longer indexing overheads to identify the subfiles is also a factor to consider when there are a large number of subfiles in any file. Because of these reasons, coded caching schemes with low-subpacketization and with good caching gains are preferable for practical applications.

Since this problem was identified, a number of papers, for instance [12]–[20] have presented new schemes for the coded caching which use smaller subpacketization than [6] at the cost of having increased rate for the same cache requirement compared to [6]. The summary of some important known schemes is given in Table I. The third column lists the cache fraction of any file $\left(\text{a fraction } \frac{M}{N} \text{ of each file is cached by a user}\right)$. Many of the schemes presented in Table I require exponential subpacketization (in $K$, for large $K$), as shown in the fourth column of Table I to achieve a constant rate (shown in the fifth column). The asymptotics of the rate are presented in the last column of Table I. A user-grouping technique combined with the scheme of [6] was used in [12] to reduce the subpacketization at the cost of rate. A variety of techniques and structures from combinatorics, coding theory, and graph theory, have also been employed to obtain many of these constructions. For instance, in [13], a special combinatorial structure called *placement delivery arrays (PDA)* was presented and used to construct coded caching schemes with reduced subpacketization than [6] (though it remained still exponential in $K$). The work [16] used resolvable designs obtained from linear block codes for the same. The papers [14], [15] used hypergraphs and bipartite graphs respectively, and showed schemes with subpacketization subexponential in $K$. The subpacketization of particular schemes of [19] have been shown to be subexponential, while some schemes of [18] have subpacketization that is linear or polynomial (in $K$) at the cost of either requiring much larger cache $M$ or much larger rate compared to [6]. Interestingly, a linear subpacketization scheme ($F = K$) was shown in [17] using a graph theoretic construction with near constant rate and small memory requirement. However, the construction in [17] holds for very large values of $K$ only. In [14], it was shown that the schemes with subpacketization linear in $K$ is impossible if we require constant rate. In [20], the authors consider caching schemes without file splitting, i.e., the scenario when $F = 1$. Constructions of low-subpacketization coded caching schemes is an active area of research, with many other recent works such as [21]–[25] presenting new or modifying existing constructions, using a variety of techniques including combinatorial designs, bipartite graphs, orthogonal arrays, covering arrays, etc.

| Scheme | Number of Users $\mathbf{K}$ | Cache Fraction $\frac{\mathbf{M}}{\mathbf{N}}$ | Subpacketization $\mathbf{F}$ | Rate $\mathbf{R}$ Expression | Asymptotics |
|---|---|---|---|---|---|
| Ali-Niesen [6] | any $K$ | $\frac{M}{N}$ for $M < N$ $\frac{MK}{N} \in \mathbb{Z}^+$ | $O\left(2^K\right)$ | $\frac{K\left(1-\frac{M}{N}\right)}{\frac{MK}{N}+1}$ | $O(1)$ |
| Ali-Niesen scheme with Grouping [12] | $K$ | Same as [6] | $O\left(e^g\right)$ | $\frac{K}{g+1}\left(1-\frac{1}{\lceil\frac{N}{M}\rceil}\right)$, where $g \in \mathbb{Z}^+$ such that $\frac{K}{g\lceil\frac{N}{M}\rceil} \in \mathbb{Z}^+$. | $O\left(\frac{K}{g}\right)$ |
| Yan et al. [13] based on Placement Delivery Arrays (PDAs) | Any $K$ | $1-\frac{1}{q}$ or $\frac{1}{q}$ | $O\left(e^K\right)$ | $\frac{K\left(1-\frac{M}{N}\right)}{\frac{MK}{N}}$ | $O(1)$ |
| Shangguan et al. [14] (PDAs based on hypergraphs) | Specific choices | $1-\frac{1}{q}$ or $\frac{1}{q}$ | $O\left(e^{\sqrt{K}}\right)$ | $R \approx (2q-1)^2$, such that $q = \frac{\lambda}{2}$, where $\lambda$ is such that $\frac{M}{N} = \frac{2\lambda-1}{\lambda^2}$ | $O(1)$ |
| Yan et al. [15] (for integers $0 < a, b < m$ and $\lambda < min\{a,b\}$ based on strong edge coloring of bipartite graph) | $\binom{m}{a}$ | $\frac{\binom{a}{\lambda}\binom{m-a}{b-\lambda}}{\binom{m}{a}}$ | $\binom{m}{b}$ | $\frac{\binom{m}{a+b-2\lambda}\binom{a+b-2\lambda}{a-\lambda}}{\binom{m}{b}}$ | — |
| Tang et al. [16] based on resolvable designs using $n$ length linear block code of rate $c$ | $nq$ (for some constant $q$) | $1-\frac{1}{q}$ or $\frac{1}{q}$ | $O\left(q^K\right)$ | $\frac{K\left(1-\frac{M}{N}\right)}{\frac{cK}{q}+1}$ | $O(1)$ |
| Scheme from [17] based on induced matchings of a Ruzsa Szemeredi graph | $K$ (necessarily extremely large) | $K^{-\epsilon}$ (some small $\epsilon$) | $K$ | $K^\delta$ (some small $\delta$) | $K^\delta$ |
| PDA scheme $P_1$ from Cheng et al. [18] $k, t \in \mathbb{Z}^+$ | $\binom{k}{t+1}$ | $1-\frac{t+1}{\binom{k}{t}}$ | $\binom{k}{t}$ | $\frac{k}{\binom{k}{t}}$ | — |
| Two PDA Schemes from [19] $q, z, m, t \in \mathbb{Z}^+, q \geq 2, z < q, t < m$ | $\binom{m}{t}q^t$ and $(m+1)q$ | $1-\left(\frac{q-z}{q}\right)^t$ and $\frac{z}{q}$ | $O\left(q^{t\sqrt{K}}\right)$ | $\left((q-z)/\lfloor\frac{q-1}{q-z}\rfloor\right)^t$ and $(q-z)/\lfloor\frac{q-1}{q-z}\rfloor$ | $O(1)$ |

The coded caching scheme proposed in [6] was extended to a variety of other settings, including device-to-device communication (D2D) networks [26], distributed computing [27], and interference management in wireless interference channels [28]. Every one of these settings can be modelled as a multi-client communication scenario with one or more transmitters, with the clients (receivers), and in some situations the transmitters as well, having cache. This enables *coded* transmissions, which generates rate advantages in all such situations. Because the fundamental scheme of [6] is adapted to each of these settings, the subpacketization issue continues to affect the adapted schemes as well. In fact, the problem is sometimes exacerbated because the special features of the setting requires a further division of the subfiles into smaller packets (for instance, the subpacketization in the D2D scheme of [26] is $\frac{MK}{N}\binom{K}{MK/N}$, thus having a multiplicative factor of $\frac{MK}{N}$ over that of the subpacketization of [6]).

In this work, we construct low-subpacketization schemes for coded caching utilizing ideas from graph theory and projective geometry over finite fields. We give the summary of the contributions of this paper in the next section.

## II. SUMMARY OF MAIN CONTRIBUTIONS

The contributions and organization of this work are as follows. In Section III, we review the formal system model for coded caching on broadcast networks from [6]. In Section III-B, we present the bipartite graph model for coded caching given in [15]. For *symmetric* coded caching schemes in which the caching is file-index invariant, the bipartite graph model captures the caching scheme and a specific class of delivery schemes in the form of a bipartite graph and its subgraphs.

The central contribution of this work is the construction of coded caching schemes which have subpacketization subexponential in the number of users $K$ (for large $K$). Using the bipartite graph framework, and utilizing some basic ideas from the

TABLE II: Summary of asymptotic behaviour (as $K$ grows large) of the coded caching schemes for the broadcast channel presented in this paper. (where $n, m \in \mathbb{Z}^+$ and $q$ is prime power.)

| Scheme name and Characteristics | Cache fraction $\left(\frac{M}{N}\right)$ | Subpacketization $(F)$ | Rate $(R)$ | Results, examples, numerical comparisons |
|---|---|---|---|---|
| **Low (subexponential) subpacketization scheme with large cache fraction** <br><br> **(Scheme A)** | $\leq constant\ (\geq 0.5)$ | $O(poly(K))$ | $\Theta(K)$ | Section IV-C Theorem 2 |
| | $1 - \Theta\left(\frac{1}{\sqrt{K}}\right)$ | $q^{O\left((\log_q K)^2\right)}$ <br><br> **(subexponential in $K$)** | $O(1)$ | Table VI <br> Examples 3,4 <br> Appendix B |
| **Low (subexponential) subpacketization scheme with small cache fraction** <br><br> **(Scheme B)** | $\leq constant$ | $q^{O\left((\log_q K)^2\right)}$ <br><br> **(subexponential in $K$)** | $\Theta\left(\frac{K}{(\log_q K)^n}\right)$ | Section V-C <br> Theorem 3 <br> Table VII <br> Table VIII <br> Table IX <br> Example 5 <br><br> Appendix D |
| **Linear subpacketization $(F = K)$ scheme with small cache fraction (Scheme C)** | $\leq constant$ | $= K$ <br><br> **(linear in $K$)** | $\Theta\left(\frac{K}{(\log_q K)^n}\right)$ | Section V-E Corollary 2) <br><br> (Theorem 3 with $n = m$) |

projective geometries over finite fields, we construct two new coded caching schemes, given in Section IV (Scheme A) and Section V (Scheme B). We briefly describe the salient features of these schemes. Scheme A, which we present in Section IV, has subexponential subpacketization, but uses high cache size, to achieve a rate upper bounded by a constant. When compared to Scheme A, Scheme B presented in Section V has equivalent subpacketization (as an asymptotic function of $K$), while having lower cache size but a higher rate. Because of the low-cache requirement which is relevant in practice, we consider Scheme B to be the more important among the two schemes presented. For ranges of users from 10s-1000s and cache-fraction in the range 0.05 to 0.35, we show via numerical examples that we obtain the practical values of subpacketization in the range of $10^2 - 10^4$, achieving the rates in the range $10 - 100$, (serving number of users in the range of $3 - 15$ per transmission). These numerical comparisons of Scheme A and Scheme B with existing schemes are shown in Table VI (Section IV), and Tables VII,VIII and IX (Section V) respectively.

Observing the subpacketization and the rate of any given scheme, as the number of users grow, gives us an understanding of the performance of the scheme, and is done in the prior literature also (see the subpacketization column of Table I, for instance). It also enables comparison with the existing schemes, in situations when parameters cannot be matched accurately. Asymptotically in $K$, both the schemes A and B achieve subpacketization $F = q^{O\left((\log_q K)^2\right)}$, Scheme A achieves this $F$ when $R = O(1)$ and $\frac{M}{N} = 1 - \Theta\left(\frac{1}{\sqrt{K}}\right)$, whereas Scheme B achieves this $F$ when $R = \Theta\left(\frac{K}{(\log_q K)^n}\right)$ and $\frac{M}{N}$ is smaller than some constant (where $n$ is an integer constant and $q$ is some prime power). So clearly, Scheme A requires a large cache size that grows with $K$, whereas in scheme $B$ the cache fraction can be maintained constant. In the regime when $\frac{M}{N}$ is smaller than

TABLE III: Summary of the subpacketization dependent lower bounds on the rate presented in this paper. Here $D = F(1 - \frac{M}{N})$ and $d = K(1 - \frac{M}{N})$.

| Corollary 3 (Section VII) Table XII (Symmetric caching with every user caching the same number of subfiles) | $R^* F \geq (K + F)\left(1 - \frac{M}{N}\right) - 1$ |
|---|---|
| Theorem 9 (Section VII) Table XII (Symmetric caching with every user caching the same number of subfiles and every subfile cached in the same number of users) | $R^* F \geq D + \left\lceil \frac{(d-1)D}{K-1} \right\rceil + \cdots + \left\lceil \frac{1}{\frac{KM}{N}+1} \left\lceil \frac{2}{\frac{KM}{N}+2} \left[ \cdots \left\lceil \frac{d-2}{K-2} \left\lceil \frac{(d-1)D}{K-1} \right\rceil \right\rceil \cdots \right] \right\rceil \right\rceil$ |

some constant, Scheme A achieves subpacketization $F = O(poly(K))$ when $R = \Theta(K)$ (thus the gain of coded caching is small and therefore this regime is not very interesting). These asymptotics are proved in respective sections, and also captured in Table II. The last column of Table II provides the location of the related results, numerical comparisons and examples in the paper.

For specific values of the scheme parameters, we get a new linear subpacketization scheme (Scheme C) in Section V-E, parametrized with two parameters, $q$ a prime power and $\lambda \in (0,1)$. For this scheme, we get the number of users $K \leq \frac{q^{2\lambda^2 q^2}}{(\lambda q)!}$, the subpacketization level $F = K$, the cache fraction $\frac{M}{N} \leq \lambda$, and the rate being $\frac{K(1-M/N)}{\gamma}$, where $\gamma \geq \frac{4^{\lambda q}}{2\sqrt{\lambda q}}$. The asymptotics of Scheme C and location of its results in the paper are captured in Table II. A generalized version of Scheme B, with one more tunable parameter, is presented in Section V-F.

In Section VI, we extend our low-subpacketization low-cache Scheme B to some other settings explored in the literature where coded caching helps. In particular, we extend our Scheme B to the distributed computing in Section VI-A (Theorem 6), and to the cache-aided interference channel setting in Section VI-B (Theorem 7). In each of these settings, our extended scheme is compared with the existing schemes numerically to illustrate our low subpacketization advantage (Tables X, XI) in Section VI.

Utilizing the perspective that is given by the bipartite graph model in Section III-B, we obtain two information-theoretic lower bounds for the rate of coded caching schemes with some fixed finite subpacketization level in Section VII. Prior literature, for instance, in [11], [18], has such lower bounds on the rate. However, not all of them take into account the finite-subpacketization constraint. The two lower bounds on the rate we obtain for given parameters $K, F, M$ and $N$, are given in Table III in this section. In this table, the parameter $D \triangleq F\left(1 - \frac{M}{N}\right)$ and $d \triangleq K\left(1 - \frac{M}{N}\right)$. Using numerical examples, we also compare the performance of these lower bounds with the bounds from [11], [18] in Table XII (Section VII). We conclude the paper with discussions regarding future work in Section VIII.

*Notations and terminology:* $\mathbb{Z}^+$ denotes the set of positive integers. We denote the set $\{1, \ldots, n\}$ by $[n]$ for some $n \in \mathbb{Z}^+$. The empty set is denoted by $\phi$. For sets $A, B$, the set of elements in $A$ but not in $B$ is denoted by $A \backslash B$. The set $A \backslash a$ denotes $A \backslash \{a\}$. The set of all $b$ sized subsets of $A$ is denoted by $\binom{A}{b}$. For $a, b \in \mathbb{Z}^+$ such that $1 \leq a \leq b$, $\binom{a}{b}$ represents the binomial coefficient. A set $\{A_1, A_2, \cdots, A_n\}$ is said to partition the set $A$ if $\bigcup_{i=1}^{n} A_i = A$ (where the union is a disjoint union). The finite field with $q$ elements is $\mathbb{F}_q$. The $k$-dim (dimensional) vector space over $\mathbb{F}_q$ is represented as $\mathbb{F}_q^k$. The dimension of a vector space $V$ over $\mathbb{F}_q$ is given as $dim(V)$. The zero vector is represented as $\mathbf{0}$. For two subspaces $V, W$, their subspace sum is denoted by $V + W$. Note that $V + W = V \oplus W$ (the direct sum) if $V \cap W = \{\mathbf{0}\}$. The subspaces $V_1, \cdots, V_n$ (all are subspaces of a vector space) are said to be linearly independent if $v_1 + \cdots + v_n = \mathbf{0}$ (for each $v_i \in V_i, i \in [n]$) holds only for $v_1 = v_2 = \cdots = v_n = \mathbf{0}$, and linearly dependent otherwise. A graph $G$ is defined by its vertex set $V(G)$ and edge set $E(G) \subseteq \binom{V(G)}{2}$. A graph $H$ is said to be a subgraph of $G$ if $V(H) \subseteq V(G)$ and $E(H) \subseteq E(G)$ such that the vertices of each edge in $E(H)$ are in $V(H)$. Further, $H$ is said to be an induced subgraph of $G$ if $E(H)$ consists of all those edges of $G$ both vertices of which are in $V(H)$. For a graph $G$, the neighbourhood of a vertex $u \in V(G)$ is defined
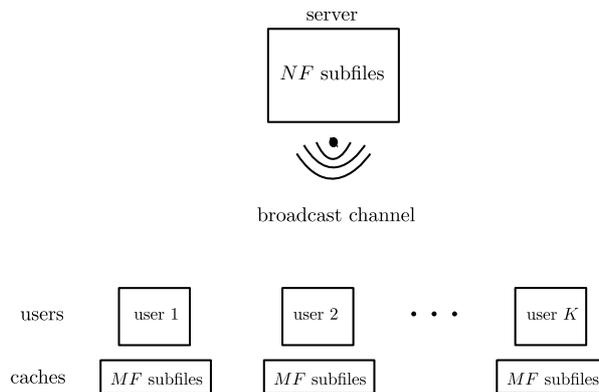
Fig. 1: Broadcast coded caching setup.

as $\mathcal{N}(u) = \{v \in V(G) : \{u, v\} \in E(G)\}$ and $|\mathcal{N}(u)|$ denotes the degree of $u$. A graph $G$ is said to be a bipartite graph if there exist $A, B \subseteq V(G)$ such that $A \cup B = V(G)$ (where the union is a disjoint union) and every edge in $E(G)$ connects a vertex in $A$ to a vertex in $B$. Further, $A$ and $B$ are called as the set of left and right vertices respectively. A bipartite graph $G$ is said to be left-regular (right-regular) if the degree of every left (right) vertex is the same. A bipartite graph $G$ is said to be bi-regular if it is both left and right regular. A graph $G$ is said to be regular if the degree of every vertex is the same. WLOG stands for "Without loss of generality".

## III. SYSTEM MODEL

In this section, we present the classical coded caching setup as given by Ali-Niesen in [6]. We then discuss the bipartite graph model for coded caching as given in [15].

### A. System Model

Consider a broadcast coded caching setup as shown in Fig. 1. Let $\mathcal{K}$ be the set of users (clients) ($|\mathcal{K}| = K$) in a system consisting of one server having a library of $N$ files, $\{W_i : i \in [N]\}$, connected to the clients via an error-free broadcast channel. We assume $K \leq N$, i.e., the number of files is larger than the number of users. This is typically the case in most research works in the coded caching literature (some exceptions exist, for instance [29]). It is also true in many practical situations, as the library (set of all files) has possibly many more files than the number of receivers in the network. Further, if the number of receivers is large, then they may be grouped into multiple groups, and coded caching may be applied to the modified system, in which each group is considered like a single user (see, for instance, [12]); in which case the setting of $K \leq N$ is more relevant.

Let $F$ be the subpacketization level, i.e., we assume each file is composed of $F$ subfiles, each taking values according to a uniform distribution from some finite Abelian group $\mathcal{A}$. The subfiles of file $W_i$ are denoted as $W_{i,f} : f \in \mathcal{F}$ for some set $\mathcal{F}$ of size $F$. Each user can store $M$ files (equivalently, $MF$ subfiles) in its cache. A *coded caching scheme* consists of two sub schemes (as in [6]), a *caching scheme* according to which subfiles of the files are placed in the user caches during periods when the traffic is low (the caching phase), and a *transmission scheme* or a *delivery scheme* that consists of broadcast transmissions from the server satisfying the demands of the clients appearing during the demand phase. We assume *symmetric caching* throughout the paper, i.e., for any $f \in \mathcal{F}$, any user either caches $W_{i,f}$ $\forall i \in [N]$ or does not cache $W_{i,f}$ $\forall i \in [N]$. Most schemes in the literature, including those for instance in [12]–[19], employ symmetric caching. The parameters $F$ and $\frac{M}{N}$ in the symmetric caching will lead to the quantity $\frac{MF}{N}$ being an integer, indicating the number of subfiles of any particular file stored in a user's cache.

During the demand phase, each client demands one file from the library. In the delivery scheme, the transmissions must be done so that the demands of the clients are all satisfied. The worst-case demand scenario corresponds to that situation in which each receiver demands a unique file. As in [6], the rate $R$ (for the worst-case demands) of such a coded caching scheme is defined as,
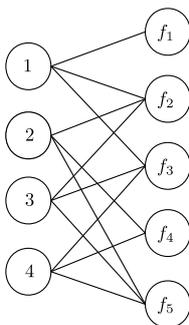
Fig. 2: The figure is a bipartite caching graph $B(K = 4, F = 5, D = 3)$ and represents a symmetric caching scheme with $4$ users, $5$ subfiles and cache fraction $\frac{M}{N} = \frac{2}{5}$. Edges indicate the missed subfiles.

$$\text{Rate } R \triangleq \frac{\substack{\text{Number of bits transmitted in the transmission} \\ \text{scheme considering worst-case demands}}}{\text{Number of bits in a file}}.$$

The delivery scheme typically involves multiple transmission rounds. Under the assumption that the size of each transmission is the same as that of the subfile size, the rate expression can be simplified as,

$$R = \frac{\substack{\text{Number of transmission rounds in the} \\ \text{transmission scheme for worst-case demands}}}{\text{Number of subfiles in a file } (F)}. \tag{1}$$

In an achievable scheme in which the delivery scheme consists of uncoded transmissions, note that the rate must be $K\left(1 - \frac{M}{N}\right)$, as the uncached fraction of each demanded file consisting of $F\left(1 - \frac{M}{N}\right)$ subfiles are sent uncoded. Note that in this scheme, each transmitted subfile is intended for one particular user only.

The ratio of the rate of the uncoded scheme to the rate $R$ of a coded caching scheme is defined as the *global caching gain* ($\gamma$) of the coded caching scheme. Thus,

$$\gamma = \frac{K(1 - M/N)}{R}. \tag{2}$$

The global caching gain $\gamma$ of a coded caching scheme also represents the average number of users served per transmission in the coded caching scheme. We are interested in designing coded caching schemes with low rate (or high gain) and low subpacketization level, which also uses low cache size at the users.

### B. Bipartite Graph based Coded Caching and Delivery based on [15]

We can visualize the symmetric caching scheme (with fully populated caches) using a bipartite graph, following [15]. We shall use this bipartite coded caching picture to obtain our coded caching schemes in Section IV and Section V, as well as to obtain lower bounds on the rate of coded caching in Section VII.

Consider a bipartite graph $B$ with $\mathcal{K}$ being the left (user) vertices and the right (subfile) vertices being $\mathcal{F}$. We then define the edges of the bipartite graph to denote the uncached subfiles of the files, i.e., for $k \in \mathcal{K}, f \in \mathcal{F}$, an edge $\{k, f\} \in E(B)$ exists if and only if the user $k$ does *not* contain in its cache the subfile $W_{i,f}, \forall i \in [N]$. Clearly, this bipartite graph is left-regular, with $F\left(1 - \frac{M}{N}\right)$ being the degree of any user vertex. Indeed any left-regular bipartite graph defines a caching scheme, which we formalize below.

**Definition 1** (Bipartite Caching Scheme, Bipartite Caching Graph). *Given a bipartite $D$-left-regular graph with $K$ left vertices and $F$ right vertices denoted by $B(K, F, D)$ (or in short, $B$), the symmetric caching scheme defined on $K$ users with subpacketization $F$ with the edges of $B$ indicating the uncached subfiles at the users, is called the $(K, F, D)$ bipartite caching scheme associated with the bipartite graph $B$. Further, $B(K, F, D)$ is known as the bipartite caching graph.*

**Remark 1.** *We observe that the bipartite caching scheme associated with the graph $B(K, F, D)$ has the cache fraction $\frac{M}{N} = 1 - \frac{D}{F}$.*

**Example 1.** *Fig. 2 shows a graph describing a* $(K = 4, F = 5, D = 3)$ *bipartite caching scheme. The cache-fraction is* $\frac{M}{N} = \frac{2}{5}$, *meaning that each receiver caches* 2 *out of the* 5 *subfiles in each file. For instance, the user* 2 *caches the subfiles* $W_{i,f_1}, W_{i,f_3}$ *and does not cache* $W_{i,f_2}, W_{i,f_4}, W_{i,f_5}, \forall i \in [N]$. *Similarly, the subfile* $W_{i,f_1}$ *is cached in the users* $2, 3, 4$, *and not at the user* 1, $\forall i \in [N]$ *(where, $N$ represents the number of files in the library).*

Most schemes in the literature can be captured via the bipartite caching model. The delivery scheme of such schemes are the so-called 'all-but-one' delivery schemes, in which each transmission is simply a sum of subfiles (one for each client in some subset of clients), with each summand subfile being demanded by a unique client in the subset, while the other summands are available in the unique clients cache. It turns out that this all-but-one delivery scheme can be captured in a graph-theoretic sense, as given in [15], an equivalent version of which we now briefly discuss.

A matching of a graph $G$ is a subset of edges with no common vertices between any two distinct edges in the subset. An *induced matching* $\mathcal{C}$ of a graph $G$ is a matching such that the induced subgraph of the vertices of $\mathcal{C}$ is $\mathcal{C}$ itself. The formal definition is as follows:

**Definition 2** (Induced Matching, Induced Matching Cover). *Consider a bipartite graph $B$. A set of edges $\mathcal{C} \subseteq E(B)$ is called an induced matching of $B$, if every $\{k_1, f_1\}, \{k_2, f_2\} \in \mathcal{C}$ satisfies $k_1 \neq k_2, f_1 \neq f_2$ and $\{k_1, f_2\}, \{k_2, f_1\} \notin E(B)$. A set of induced matchings $\{\mathcal{C}_i : i \in [S]\}$ is called an induced matching cover of $B$ if it partitions $E(B)$.*

Now, an induced matching cover $\{\mathcal{C}_i : i \in [S]\}$ of $B$ is equivalent to the set of color classes of an $S$-*strong-edge-coloring* of $B$. The work [15] connects the strong edge coloring of bipartite graphs to coded caching. For more details on the definition of strong edge coloring we refer the reader to [15].

Let $W_{d_k}$ denote the demanded file of the user $k \in \mathcal{K}$, in the demand phase. For an induced matching $\mathcal{C}$ of the bipartite caching graph $B$ consisting of the edges $\{\{k_i, f_i\} : i \in [g]\}$ (where $g$ represents the number of edges in the induced matching $\mathcal{C}$), consider the associated transmission

$$Y_{\mathcal{C}} = \sum_{i=1}^{g} W_{d_{k_i}, f_i}. \tag{3}$$

As $\mathcal{C}$ is an induced matching, $W_{d_{k_i}, f_i}$ is a subfile unavailable but demanded at the user $k_i$. By the same reason, each user $k_i$ has all the subfiles in (3) in its cache except for $W_{d_{k_i}, f_i}$, hence user $k_i$ can decode $W_{d_{k_i}, f_i}, \forall i \in [g]$. If $\{\mathcal{C}_i, i \in [S]\}$ is an induced matching cover then it is easy to see that the transmissions $Y_{\mathcal{C}_i} : i \in [S]$ (constructed as in (3)) corresponding to $\mathcal{C}_i : i \in [S]$ satisfy the demands of all the users, as all the edges of the bipartite caching graph are 'covered' by the induced matching cover. Therefore, we have obtained a valid delivery scheme. We refer to this delivery scheme as the *induced-matching based delivery scheme*. Here, the parameter $S$ represents the number of induced matchings in the induced matching cover or equivalently the number of transmissions in the associated delivery scheme. By (1), the rate of this transmission scheme is $\frac{S}{F}$. Further, if $|\mathcal{C}_i| = g, \forall i \in [S]$ then the rate will be $\frac{S}{F} = \frac{K(1-M/N)}{g}$ and the coded caching gain is $g$. Thus, each transmission corresponding to each induced matching will serve exactly $g$ users.

The above discussion of the bipartite caching scheme and the induced-matching based delivery scheme is summarized into the following theorem, which will be used to derive the coded caching parameters from the bipartite caching graphs that we construct in Section IV and Section V. The equivalent description of this theorem can be found in [15] in the language of strong edge coloring.

---

**Theorem 1.** *Consider a bipartite caching graph $B(K, F, D)$ with an induced matching cover $\{\mathcal{C}_i : i \in [S]\}$ such that $|\mathcal{C}_i| = g, \forall i \in [S]$. Then there is a coded caching scheme for a broadcast system with $K$ users, each with cache size $M$, with the number of files $N \geq K$, consisting of the caching scheme defined by $B(K, F, D)$ with subpacketization $F$, cache fraction $\frac{M}{N} = 1 - \frac{D}{F}$, and the associated delivery scheme based on the induced matching cover $\{\mathcal{C}_i : i \in [S]\}$ having rate $R = \frac{S}{F}$ and global caching gain $\gamma = g$.*

---

The following example illustrates an induced matching cover for the bipartite caching scheme presented in Example 1.

**Example 2** (Continuation of Example 1)**.** *Consider the following subsets of edges of the bipartite caching graph presented in Fig. 2.* $\mathcal{C}_1 = \{\{1, f_1\}, \{3, f_5\}\}$ , $\mathcal{C}_2 = \{\{1, f_2\}, \{4, f_5\}\}$ , $\mathcal{C}_3 = \{\{1, f_3\}, \{2, f_5\}\}$ , $\mathcal{C}_4 = \{\{2, f_2\}, \{4, f_3\}\}$ , $\mathcal{C}_5 = \{\{2, f_4\}, \{3, f_3\}\}$ , $\mathcal{C}_6 = \{\{3, f_2\}, \{4, f_4\}\}$. *Further* $\bigcup_{i=1}^{6} \mathcal{C}_i = E(B)$, *(where the union is a disjoint union). It is easy to see that each* $\mathcal{C}_i, \forall i \in [6]$ *is an induced matching and* $\{\mathcal{C}_i : i \in [6]\}$ *is an induced matching cover by Definition 2. Therefore, by Theorem 1, we have a delivery scheme corresponding to the induced matchings of B, with rate* $R = \frac{6}{5}$ *and global caching gain* $\gamma = 2$. *In Example 6, we show that this rate is optimal for the symmetric caching scheme defined by* $B(4, 5, 3)$.

## IV. A New Low Subpacketization Scheme with Large Cache Fraction (Scheme A)

In this section and in Section V, we present new coded caching schemes using the bipartite graph framework, which we have recollected in Section III-B, via tools from the projective geometry. Intuitively, these schemes combine ideas from the base-line coded caching scheme of [6], which is based on sets and set-containment, with ideas from the projective geometry, i.e., subspaces and subspace-containment (containment or the lack thereof of smaller subspaces within larger ones).

In this section we present a construction of coded caching scheme (Scheme A) which achieves a subpacketization which is subexponential in $K$, with cache fraction $(\frac{M}{N}) \geq 0.5$. Example 3 is a motivating example corresponding to the construction in Scheme A.

The construction of Scheme A (also Scheme B in Section V) is based on the bipartite graph approach recollected in Section III-B. We construct the bipartite caching graph by giving user vertices and subfile vertices, and then identify an induced matching cover in it. Then by using Theorem 1, we obtain the parameters of the coded caching scheme.

As our constructions use some simple results from the projective geometry, we first review some basic concepts from the projective geometry over finite fields and develop some mathematical terminology.

### A. Review of the Projective Geometries over Finite Fields [30]

Let $k, q \in \mathbb{Z}^+$ such that $q$ is a prime power. Let $\mathbb{F}_q^k$ be a $k$-dim vector space over a finite field $\mathbb{F}_q$. Let '**0**' represent the zero vector of $\mathbb{F}_q^k$. Consider the set of equivalence classes of $\mathbb{F}_q^k \setminus \{\mathbf{0}\}$ under the equivalence relation $\sim$ defined by $x \sim y$ if there is a nonzero element $\alpha \in \mathbb{F}_q$ such that $x = \alpha y$. The $(k-1)$-dim *projective space* over $\mathbb{F}_q$ is denoted by $PG_q(k-1)$ and is defined as the set of these equivalence classes. For $m \in [k]$, let $PG_q(k-1, m-1)$ denote the set of all $m$-dim subspaces of $\mathbb{F}_q^k$. From Chapter 3 in [30] it is known that $|PG_q(k-1, m-1)|$ is equal to the *Gaussian(or q)-binomial coefficient* $\begin{bmatrix} k \\ m \end{bmatrix}_q$ where,

$$\begin{bmatrix} k \\ m \end{bmatrix}_q \triangleq \frac{(q^k - 1) \dots (q^{k-m+1} - 1)}{(q^m - 1) \dots (q - 1)}. \text{ (where } k \geq m)$$

In fact, $\begin{bmatrix} k \\ m \end{bmatrix}_q$ gives the number of $m$-dim subspaces of any $k$-dim vector space over $\mathbb{F}_q$. Further, by definition, $\begin{bmatrix} k \\ 0 \end{bmatrix}_q = 1$. In any Gaussian binomial coefficient $\begin{bmatrix} a \\ b \end{bmatrix}_q$ given in this paper we assume that $a, b \in \mathbb{Z}^+$ and $1 \leq b \leq a$.

The following known results from [30] are used to describe our schemes (Scheme A and Scheme B).

**Lemma 1.** *[30] Consider a $k$-dim vector space $\mathbb{F}_q^k$. Let $1 \leq r, s, l < k$.*

A1: $\begin{bmatrix} k \\ r \end{bmatrix}_q = \begin{bmatrix} k \\ k - r \end{bmatrix}_q$.

A2: *The number of distinct $r$-dim subspaces of $\mathbb{F}_q^k$ containing a fixed $l$-dim subspace is* $\begin{bmatrix} k - l \\ r - l \end{bmatrix}_q$.

A3: *The number of distinct $r$-dim subspaces of $\mathbb{F}_q^k$ intersecting a fixed $s$-dim subspace in some $l$-dim subspace is*

$$q^{(r-l)(s-l)} \begin{bmatrix} k - s \\ r - l \end{bmatrix}_q \begin{bmatrix} s \\ l \end{bmatrix}_q.$$

We are essentially interested in finding out some asymptotic results of the schemes which we develop in next sections. For this reason, we use the following simple upper and lower bounds on Gaussian binomial coefficients and their relationships.

**Lemma 2.** *For non-negative integers $a, b, f$, for $q$ being some prime power,*

$$q^{(a-b)b} \leq \begin{bmatrix} a \\ b \end{bmatrix}_q \leq q^{(a-b+1)b} \tag{4}$$

$$q^{(a-f-1)b} \leq \frac{\begin{bmatrix} a \\ b \end{bmatrix}_q}{\begin{bmatrix} f \\ b \end{bmatrix}_q} \leq q^{(a-f+1)b} \tag{5}$$

$$q^{(a-f-b-1)\delta} \leq \frac{\begin{bmatrix} a \\ b \end{bmatrix}_q}{\begin{bmatrix} a \\ f \end{bmatrix}_q} \leq q^{(a-f-b+1)\delta}. \tag{6}$$

*where, $\delta = max(b, f) - min(b, f)$.*

*Proof:* The first lower bound for $\begin{bmatrix} a \\ b \end{bmatrix}_q$ is well known from the combinatorics literature (see, for instance, [31]). All the other bounds are proved by definition of the Gaussian binomial coefficient and by noting that $q^a - 1 \geq q^{a-1}$ (since $q \geq 2$), and $q^a - 1 \leq q^a$. This completes the proof. ∎

### B. An Illustrative Example

We now show an example which we shall see illustrates the idea behind the construction in Section IV-C.

**Example 3.** *Consider a caching system with $K = 7, F = 7, \frac{M}{N} = \frac{4}{7}$. To present this system, we need to provide the indexing for the users $(\mathcal{K})$ and the subfiles $(\mathcal{F})$. For this purpose, we consider some quantities from the projective geometry over finite fields. Consider a 3-dim vector space $(\mathbb{F}_2^3)$ over a binary field $\mathbb{F}_2$. Consider the 1-dim subspaces of $\mathbb{F}_2^3$ which are listed as follows,*

$$V_1 = span\{(0, 0, 1)\}$$
$$V_2 = span\{(0, 1, 0)\}$$
$$V_3 = span\{(1, 0, 0)\}$$
$$V_4 = span\{(1, 1, 0)\}$$
$$V_5 = span\{(1, 0, 1)\}$$
$$V_6 = span\{(0, 1, 1)\}$$
$$V_7 = span\{(1, 1, 1)\}.$$

*Let $\mathbb{V} = \{V_1, V_2, V_3, V_4, V_5, V_6, V_7\}$. Consider the 2-dim subspaces of $\mathbb{F}_2^3$ which are listed as follows,*

$$X_1 = \{(0,0,1),(0,1,0),(0,1,1),(0,0,0)\}$$
$$X_2 = \{(0,0,1),(1,0,0),(1,0,1),(0,0,0)\}$$
$$X_3 = \{(0,1,0),(1,1,0),(1,0,0),(0,0,0)\}$$
$$X_4 = \{(0,0,1),(1,1,0),(1,1,1),(0,0,0)\}$$
$$X_5 = \{(0,1,0),(1,0,1),(1,1,1),(0,0,0)\}$$
$$X_6 = \{(1,0,0),(0,1,1),(1,1,1),(0,0,0)\}$$
$$X_7 = \{(1,1,0),(0,1,1),(1,0,1),(0,0,0)\}.$$

Let $\mathbb{X} = \{X_1, X_2, X_3, X_4, X_5, X_6, X_7\}$. *We now proceed to describe the caching phase and the delivery phase.*

**Caching phase:** *Let $\underline{\mathcal{K} = \mathbb{V}}$ and $\underline{\mathcal{F} = \mathbb{X}}$. During the caching phase, every file $(W_i, i \in [N])$ is divided into $F = 7$ subfiles. The subfiles of $W_i$ are denoted as $W_{i,X}, \forall X \in \mathcal{F}$. The caching scheme is,*

- *For each $i \in [N]$, the user $V_l \in \mathcal{K}$ caches the subfile $W_{i,X}$ if $V_l$ is not a subspace of $X$.*

*Following this rule, we have the cached and uncached subfile indices of each user as shown in Table IV. The first two columns of Table IV, provides users and indices of cached subfiles respectively. Therefore, every user caches 4 out of the 7 subfiles in every file. Hence, the cache fraction $\frac{M}{N} = \frac{4}{7}$.*

TABLE IV: Indices of the cached and uncached subfiles for the caching scheme presented in Example 3.

| Users | Indices of cached subfiles | Indices of uncached (equivalently demanded) subfiles |
|---|---|---|
| $V_1$ | $X_3, X_5, X_6, X_7$ | $X_1, X_2, X_4$ |
| $V_2$ | $X_2, X_4, X_6, X_7$ | $X_1, X_3, X_5$ |
| $V_3$ | $X_1, X_4, X_5, X_7$ | $X_2, X_3, X_6$ |
| $V_4$ | $X_1, X_2, X_5, X_6$ | $X_3, X_4, X_7$ |
| $V_5$ | $X_1, X_3, X_4, X_6$ | $X_2, X_5, X_7$ |
| $V_6$ | $X_2, X_3, X_4, X_5$ | $X_1, X_6, X_7$ |
| $V_7$ | $X_1, X_2, X_3, X_7$ | $X_4, X_5, X_6$ |

**Delivery phase:** *Let demand of an arbitrary user $V_l \in \mathcal{K}$ be $W_{d_{V_l}}$. The demanded (uncached) subfile indices corresponding to each user are given in the last column of Table IV. To satisfy the demands of the users, the server transmits the following 7 transmissions.*

$$W_{d_{V_6},X_1} + W_{d_{V_5},X_2} + W_{d_{V_7},X_4}$$
$$W_{d_{V_1},X_1} + W_{d_{V_4},X_3} + W_{d_{V_7},X_5}$$
$$W_{d_{V_1},X_2} + W_{d_{V_2},X_3} + W_{d_{V_7},X_6}$$
$$W_{d_{V_3},X_3} + W_{d_{V_1},X_4} + W_{d_{V_6},X_7}$$
$$W_{d_{V_3},X_2} + W_{d_{V_2},X_5} + W_{d_{V_4},X_7}$$
$$W_{d_{V_2},X_1} + W_{d_{V_3},X_6} + W_{d_{V_5},X_7}$$
$$W_{d_{V_4},X_4} + W_{d_{V_5},X_5} + W_{d_{V_6},X_6}.$$

*Each transmission is a linear combination of 3 demanded subfiles. It is easy to see (using Table IV) that any user decodes their demanded subfiles. For instance, the user $V_1$ decodes $W_{d_{V_1},X_1}$ from the second transmission as it contains $W_{d_{V_4},X_3}$*

and $W_{d_{V_7}, X_5}$ in its cache. Similarly, the user $V_1$ decodes $W_{d_{V_1}, X_2}$ from the third transmission and $W_{d_{V_1}, X_4}$ from the fourth transmission. Thus, each transmission serves 3 users. Hence, the global caching gain is $\gamma = 3$. By using (1), the rate of the scheme is $R = \frac{7}{7} = 1$.

**Remark 2.** *The caching scheme of the Example 3 is obtained based on the idea of 'subspace containment', which we shall see motivates the caching scheme of our new construction in this section. We shall show in Example 4 in Section IV-C that the delivery scheme of Example 3 is obtained from an induced matching cover of the bipartite caching graph of the new construction.*

We are now ready to construct Scheme A, which is the main result of this section.

### C. Construction of Scheme A

We will construct Scheme A by constructing a bipartite caching graph $B(K, F, D)$ and identify an induced matching cover in it. Then by using Theorem 1, we obtain the coded caching parameters $(K, F, \frac{M}{N}, R, \gamma)$ of Scheme A.

Let $k, m, t, q$ be positive integers such that $m + t \leq k$ and $q$ is some prime power. Consider a $k$-dim vector space $\mathbb{F}_q^k$. Consider the following sets of subspaces which are used to index our user vertices, subfile vertices and the induced matchings of the bipartite caching graph $B$ which we construct.

$$\mathbb{V} \triangleq PG_q(k-1, t-1). \text{ (set of all } t\text{-dim subspaces)}$$
$$\mathbb{X} \triangleq PG_q(k-1, m+t-1). \text{(set of all } (m+t)\text{-dim subspaces)}$$
$$\mathbb{T} \triangleq PG_q(k-1, m-1). \text{ (set of all } m\text{-dim subspaces)}$$

Construct a bipartite graph $B$ with left (user) vertex set $\underline{\mathcal{K} = \mathbb{V}}$ and right (subfile) vertex set $\underline{\mathcal{F} = \mathbb{X}}$. Define the edge set of $B$ as,

$$E(B) \triangleq \{\{V, X\} : V \in \mathbb{V}, X \in \mathbb{X}, V \subseteq X\}.$$

We now find the values of $K, F$ and the degree of any user vertex (left degree) $D$.

**Lemma 3.** *The following relationships hold for the construction we have presented.*

$$K = |\mathbb{V}| = \begin{bmatrix} k \\ t \end{bmatrix}_q, \qquad F = |\mathbb{X}| = \begin{bmatrix} k \\ m+t \end{bmatrix}_q,$$

$$D \triangleq |\mathcal{N}(V)| = \begin{bmatrix} k-t \\ m \end{bmatrix}_q,$$

*where the last relationship holds for any $V \in \mathbb{V}$.*

*Proof:* By using the ideas presented in Section IV-A we can write,

$$K = |\mathbb{V}| = |PG_q(k-1, t-1)| = \begin{bmatrix} k \\ t \end{bmatrix}_q.$$

$$F = |\mathbb{X}| = |PG_q(k-1, m+t-1)| = \begin{bmatrix} k \\ m+t \end{bmatrix}_q.$$

Now, we will find $|\mathcal{N}(V)|$. Consider an arbitrary $V \in \mathbb{V}$. It is easy to see that $\mathcal{N}(V) = \{X \in \mathbb{X} : V \subseteq X\}$. Therefore, finding $|\mathcal{N}(V)|$ is equivalent to counting the number of $(m+t)$-dim subspaces of $\mathbb{F}_q^k$ containing the fixed $t$-dim subspace $V$. By applying A2 of Lemma 1 we get,

$$D = |\mathcal{N}(V)| = \begin{bmatrix} k-t \\ (m+t)-t \end{bmatrix}_q = \begin{bmatrix} k-t \\ m \end{bmatrix}_q.$$

This completes the proof. ■

Note that, by Lemma 3, $B$ is a $D$-left regular bipartite graph with $K$ left vertices and $F$ right vertices. Therefore, by Definition 1, $B(K, F, D)$ is a valid bipartite caching graph.

**Remark 3.** *Similar to left degree $D$, it is easy to see that the degree of any right (subfile) vertex $X \in \mathbb{X}$ is $|\mathcal{N}(X)| = |PG_q(m + t - 1, t - 1)| = \begin{bmatrix} m + t \\ t \end{bmatrix}_q$. Therefore, $B$ is a bi-regular bipartite caching graph.*

We now show that $B$ has an induced matching cover $\{\mathcal{C}_i : i \in [S]\}$ such that $|\mathcal{C}_i| = g, \forall i \in [S]$, for some $g, S \in \mathbb{Z}^+$.

**Induced matching cover:** The induced matchings of $B$, that we wish to obtain, is based on a relabeling of the edges $B$ based on $m$-dim subspaces of $\mathbb{F}_q^k$. Towards that end, we first require the following lemmas (Lemma 4 and Lemma 5) using which we can find 'matching' labels to the $t$-dim and $m$-dim subspaces of some $X \in \mathbb{X}$. Subsequently, using Lemma 6 and Lemma 7, we show the induced matching cover of $B$.

**Lemma 4.** *Consider some element $X \in \mathbb{X}$. Let $\left\{ V_i, i = 1, \ldots, \begin{bmatrix} m + t \\ t \end{bmatrix}_q \right\}$ denote the $t$-dim subspaces of $X$ taken in some fixed order. Then the set of $m$-dim subspaces of $X$ can be written as an indexed set as $\left\{ T_{V_i, X}, \ i = 1, \ldots, \begin{bmatrix} m + t \\ m \end{bmatrix}_q \right\}$ such that $T_{V_i, X} \oplus V_i = X, \forall i$ (where $\oplus$ denotes direct sum). Moreover, such an indexed set can be found in operations polynomial in $\begin{bmatrix} m + t \\ t \end{bmatrix}_q$.*

*Proof:* See Appendix A. ■

For a $t$-dim subspace $V_i$ contained in an $(m + t)$-dim subspace $X$, let $T_{V_i, X}$ (the $m$-dim subspace as obtained in Lemma 4 such that $T_{V_i, X} \oplus V_i = X$) be called *the matching subspace of $V_i$ in $X$*. Using these matching subspaces, we can obtain an alternate labeling scheme for the edges of our bipartite caching graph $B$. The alternate labels are given as follows:

- Let the alternate label for $\{V, X\}$ be $\{V, T_{V,X}\}$, where $T_{V,X}$ is the $m$-dim matching subspace of $V$ in $X$ obtained using Lemma 4.

The following lemma ensures that the alternative labeling given above is indeed a valid labelling, i.e., it uniquely identifies the edges of $B$.

**Lemma 5.** *No two edges of $B$ have the same alternate label.*

*Proof:* If $\{V_1, X_1\}, \{V_2, X_2\} \in E(B)$ have the same alternate label $\{V, T_{V,X}\}$, then clearly $V_1 = V_2 = V$. Moreover, we should also, by definition of the alternate labels, have that $X_1 = T_{V,X} \oplus V = X_2$. Therefore $\{V_1, X_1\} = \{V_2, X_2\}$. This completes the proof. ■

We are now in a position to present the induced matching cover of $B$. Our induced matchings (defined in Lemma 6) are represented in terms of the alternate labels given to the edges of $B$. We first show the structure of one such induced matching.

**Lemma 6.** *For an $m$-dim subspace $T \in \mathbb{T}$, consider the subset of $E(B)$ (identified by their alternate labels) as follows:*

$$\mathcal{C}_T \triangleq \{\{V, T\} \in E(B) : V \in \mathbb{V}\}.$$

*Then $\mathcal{C}_T$ is a $\begin{bmatrix} k - m \\ t \end{bmatrix}_q$ -sized induced matching of $B$.*

*Proof:* Firstly, we observe that $\mathcal{C}_T$ is a well-defined set because the $T$ is an $m$-dim subspace of precisely $\begin{bmatrix} k - m \\ (m + t) - m \end{bmatrix}_q$ subspaces of dimension $(m + t)$ by A2 of Lemma 1. Note that $\{V, T\}$ is the alternate label for $\{V, T \oplus V\} \in E(B)$. Also, we can observe that for distinct $\{V_1, T\}, \{V_2, T\} \in \mathcal{C}_T$, we must have $V_1 \oplus T \neq V_2 \oplus T$. This is due to the fact that each

$m$-dim subspace within an $(m + t)$-dim subspace $X$ is matched to a unique $t$-dim subspace of $X$. Hence, by Lemma 4 and our alternate labeling scheme, we should have $|\mathcal{C}_T| = \begin{bmatrix} k - m \\ (m + t) - m \end{bmatrix}_q = \begin{bmatrix} k - m \\ t \end{bmatrix}_q$.

We now show that $\mathcal{C}_T$ forms an induced matching of $B$. We do this using Definition 2. Consider two distinct and arbitrary edges $\{V_1, T\}, \{V_2, T\} \in \mathcal{C}_T$. These are the alternate labels for $\{V_1, T \oplus V_1\}, \{V_2, T \oplus V_2\}$ respectively. We have that $V_1 \neq V_2$, and have already checked that $T \oplus V_1 \neq T \oplus V_2$. Further, note that $V_1 \not\subset T \oplus V_2$. This is because if $V_1 \subset T \oplus V_2$, then, $T \oplus V_1 = T \oplus V_2$, which is not true by the definition of $\mathcal{C}_T$. Thus $\{V_1, T \oplus V_2\} \notin E(B)$. Similarly, $\{V_2, T \oplus V_1\} \notin E(B)$. By invoking the Definition 2, it is clear that $\mathcal{C}_T$ is an induced matching of $B$. This completes the proof. ■

We now show that the set of induced matchings $\{\mathcal{C}_T : T \in \mathbb{T}\}$ partition $E(B)$.

**Lemma 7.**

$$\bigcup_{T \in \mathbb{T}} \mathcal{C}_T = E(B),$$

*where the above union is a disjoint union (the induced matchings $\mathcal{C}_T$ are as defined in Lemma 6).*

*Proof:* It should be clear from our alternate labeling scheme and the definition of $\mathcal{C}_T$ that any edge $\{V, X\} \in E(B)$ (which gets some alternate label $\{V, T_{V,X}\}$) appears in the induced matching of $B$ given by $\mathcal{C}_{T_{V,X}}$. Furthermore, by definition $\mathcal{C}_{T_1}$ and $\mathcal{C}_{T_2}$ are disjoint for any distinct $T_1$ and $T_2$ in $\mathbb{T}$. This completes the proof. ■

Therefore, $\{\mathcal{C}_T : T \in \mathbb{T}\}$ is an induced matching cover of $B$. In the light of the construction of the bipartite caching graph $B$ of this section and the induced matching cover of $B$, we are now ready to present our coded caching scheme, Scheme A, in the following theorem.

---

**Theorem 2.** *(Scheme A) Let $k, m, t$ be positive integers such that $m + t \leq k$ and $q$ be any prime power. The bipartite graph $B$ constructed in Section IV-C is a $B(K, F, D)$ bipartite caching graph with an induced matching cover having $S = \begin{bmatrix} k \\ m \end{bmatrix}_q$ induced matchings, each having $g = \begin{bmatrix} k - m \\ t \end{bmatrix}_q$ edges and defines a coded caching scheme with,*

$$K = \begin{bmatrix} k \\ t \end{bmatrix}_q, \qquad\qquad F = \begin{bmatrix} k \\ m + t \end{bmatrix}_q,$$

$$\frac{M}{N} = 1 - \frac{\begin{bmatrix} k - t \\ m \end{bmatrix}_q}{\begin{bmatrix} k \\ m + t \end{bmatrix}_q}, \qquad R = \frac{\begin{bmatrix} k \\ m \end{bmatrix}_q}{\begin{bmatrix} k \\ m + t \end{bmatrix}_q},$$

*Global caching gain $\gamma = \begin{bmatrix} k - m \\ t \end{bmatrix}_q.$*

---

*Proof:* From Lemma 3, we get the expressions of $K, F$ and $D$. By Lemma 6 and Lemma 7, the size of the induced matchings of $B$ is $g = |\mathcal{C}_T| = \begin{bmatrix} k - m \\ t \end{bmatrix}_q$ for any $T \in \mathbb{T}$ and they partition the edge set $E(B)$. Further, by Lemma 6, the number of induced matchings in the induced matching cover is $S = |\mathbb{T}| = |PG_q(k - 1, m - 1)| = \begin{bmatrix} k \\ m \end{bmatrix}_q$. Hence, the bipartite graph $B$ satisfies all the conditions in Theorem 1. Therefore, there exists a coded caching scheme with $K$ users, subpacketization $F$,

$$\frac{M}{N} = 1 - \frac{D}{F} = 1 - \frac{\begin{bmatrix} k - t \\ m \end{bmatrix}_q}{\begin{bmatrix} k \\ m + t \end{bmatrix}_q},$$
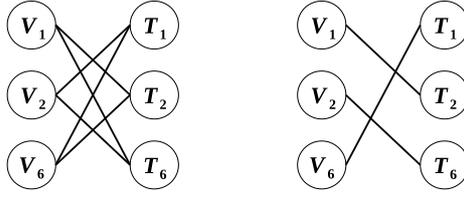
Fig. 3: Illustration of the relabeling procedure presented in Example 4: The left figure is a bi-regular bipartite graph $(B_1)$ with the left and right vertices being the $t$-dim and $m$-dim subspaces of $X_1$ respectively. Note that $\{V_i, T_j\} \in E(B)$ if $V_i \oplus T_j = X_1$ where $i, j \in \{1, 2, 6\}$. The right figure is a perfect matching of $B_1$. Hence, the relabels of $\{V_1, X_1\}, \{V_2, X_1\}, \{V_6, X_1\}$ are $\{V_1, T_2\}, \{V_2, T_6\}, \{V_6, T_1\}$ respectively.

$$R = \frac{S}{F} = \frac{\begin{bmatrix} k \\ m \end{bmatrix}_q}{\begin{bmatrix} k \\ m+t \end{bmatrix}_q},$$

Global caching gain $\gamma = g = \begin{bmatrix} k - m \\ t \end{bmatrix}_q$.

This completes the proof. ∎

**Example 4.** *(Continuation of Example 3) Example 3 gives us an illustration of our Scheme A for the values of $t = 1, m = 1, k = 3$ and $q = 2$. It is not difficult to see that the bipartite caching graph $B(K, F, D)$ obtained according to Scheme A gives us the caching scheme as given in Example 3. We now illustrate how we obtained the delivery scheme as shown in Example 3, using the induced matching cover of $B$ obtained through perfect matchings of the bipartite graph constructed in Appendix A (proof of Lemma 4) and our alternate labelling scheme.*

*In the bipartite graph terminology we have recollected, the edge set of the bipartite caching graph $B$ can be inferred from Table IV. For instance, the edges incident at the user vertex $V_1$ are $\{V_1, X_1\}, \{V_1, X_2\}, \{V_1, X_4\}$ and the edges incident at the subfile vertex $X_1$ are $\{V_1, X_1\}, \{V_2, X_1\}, \{V_6, X_1\}$. First, we relabel the edges of $B$ as per Lemma 4. Let $\mathbb{T}$ be the set of all $m$-dim subspaces of $\mathbb{F}_2^3$. Since $t = m = 1$, we can consider $\mathbb{T} = \mathbb{V}$ with $T_i = V_i, i \in [7]$. Consider an $(m + t)$-dim subspace $X_1$. This can be written as $X_1 = V_1 \oplus T_2 = V_1 \oplus T_6 = V_2 \oplus T_1 = V_2 \oplus T_6 = V_6 \oplus T_1 = V_6 \oplus T_2$.*

*The relabeling procedure is illustrated in Fig. 3. The left figure represents a bipartite graph with the $t$-dim ($m$-dim) subspaces of $X_1$ as the left (right) vertices. The edges denote the subspaces in direct sum. This bipartite graph is not the caching graph, but it corresponds to the one defined in Appendix A used to find the relabeling. The right figure of Fig. 3 denotes a perfect matching of the bipartite graph. In the same way, for each $(m + t)$-dim subspace $X \in \mathbb{X}$, a perfect matching is obtained. Following the edges of the perfect matchings gives us the new labels for the elements of $E(B)$, which are presented in Table V.*

*Now we describe how we can get the delivery scheme using the induced matchings of the bipartite caching graph $B$, which are indexed using $\mathbb{T}$. For each $m$-dim subspace $\mathbb{T}$ we can define an induced matching of $B$ as per Lemma 6 using the alternate labels for the edges of $B$. For instance, the induced matching corresponding to $T_1$ is $\mathcal{C}_{T_1} = \{\{V_6, T_1\}, \{V_5, T_1\}, \{V_7, T_1\}\} = \{\{V_6, X_1\}, \{V_5, X_2\}, \{V_7, X_4\}\}$. Therefore, the transmission corresponding to $\mathcal{C}_{T_1}$ is $W_{d_{V_6}, X_1} + W_{d_{V_5}, X_2} + W_{d_{V_7}, X_4}$ as mentioned in Example 3. Similarly, other transmissions can be obtained from the induced matchings $\mathcal{C}_T : T \in \mathbb{T}$, corresponding to which the transmissions are shown in Example 3.*

### D. Asymptotic Analysis and Numerical Comparisons of Scheme A with Ali-Niesen Scheme

In Appendix B, we provide the asymptotic analysis (as $K$ grows) for the scheme presented in Theorem 2. We have provided two cases. In case 1 of Appendix B, the cache fraction is upper bounded by a constant ($\frac{M}{N} \leq constant$) and our scheme has subpacketization, $F = O(poly(K))$. The rate, however, increases linearly with $K$ i.e., $R = \Theta(K)$ (similar to the uncoded caching rate); hence this regime does not have much significance. In case 2 of Appendix B, we keep the rate upper bounded

TABLE V: Old and new labels of the edges in $E(B)$ through perfect matchings, for the bipartite caching graph $B$ presented in Example 4.

| Old labels | New labels | Old labels | New labels |
|------------|------------|------------|------------|
| $\{V_1, X_1\}$ | $\{V_1, T_2\}$ | $\{V_7, X_4\}$ | $\{V_7, T_1\}$ |
| $\{V_2, X_1\}$ | $\{V_2, T_6\}$ | $\{V_2, X_5\}$ | $\{V_2, T_5\}$ |
| $\{V_6, X_1\}$ | $\{V_6, T_1\}$ | $\{V_5, X_5\}$ | $\{V_5, T_7\}$ |
| $\{V_1, X_2\}$ | $\{V_1, T_3\}$ | $\{V_7, X_5\}$ | $\{V_7, T_2\}$ |
| $\{V_3, X_2\}$ | $\{V_3, T_5\}$ | $\{V_3, X_6\}$ | $\{V_3, T_6\}$ |
| $\{V_5, X_2\}$ | $\{V_5, T_1\}$ | $\{V_6, X_6\}$ | $\{V_6, T_7\}$ |
| $\{V_2, X_3\}$ | $\{V_2, T_3\}$ | $\{V_7, X_6\}$ | $\{V_7, T_3\}$ |
| $\{V_3, X_3\}$ | $\{V_3, T_4\}$ | $\{V_4, X_7\}$ | $\{V_4, T_5\}$ |
| $\{V_4, X_3\}$ | $\{V_4, T_2\}$ | $\{V_5, X_7\}$ | $\{V_5, T_6\}$ |
| $\{V_1, X_4\}$ | $\{V_1, T_4\}$ | $\{V_6, X_7\}$ | $\{V_6, T_4\}$ |
| $\{V_4, X_4\}$ | $\{V_4, T_7\}$ | | |

TABLE VI: Comparison of the coded caching scheme in Theorem 2 (Scheme A) with the scheme in [12] (Ali-Niesen scheme with grouping). We match the cache fraction, the gain, and number of users as closely as possible, and compare the subpacketization level.

| Number of users | | Cache fraction | Global caching gain | Subpacketization | |
|------------------|------------------|----------------|---------------------|------------------|------------|
| $(k, m, t, q)$ $K_1$ (Theorem 2) | $(K', l)$ $K_2$ [12] | $\frac{M}{N}$ | $\gamma$ | $F_1$ (Theorem 2) | $F_2$ [12] |
| $(8, 3, 1, 2)$ 255 | $(60, 4)$ 240 | $\frac{16}{17}$ | 31 | 200787 | $10^{17}$ |
| $(6, 3, 2, 3)$ 11011 | $(24, 459)$ 11016 | $\frac{81}{91}$ | 13 | 364 | 2704156 |
| $(6, 3, 2, 2)$ 651 | $(12, 54)$ 648 | $\frac{16}{21}$ | 7 | 63 | 924 |
| $(7, 4, 1, 2)$ 127 | $(12, 10)$ 120 | $\frac{96}{127}$ | 7 | 2667 | 924 |

by a constant ($R = O(1)$) then our scheme has subexponential subpacketization, $F = q^{O\left((\log_q K)^2\right)}$ and cache fraction $\frac{M}{N} = 1 - \Theta\left(\frac{1}{\sqrt{K}}\right)$. Hence, the drawback of Scheme A is that it requires higher cache fraction even though it has subexponential subpacketization (when rate is upper bounded by a constant). The asymptotics of Scheme A obtained in this section are summarized in the first 2 rows of Table II in Section II.

In Table VI, we compare numerically the scheme in Theorem 2 with the scheme in [12] (Section V-A in [12], this is also given in row 2 of Table I), which is a modified version of the coded caching scheme in [6] with user grouping. The scheme in [12] is parameterized by the cache fraction $\frac{M}{N}$, global caching gain $\gamma$ and number of user groups $l$, and gives a scheme with the number of users $K_2 = K'l$ and subpacketization $F_2 = \binom{K'}{\gamma-1}$, where $K' = (\gamma - 1)\lceil\frac{N}{M}\rceil$. The number of users and subpacketization corresponding to Theorem 2, are labelled as $K_1$ and $F_1$ respectively. From the table it is clear that Scheme A performs better than [12], for most of the cases, in terms of the subpacketization.

We also observe from the construction of Scheme A as well as Table VI that the cache fraction of our scheme is at least 0.5 in general. To overcome the drawback of high cache requirement of Scheme A, we propose a new scheme (Scheme B) in Section V. Scheme B also uses ideas from the projective geometry.

## V. A New Low Subpacketization Scheme with Small Cache Fraction (Scheme B)

In this section, we present a coded caching scheme (Scheme B), which achieves subexponential (in $K$) subpacketization when memory is upper bounded by a constant. In Section V-E, we give Scheme C, which is a special case of Scheme B, and achieves linear subpacketization and is comparatively more flexible than the existing linear subpacketization schemes in the literature with non-trivial caching gain. Finally, in Section V-F, we show a generalized version of Scheme B which adds one more tunable parameter to the construction of Scheme B.

The construction of Scheme B uses similar techniques (bipartite graph and projective geometry) as that of Scheme A. Before presenting the new construction of bipartite caching graph $B(K, F, D)$ (which gives Scheme B), we first present some simple results using the projective geometry over finite fields which are used in this section.

### A. An useful lemma about sets of subspaces

Let $\mathbb{T} \triangleq PG_q(k - 1, 0)$. Let $\theta(k)$ denote the number of distinct 1-dim subspaces of $\mathbb{F}_q^k$. Therefore,

$$\theta(k) = |\mathbb{T}| = \begin{bmatrix} k \\ 1 \end{bmatrix}_q = \frac{q^k - 1}{q - 1}.$$

The following lemma will be used repeatedly in this section.

**Lemma 8.** *Let $k, a, b \in \mathbb{Z}^+$ such that $1 \leq a + b \leq k$. Consider a $k$-dim vector space $V$ over $\mathbb{F}_q$ and a fixed $a$-dim subspace $A$ of $V$. The number of distinct (un-ordered) $b$-sized sets $\{T_1, T_2, \cdots, T_b\} \subseteq \mathbb{T}$ and $A \oplus T_1 \oplus T_2 \oplus \cdots \oplus T_b \in PG_q(k-1, a+b-1)$ is $\frac{1}{b!}\prod_{i=0}^{b-1}(\theta(k) - \theta(a + i))$.*

*Proof:* First we find the number of $T_1 \in \mathbb{T}$ such that $A \oplus T_1$ is an $(a + 1)$-dim subspace of $V$. To pick such a $T_1$ we define, $T_1 = span(\mathbf{t_1})$ for some $\mathbf{t_1} \in V \setminus A$. Such a $\mathbf{t_1}$ can be picked in $(q^k - q^a)$ ways. However, for one such fixed $\mathbf{t_1}$, there exist $(q - 1)$ number of $\mathbf{t_1'} \in V \setminus A$ such that $span(\mathbf{t_1}) = span(\mathbf{t_1'}) = T_1$. Thus, the required number of unique $T_1 \in \mathbb{T}$ is $\frac{q^k - q^a}{q - 1} = \theta(k) - \theta(a)$. Similarly, for every such $T_1$ we can select $T_2$ with the condition that $A \oplus T_1 \oplus T_2$ is $(a+2)$-dim subspace of $V$ in $(\theta(k) - \theta(a+1))$ ways. So the number of distinct ordered sets $\{T_1, T_2\}$ is $(\theta(k) - \theta(a))(\theta(k) - \theta(a+1))$. By induction the number of distinct ordered sets $\{T_1, T_2, \cdots, T_b\}$, such that $A \bigoplus_{i=1}^{b} T_i$ is an $(a+b)$-dim subspace of $V$ is $\prod_{i=0}^{b-1}(\theta(k) - \theta(a+i))$. We know that the number of permutations of a $b$-sized set is $b!$. Therefore, the number of distinct (un-ordered) sets satisfying the required conditions is $\frac{1}{b!}\prod_{i=0}^{b-1}(\theta(k) - \theta(a + i))$. This completes the proof. ∎

We also use the following corollary to Lemma 8. The proof of this follows from Lemma 8 (by taking $k = a, b = 1, a = a - 1$ in Lemma 8).

**Corollary 1.** *Consider two subspaces $A', A$ of a $k$-dim vector space $V$ over $\mathbb{F}_q$ such that $dim(A') = a - 1, dim(A) = a$ and $A' \subset A$. The number of distinct $T \in \mathbb{T}$ such that $A' \oplus T = A$ is $\frac{1}{1!}(\theta(a) - \theta(a-1)) = \frac{q^a - q^{a-1}}{q-1} = q^{a-1}$.*

We now give an example of another coded caching scheme, which we shall see in Section V-C to be illustrative of Scheme B of this paper.

*B. An Illustrative Example*

**Example 5.** *Consider a caching system with $K = 7, F = 21, \frac{M}{N} = 0.4285$. Similar to Example 3, we need to provide the indexing for the users $(\mathcal{K})$ and the subfiles $(\mathcal{F})$. For this purpose, we consider some quantities from the projective geometry over finite fields. Consider a 3-dim vector space $(\mathbb{F}_2^3)$. Consider the 1-dim subspaces of $\mathbb{F}_2^3$ which are listed as follows,*

$$T_1 = span\{(0,0,1)\}$$
$$T_2 = span\{(0,1,0)\}$$
$$T_3 = span\{(1,0,0)\}$$
$$T_4 = span\{(1,1,0)\}$$
$$T_5 = span\{(1,0,1)\}$$
$$T_6 = span\{(0,1,1)\}$$
$$T_7 = span\{(1,1,1)\}.$$

*Let $\mathbb{X} = \{T_1, T_2, T_3, T_4, T_5, T_6, T_7\}$. We know that any two distinct 1-dim subspaces are linearly independent. Let $\mathbb{Y}$ be the set of all 2-sized sets of linearly independent 1-dim subspaces of $\mathbb{F}_2^3$ i.e., $\mathbb{Y} = \{\{T_i, T_j\}, i, j \in \{1, 2, \cdots, 7\}, i \neq j\}$. Therefore $|\mathbb{Y}| = 21$.*

*Let $\mathbb{Z}_1$ be the set of all 3-sized sets of linearly dependent 1-dim subspaces of $\mathbb{F}_2^3$. It is easy to see that*

$$\mathbb{Z}_1 = \{\{T_1, T_2, T_6\}, \{T_2, T_3, T_4\}, \{T_1, T_3, T_5\}, \{T_1, T_4, T_7\},$$
$$\{T_2, T_5, T_7\}, \{T_3, T_6, T_7\}, \{T_4, T_5, T_6\}\}.$$

*Let $\mathbb{Z}$ be the set of all 3-sized sets of linearly independent 1-dim subspaces of $\mathbb{F}_2^3$. Therefore $\mathbb{Z} = \binom{\mathbb{X}}{3} \setminus \mathbb{Z}_1$. Therefore $|\mathbb{Z}| = 28$. We now proceed to describe the caching phase and delivery phase.*

**Caching phase:** *Let the user set be $\underline{\mathcal{K} = \mathbb{X}}$ and the set of subfiles be $\underline{\mathcal{F} = \mathbb{Y}}$. During the caching phase, every file $(W_i, i \in [N])$ is divided into $F = |\mathbb{Y}| = 21$ subfiles. The subfiles of $W_i$ are denoted as $W_{i,Y}, \forall Y \in \mathcal{F}$. The caching scheme is,*

- *For each $i \in [N]$, the user $T_l \in \mathcal{K}$ caches the subfile $W_{i,Y}$ if $\{T_l\} \cup Y \notin \mathbb{Z}$.*

*For instance, the user $T_1$ caches the subfiles $W_{i,\{T_1,T_2\}}, W_{i,\{T_1,T_3\}}, W_{i,\{T_1,T_4\}}, W_{i,\{T_1,T_5\}}, W_{i,\{T_1,T_6\}}, W_{i,\{T_1,T_7\}}$ and $W_{i,\{T_2,T_6\}}, W_{i,\{T_3,T_5\}}, W_{i,\{T_4,T_7\}}$ for every $i \in [N]$. It is easy to see that, every user caches 9 subfiles of every file. Hence, the cache fraction $\frac{M}{N} = \frac{9}{21} = 0.4285$.*

**Delivery phase:** *Let demand of an arbitrary user $T_l \in \mathcal{K}$ be $W_{d_{T_l}}$. Note that the subfiles requested by the user $T_l$ are precisely $\{W_{d_{T_l}, Z \setminus T_l} : \forall Z \in \mathbb{Z}\}$. The transmission scheme is,*

- *For each $Z \in \mathbb{Z}$, the server makes the transmission $\bigoplus_{T_l \in Z} W_{d_{T_l}, Z \setminus T_l}$.*

*For instance, the transmission corresponding to $\{T_1, T_2, T_3\} \in \mathbb{Z}$ is $W_{d_{T_1},\{T_2,T_3\}} + W_{d_{T_2},\{T_1,T_3\}} + W_{d_{T_3},\{T_1,T_2\}}$. It is clear that, from this transmission, the user $T_1$ decodes $W_{d_{T_1},\{T_2,T_3\}}$, the user $T_2$ decodes $W_{d_{T_2},\{T_1,T_3\}}$, and the user $T_3$ decodes $W_{d_{T_3},\{T_1,T_2\}}$. Since one transmission is made for each $Z \in \mathbb{Z}$, all the user demands will be satisfied. As 3 users are served in any transmission, the global caching gain is 3. By using (1), the rate of the scheme is $R = \frac{|\mathbb{Z}|}{|\mathcal{F}|} = \frac{28}{21} = 1.33$.*

We are now ready to construct Scheme B, which is the main result of this section.

## C. Construction of Scheme B

We now proceed to develop some notations which are used to label our user vertices, subfile vertices and induced matchings of the bipartite caching graph that we construct.

Let $k, m, n, q$ be positive integers such that $n + m \le k$ and $q$ is some prime power. Consider a $k$-dim vector space $\mathbb{F}_q^k$ and the following sets of subspaces.

$$\mathbb{T} \triangleq PG_q(k-1, 0). \text{ (set of all 1-dim subspaces)}$$

$$\mathbb{R} \triangleq PG_q(k-1, n-1). \text{ (set of all } n\text{-dim subspaces)}$$

$$\mathbb{S} \triangleq PG_q(k-1, m-1). \text{ (set of all } m\text{-dim subspaces)}$$

$$\mathbb{U} \triangleq PG_q(k-1, n+m-1). \text{ (set of all } (n+m)\text{-dim subspaces)}$$

Now, consider the following sets, which are used to present our bipartite caching graph.

$$\mathbb{X} \triangleq \left\{ \{T_1, T_2, \cdots, T_n\} : \forall T_i \in \mathbb{T}, \sum_{i=1}^{n} T_i \in \mathbb{R} \right\}. \tag{7}$$

$$\mathbb{Y} \triangleq \left\{ \{T_1, T_2, \cdots, T_m\} : \forall T_i \in \mathbb{T}, \sum_{i=1}^{m} T_i \in \mathbb{S} \right\}. \tag{8}$$

$$\mathbb{Z} \triangleq \left\{ \{T_1, \cdots, T_{n+m}\} : \forall T_i \in \mathbb{T}, \sum_{i=1}^{n+m} T_i \in \mathbb{U} \right\}. \tag{9}$$

Thus, $\mathbb{X}$ is the set of all $n$-sized sets of 1-dim subspaces such that their sum is an $n$-dim subspace. Intuitively, each $T_i, i \in [n]$ in $\{T_1, T_2, \cdots, T_n\}$ has a different extra dimension. So, the dimension of the subspace $\sum_{i=1}^{n} T_i$ is $n$. Similarly, we have $\mathbb{Y}$ and $\mathbb{Z}$.

We now proceed to construct the bipartite caching graph $B(K, F, D)$. Construct a bipartite graph $B$ with left (user) vertex set $\underline{\mathcal{K}} = \mathbb{X}$ and right (subfile) vertex set $\underline{\mathcal{F}} = \mathbb{Y}$. Define the edge set of $B$ as,

$$E(B) \triangleq \{\{X, Y\} : X \in \mathbb{X}, Y \in \mathbb{Y}, X \cup Y \in \mathbb{Z}\}.$$

We now find the values of $K, F$ and the degree of any user vertex (left degree) $D$.

**Lemma 9.** *From the given construction, we have the following.*

$$K = |\mathbb{X}| = \frac{1}{n!} \, q^{\frac{n(n-1)}{2}} \prod_{i=0}^{n-1} \begin{bmatrix} k-i \\ 1 \end{bmatrix}_q,$$

$$F = |\mathbb{Y}| = \frac{1}{m!} \, q^{\frac{m(m-1)}{2}} \prod_{i=0}^{m-1} \begin{bmatrix} k-i \\ 1 \end{bmatrix}_q,$$

$$D \triangleq |\mathcal{N}(X)| = \frac{q^{nm}}{m!} \, q^{\frac{m(m-1)}{2}} \prod_{i=0}^{m-1} \begin{bmatrix} k-n-i \\ 1 \end{bmatrix}_q.$$

*where the last equation holds for any $X \in \mathbb{X}$.*

*Proof:* See Appendix C. ∎

Note that, by Lemma 9, $B$ is a $D$-left regular bipartite graph with $K$ left vertices and $F$ right vertices. Therefore, by Definition 1, $B(K, F, D)$ is a valid bipartite caching graph.

**Remark 4.** *It is also easy to show that the degree of any right (subfile) vertex $Y \in \mathbb{Y}$ is $|\mathcal{N}(Y)| = \dfrac{q^{nm}}{n!} \, q^{\frac{n(n-1)}{2}} \prod_{i=0}^{n-1} \begin{bmatrix} k-m-i \\ 1 \end{bmatrix}_q$ (the proof is similar to that of $D$ in Appendix C with an interchange of $n$ and $m$). Therefore, $B$ is a bi-regular bipartite caching graph.*

We now show that $B$ has an induced matching cover $\{\mathcal{C}_i : i \in [S]\}$ such that $|\mathcal{C}_i| = g, \forall i \in [S]$, for some $g, S \in \mathbb{Z}^+$.

**Induced matching cover:** We first describe an induced matching of $B$ and show that such equal-sized induced matchings partition $E(B)$. This will suffice to show the delivery scheme as per Theorem 1.

We now present an induced matching of size $\binom{n+m}{n}$ in $B$ (where $\binom{a}{b}$ represents binomial coefficient). Recall the definition of $\mathbb{Z}$ from (9).

**Lemma 10.** *Consider $Z = \{T_1, T_2, \cdots, T_{n+m}\} \in \mathbb{Z}$. Then $\mathcal{C}_Z \triangleq \left\{ \{X, Z \setminus X\} : X \in \binom{Z}{n} \right\} \subseteq E(B)$ is an induced matching of size $\binom{n+m}{m}$ in $B$.*

*Proof:* First note that $\mathcal{C}_Z$ is well-defined as $Z \in \mathbb{Z}$. Consider an arbitrary $X \subset Z$ such that $X \in \mathbb{X}$. It is clear that $\{X, Z \setminus X\} \in E(B)$. Consider two distinct edges $\{X_1, Z \setminus X_1\}, \{X_2, Z \setminus X_2\} \in \mathcal{C}_Z$ (where $X_1 \neq X_2$). It is clear that $Z \setminus X_1 \neq Z \setminus X_2$. WLOG, let $T_a \notin X_1$ and $T_a \in X_2$ (as $X_1 \neq X_2$, such a $T_a$ exists). Then, $T_a \notin X_1 \cup (Z \setminus X_2)$ and hence $X_1 \cup (Z \setminus X_2) \subsetneq Z$. By the definition of $\mathbb{Z}$ in (9), it is clear that $X_1 \cup (Z \setminus X_2) \notin \mathbb{Z}$. Therefore, we have $\{X_1, Z \setminus X_2\} \notin E(B)$, by the definition of $E(B)$. Similarly $\{X_2, Z \setminus X_1\} \notin E(B)$. By invoking Definition 2, it is clear that $\mathcal{C}_Z$ is an induced matching of $B$. It is easy to see that $|\mathcal{C}_Z| = \binom{n+m}{n}$. This completes the proof. ∎

We now show that the set of induced matchings $\{\mathcal{C}_Z : Z \in \mathbb{Z}\}$ partition $E(B)$.

**Lemma 11.** *The induced matchings $\{\mathcal{C}_Z : Z \in \mathbb{Z}\}$ as defined in Lemma 10 partition the edge set $E(B)$.*

*Proof:* Consider $Z, Z' \in \mathbb{Z}$ such that $Z \neq Z'$. By definition of $\mathcal{C}_Z, \mathcal{C}_{Z'}$, we have $\mathcal{C}_Z \cap \mathcal{C}_{Z'} = \phi$. Now consider an arbitrary edge $\left\{ \{T_1, T_2, \cdots, T_n\}, \{T_{n+1}, T_{n+2} \cdots, T_{n+m}\} \right\} \in E(B)$. By the construction of $B$, $\{T_1, T_2, \cdots, T_n\} \cup \{T_{n+1}, T_{n+2} \cdots, T_{n+m}\} \in \mathbb{Z}$. Therefore, the edge $\left\{ \{T_1, T_2, \cdots, T_n\}, \{T_{n+1}, T_{n+2} \cdots, T_{n+m}\} \right\}$ lies in the unique induced matching, $\mathcal{C}_{\{T_1, T_2, \cdots, T_{n+m}\}}$ (defined as in Lemma 10). This completes the proof. ∎

Therefore, $\{\mathcal{C}_Z : Z \in \mathbb{Z}\}$ is an induced matching cover of $B$. We are now ready to present our coded caching scheme using the bipartite caching graph constructed above.

---

**Theorem 3. (Scheme B)** *Let $k, n, m, q$ be positive integers such that $n + m \leq k$ and $q$ be some prime power. The bipartite graph $B$ constructed in Section V-C is a $B(K, F, D)$ bipartite caching graph with an induced matching cover having $S = \frac{1}{(n+m)!} q^{\frac{(n+m)(n+m-1)}{2}} \prod_{i=0}^{n+m-1} \begin{bmatrix} k-i \\ 1 \end{bmatrix}_q$ induced matchings, each having $g = \binom{n+m}{n}$ edges, and defines a coded caching scheme with,*

$$K = \frac{1}{n!} q^{\frac{n(n-1)}{2}} \prod_{i=0}^{n-1} \begin{bmatrix} k-i \\ 1 \end{bmatrix}_q,$$

$$F = \frac{1}{m!} q^{\frac{m(m-1)}{2}} \prod_{i=0}^{m-1} \begin{bmatrix} k-i \\ 1 \end{bmatrix}_q,$$

$$\frac{M}{N} = 1 - q^{nm} \prod_{i=0}^{m-1} \frac{\begin{bmatrix} k-n-i \\ 1 \end{bmatrix}_q}{\begin{bmatrix} k-i \\ 1 \end{bmatrix}_q},$$

$$R = \frac{m! \, q^{nm}}{(n+m)!} q^{\frac{n(n-1)}{2}} \prod_{i=0}^{n-1} \begin{bmatrix} k-m-i \\ 1 \end{bmatrix}_q,$$

*Global caching gain* $\gamma = \binom{n+m}{n}.$

---

*Proof:* From Lemma 9, we get the expressions of $K, F$ and $D$. By Lemma 10 and Lemma 11, the size of the induced matchings of $B$ is $g = |\mathcal{C}_Z| = \binom{n+m}{m}$ for any $Z \in \mathbb{Z}$ and they partition the edge set $E(B)$. Further, by Lemma 10, the number

of induced matchings in the induced matching cover is $S = |\mathbb{Z}|$. Finding $|\mathbb{Z}|$ is same as that of finding $K$ (replace $n$ with $n + m$ in the computation of $K$ in Appendix C). Therefore,

$$S = \frac{1}{(n+m)!} \prod_{i=0}^{n+m-1} (\theta(k) - \theta(i)) \tag{10}$$

$$= \frac{1}{(n+m)!} q^{\frac{(n+m)(n+m-1)}{2}} \prod_{i=0}^{n+m-1} \begin{bmatrix} k-i \\ 1 \end{bmatrix}_q.$$

Hence the bipartite graph $B$ satisfies all the conditions in Theorem 1. Therefore, there exists a coded caching scheme with $K$ users, subpacketization $F$,

$$1 - \frac{M}{N} = \frac{D}{F} = \frac{\frac{1}{m!} \prod_{i=0}^{m-1} (\theta(k) - \theta(n+i))}{\frac{1}{m!} \prod_{i=0}^{m-1} (\theta(k) - \theta(i))} \tag{11}$$

$$= \prod_{i=0}^{m-1} \frac{q^k - q^{n+i}}{q^k - q^i} = \prod_{i=0}^{m-1} q^n \frac{q^{k-n-i} - 1}{q^{k-i} - 1},$$

(where $D, F$ expressions are from Appendix C).

Therefore,

$$\frac{M}{N} = 1 - q^{nm} \prod_{i=0}^{m-1} \frac{\begin{bmatrix} k-n-i \\ 1 \end{bmatrix}_q}{\begin{bmatrix} k-i \\ 1 \end{bmatrix}_q}.$$

$$R = \frac{S}{F} = \frac{\frac{1}{(n+m)!} \prod_{i=0}^{n+m-1} (\theta(k) - \theta(i))}{\frac{1}{m!} \prod_{i=0}^{m-1} (\theta(k) - \theta(i))}$$

$$= \frac{m!}{(n+m)!} \prod_{i=m}^{n+m-1} (\theta(k) - \theta(i))$$

$$= \frac{m!}{(n+m)!} \prod_{i=0}^{n-1} (\theta(k) - \theta(m+i))$$

$$= \frac{m! \, q^{nm}}{(n+m)!} q^{\frac{n(n-1)}{2}} \prod_{i=0}^{n-1} \begin{bmatrix} k-m-i \\ 1 \end{bmatrix}_q.$$

(where $S$ expression is from (10) and $F$ expression is from Appendix C). Finally, we have that the global caching gain $\gamma = g$, which completes the proof. ∎

We regard the coded caching scheme (Scheme B) presented in Theorem 3 as the main scheme of this work. We present Algorithm 1 in which the caching and delivery scheme of Scheme B are captured. The coded caching scheme proposed in Theorem 3, does not exist for all $K$ (similar to most of the low subpacketization coded caching schemes in the literature). So based on the design parameters (desired number of users, cache size) we add some dummy users and treat some fraction of the available cache as unused cache. This is done as follows. For the given number of users $(K')$, cache size $(M')$ and number of files $(N)$, select the appropriate parameters $k, n, m, q$ which give a coded caching scheme according to Theorem 3 with parameters $K, \frac{M}{N}, F, R$ such that $(K - K')$ and $(\frac{M'}{N} - \frac{M}{N})$ are non-negative and as small as possible (we treat the extra users $K - K'$ as dummy users and the extra cache $M' - M$ is left unused). Now construct a bipartite caching graph $B(K, F, D)$ as mentioned in Section V-C and find $\mathbb{X}$ (user indices), $\mathbb{Y}$ (subfile indices), $\mathbb{Z}$ (indices of induced matchings of $B$, equivalently indices of transmissions) by using (7),(8),(9).

---

**Algorithm 1** Coded caching scheme proposed in Theorem 3

---

1: **procedure** PLACEMENT PHASE
2:     **for** each $i \in [N]$ **do**
3:         Split $W_i$ into $\{W_{i,Y} : Y \in \mathbb{Y}\}$.
4:     **end for**
5:     **for** each $X \in \mathbb{X}$ **do**
6:         user $X$ caches the subfiles $W_{i,Y}, \forall i \in [N], \forall Y \in \mathbb{Y}$ such that $X \cup Y \notin \mathbb{Z}$.
7:     **end for**
8: **end procedure**
9: **procedure** DELIVERY PHASE( demand of user $X$ is represented as $W_{d_X}, \forall X \in \mathbb{X}$)
10:     **for** each $Z = \{V_1, V_2, \cdots, V_{n+m}\} \in \mathbb{Z}$ **do**
11:         Server transmits $\bigoplus\limits_{X \in \binom{Z}{n}} W_{d_X, Z \setminus X}$.
12:     **end for**
13: **end procedure**

---

## D. Asymptotic Analysis and Numerical Comparisons of Scheme B with the State of the Art

In Appendix D, we provide the asymptotic analysis for the scheme presented in Theorem 3. When $\frac{M}{N}$ is upper bounded by a constant ($\frac{M}{N} \leq constant$) and as $K$ grows large, from Appendix D we see that our scheme has subexponential subpacketization i.e., $F = q^{O((\log_q K)^2)}$ and rate $R = \Theta\left(\frac{K}{(\log_q K)^n}\right)$. Hence, Scheme B overcomes the drawback of Scheme A (high cache requirement), but with higher rate. The asymptotics of Scheme B obtained in this section are summarized in the third row of Table II in Section II.

We now compare our scheme with some schemes from Table I. We first discuss asymptotic comparison in subpacketization as the number of users $K$ increases. In Table I, we see that the subpacketization levels of several schemes are exponential in $K^{\frac{1}{r}}$ for some positive integer $r$. Comparatively, our scheme achieves subpacketization exponential in $O((\log_q K)^2)$, which is an improvement. Matching our scheme's parameters with those from [15] is hard. A special case of the general scheme in [15] is discussed in that work, which achieves subpacketization exponential in $\sqrt{K}$. Our Scheme B thus improves over this special case. The scheme of [17] achieves linear subpacketization, but it works only for an extremely large number of users and hence we do not compare with this numerically. The PDA scheme $P_1$ from [18] as given in Table I achieves a subpacketization that is smaller than $K$, but uses a large cache fraction in general ($\geq 0.5$). Hence, we do not compare with this scheme also.

We now come to the numerical comparisons. In Table VII, we compare numerically the scheme in Theorem 3 with a modified version of Ali-Niesen scheme with user grouping from [12] (Section V-A in [12], which is also given in row 2 of Table I). The scheme in [12] is parameterized by the cache fraction $\frac{M}{N}$, global caching gain $\gamma$ and number of user groups $l$, and gives a scheme with the number of users $K_2 = K'l$ and subpacketization $F_2 = \binom{K'}{\gamma-1}$, where $K' = (\gamma-1)\lceil\frac{N}{M}\rceil$. The number of users and subpacketization corresponding to Theorem 3, are labelled as $K_1$ and $F_1$ respectively. From the table it is clear that Scheme B performs better than [12], for most of the cases in terms of the subpacketization.

In Table VIII and Table IX, we compare numerically the scheme in Theorem 3 with the schemes in [13] and [16] respectively, for some choices of $K$ and $\frac{M}{N}$. We chose these two schemes for the reason that these two schemes show a large improvement in the subpacketization level without compromising much on the coding gain compared to the basic scheme of [6]. For instance, the PDA based scheme in [13] achieves lower subpacketization (though not in the asymptotic sense) over [6], while having a global caching gain only one less than the scheme of [6].

We label the parameters of our scheme in Theorem 3 as $K_1, \left(\frac{M}{N}\right)_1, F_1, \gamma_1$. The parameters of the scheme presented in [13] are $K_2 = q(m+1), \left(\frac{M}{N}\right)_2 = \frac{1}{q}, F_2 = q^m, \gamma_2 = \frac{K(1-\frac{M}{N})}{q-1} = m+1$ where $q(\geq 2), m \in \mathbb{Z}^+$. The parameters of the scheme presented in [16] are labeled as $K_3 = nq, \left(\frac{M}{N}\right)_3 = \frac{1}{q}, F_3 = q^k z, \gamma_3 = k+1$ where $n, k, q, z$ parameters are defined as per [16]. As it is difficult to match exact parameters, we try to match the $K$ and $\frac{M}{N}$ values as closely as possible between our scheme and these two schemes, while comparing the rate and subpacketization. Also, as the subpacketization can take large values, we approximate it to the nearest positive power of 10.

Throughout, we notice that our Scheme B has parameters for small cache sizes, offering large reductions in subpacketization

TABLE VII: Comparison of the coded caching scheme in Theorem 3 (Scheme B) with the scheme in [12] (Ali-Niesen scheme with grouping). We match the cache fraction, global caching gain, and number of users, as closely as possible, and compare the subpacketization level.

| Number of users | | Cache fraction | Global caching gain | Subpacketization | |
|---|---|---|---|---|---|
| $(k,n,m,q)$ $K_1$ (Theorem 3) | $(K',l)$ $K_2$ [12] | $\frac{M}{N}$ | $\gamma$ | $F_1$ (Theorem 3) | $F_2$ [12] |
| $(7,2,4,2)$ 8001 | $(56,143)$ 8008 | $\frac{125}{381}$ | 15 | $10^7$ | $10^{12}$ |
| $(7,2,2,2)$ 8001 | $(75,107)$ 8025 | $\frac{187}{2667}$ | 6 | 8001 | $10^7$ |
| $(5,2,2,2)$ 465 | $(20,23)$ 460 | $\frac{43}{155}$ | 6 | 465 | 15504 |
| $(4,2,1,2)$ 105 | $(10,10)$ 100 | $\frac{1}{5}$ | 3 | 15 | 45 |
| $(4,1,2,3)$ 40 | $(20,2)$ 40 | $\frac{1}{10}$ | 3 | 780 | 190 |

TABLE VIII: Comparison of the coded caching scheme in Theorem 3 (Scheme B) with the scheme in [13]. (inf represents $> 10^{307}$). (We try to match the $K$ and $\frac{M}{N}$ values of Theorem 3 with that of [13] as closely as possible by choosing appropriate parameters, and compare the subpacketization and gain).

| Number of users | | Cache fraction | | Subpacketization | | Global caching gain | |
|---|---|---|---|---|---|---|---|
| $(k,n,m,q)$ $K_1$ (Theorem 3) | $(q,m)$ $K_2$ [13] | $\left(\frac{M}{N}\right)_1$ (Theorem 3) | $\left(\frac{M}{N}\right)_2$ [13] | $F_1$ (Theorem 3) | $F_2$ [13] | $\gamma_1$ (Theorem 3) | $\gamma_2$ [13] |
| $(4,2,2,2)$ 105 | $(2,51)$ 104 | 0.54 | 0.50 | 105 | $10^{15}$ | 6 | 52 |
| $(5,2,2,2)$ 465 | $(4,116)$ 468 | 0.28 | 0.25 | 465 | $10^{69}$ | 6 | 117 |
| $(4,2,2,3)$ 780 | $(3,259)$ 780 | 0.38 | 0.33 | 780 | $10^{123}$ | 6 | 260 |
| $(7,2,4,2)$ 8001 | $(3,2666)$ 8001 | 0.33 | 0.33 | $10^7$ | inf | 15 | 2667 |
| $(7,2,2,2)$ 8001 | $(14,571)$ 8008 | 0.07 | 0.07 | 8001 | inf | 6 | 572 |

TABLE IX: Comparison of the coded caching scheme in Theorem 3 (Scheme B) with the scheme in [16]. (We try to match the $K$ and $\frac{M}{N}$ values of Theorem 3 with that of [16] as closely as possible by choosing appropriate parameters, and compare the subpacketization and gain).

| Number of users | | Cache fraction | | Subpacketization | | Global caching gain | |
|---|---|---|---|---|---|---|---|
| $(k, n, m, q)$ $K_1$ (Theorem 3) | $(k, n, z, q)$ $K_3$ [16] | $\left(\frac{M}{N}\right)_1$ (Theorem 3) | $\left(\frac{M}{N}\right)_3$ [16] | $F_1$ (Theorem 3) | $F_3$ [16] | $\gamma_1$ (Theorem 3) | $\gamma_3$ [16] |
| $(3, 1, 2, 5)$ 31 | $(7, 12, 2, 3)$ 36 | 0.2 | 0.3 | 465 | 4374 | 3 | 8 |
| $(4, 1, 2, 3)$ 40 | $(8, 12, 3, 5)$ 60 | 0.1 | 0.2 | 780 | $10^6$ | 3 | 9 |
| $(4, 2, 1, 2)$ 105 | $(9, 12, 3, 11)$ 132 | 0.2 | 0.1 | 15 | $10^9$ | 3 | 10 |
| $(5, 2, 1, 2)$ 465 | $(8, 12, 3, 29)$ 348 | 0.09 | 0.03 | 31 | $10^{12}$ | 3 | 9 |

compared to the other two existing schemes in the literature, while having smaller caching gains (and equivalently, higher rate). Since many of the subpacketization values are of the order of $10^2 - 10^4$, we consider our Scheme B to be of importance in practice.

We finally remark that we can also achieve subpacketization $F = K$ (we discuss this in the next subsection) and even $F < K$ by choosing the parameters appropriately. For instance, for $k = 4, n = 2, m = 1, q = 4$ (in Theorem 3), we have $K = 3570, F = 85$, with $\frac{M}{N} = 0.0588$, and gain 3.

### E. Scheme C: A Flexible Scheme with Linear Subpacketization (a special case of Scheme B)

One of the most interesting regimes for the coded caching problem is that of linear subpacketization, i.e., the case when $F = O(K)$. It is known from [14] via the theory of hypergraphs that linear subpacketization is not sufficient for achieving constant rate. A well known result from [17] shows that there exist coded caching schemes which have $F = K$ and achieve a rate of $K^\delta$ (for small $\delta$), and for a small cache fraction. However, the number of users required for this construction is extremely high. The Ali-Niesen scheme [6] itself achieves $F = K$ when $\frac{M}{N} = \frac{1}{K}$, but the rate is then $R = \frac{K-1}{2}$. A similar scheme with the same parameters is known from [18] (see Section V-B in [18]). Thus, most of the existing schemes with linear subpacketization either require an extremely large number of users to exist, or have a very small caching gain.

The following corollary to Theorem 3 gives a new linear subpacketization scheme (Scheme C) obtained using our projective geometry based technique. Note that the scheme is parametrized by the prime power $q$ (finite-field size) and the cache fraction $\lambda$, and hence is flexible for different numbers of users and cache fraction.

**Corollary 2. (Scheme C)** *For $q$ being some prime power and $\lambda \in (0, 1)$ such that $\lambda q$ is a positive integer, then there exists a linear subpacketization coded caching scheme with $F = K \leq \frac{q^{2\lambda^2 q^2}}{(\lambda q)!}$, cache fraction $\frac{M}{N} \leq \lambda$, and global caching gain $\gamma \geq \frac{4^{\lambda q}}{2\sqrt{\lambda q}}$ $\left( \text{with the rate achieved being } \frac{K(1 - \frac{M}{N})}{\gamma} \right)$.*

*Proof:* We choose some specific values for the parameters for our construction in Section V-C to prove this result. From (26) in Appendix D, we have that $\frac{M}{N} \leq \frac{n}{q^{\alpha - n + 1}} \leq \frac{n}{q}$. We choose $n = \lambda q$ (note that $n$ must be an integer and $q$ is a prime power and hence we have our constraints on $\lambda$).

From the expressions in Theorem 3, we have that $K = F$ when $m = n$. We choose the least valid value of $k$, i.e., $k = n + m = 2n$. We thus have with these parameters,

$$K = F \overset{(27)}{\leq} \frac{q^{kn}}{n!} = \frac{q^{2n^2}}{n!} = \frac{q^{2\lambda^2 q^2}}{(\lambda q)!},$$

as stated in the statement of the corollary. We now come to the global caching gain. From Theorem 3, we have that the global caching gain of the scheme is

$$\gamma = \binom{2n}{n} \geq \frac{4^n}{\sqrt{4n}} = \frac{4^{\lambda q}}{2\sqrt{\lambda q}},$$

where the above inequality is a well-known inequality for the middle binomial coefficient that holds for $n \geq 1$. This completes the proof. ∎

## F. Generalization of Scheme B

We can further generalize the scheme presented in Section V-C (Scheme B). In Scheme B, the user vertices, the subfile vertices and the induced matchings are the sets of linearly independent 1-dim subspaces, whereas in the generalized scheme which we present now, these are the sets of linearly independent $l$-dim subspaces.

Let $k, m, n, l, q \in \mathbb{Z}^+$ such that $nl + ml \leq k$ and $q$ is some prime power. Consider a $k$-dim vector space $\mathbb{F}_q^k$ and the following sets of subspaces.

$$\mathbb{L} \triangleq PG_q(k - 1, l - 1). \text{ (set of all } l\text{-dim subspaces)}$$

$$\mathbb{R} \triangleq PG_q(k - 1, nl - 1). \text{ (set of all } nl\text{-dim subspaces)}$$

$$\mathbb{S} \triangleq PG_q(k - 1, ml - 1). \text{ (set of all } ml\text{-dim subspaces)}$$

$$\mathbb{U} \triangleq PG_q(k - 1, (nl + ml) - 1). \text{ (set of all } (nl+ml)\text{-dim spaces)}$$

Similar to Scheme B, consider the following sets, which are used to present the new bipartite caching graph.

$$\mathbb{X} \triangleq \left\{ \{L_1, L_2, \cdots, L_n\} : \forall L_i \in \mathbb{L}, \sum_{i=1}^{n} L_i \in \mathbb{R} \right\}. \tag{12}$$

$$\mathbb{Y} \triangleq \left\{ \{L_1, L_2, \cdots, L_m\} : \forall L_i \in \mathbb{L}, \sum_{i=1}^{m} L_i \in \mathbb{S} \right\}. \tag{13}$$

$$\mathbb{Z} \triangleq \left\{ \{L_1, \cdots, L_{n+m}\} : \forall L_i \in \mathbb{L}, \sum_{i=1}^{n+m} L_i \in \mathbb{U} \right\}. \tag{14}$$

i.e., $\mathbb{X}$ is the set of all $n$-sized sets of linearly independent $l$-dim subspaces, hence their sum (or direct sum) is an $nl$-dim subspace. Similarly, we have $\mathbb{Y}$ and $\mathbb{Z}$.

Construct a bipartite caching graph $B$, similar to that of Scheme B, with the left (user) vertex set $\underline{\mathcal{K}} = \mathbb{X}$ and right (subfile) vertex set $\underline{\mathcal{F}} = \mathbb{Y}$. Define the edge set of $B$ as,

$$E(B) \triangleq \{\{X, Y\} : X \in \mathbb{X}, Y \in \mathbb{Y}, X \cup Y \in \mathbb{Z}\}.$$

The parameters of $B$ such as the number of user vertices $(K)$, the number of subfile vertices $(F)$ and the degree of any user vertex (left degree) $D$ are given in the following lemma.

**Lemma 12.** *From the given construction, we have the following.*

$$K = |\mathbb{X}| = \frac{(l!)^n \; q^{\frac{n(n-1)l^2}{2}}}{(nl)! \; (n!)} \; \frac{\prod_{i=0}^{nl-1} \begin{bmatrix} k-i \\ 1 \end{bmatrix}_q}{\left( \prod_{i=0}^{l-1} \begin{bmatrix} l-i \\ 1 \end{bmatrix}_q \right)^n} \; y_1,$$

$$F = |\mathbb{Y}| = \frac{(l!)^m \; q^{\frac{m(m-1)l^2}{2}}}{(ml)! \; (m!)} \; \frac{\prod_{i=0}^{ml-1} \begin{bmatrix} k-i \\ 1 \end{bmatrix}_q}{\left( \prod_{i=0}^{l-1} \begin{bmatrix} l-i \\ 1 \end{bmatrix}_q \right)^m} \; y_2,$$

$$D = |\mathcal{N}(X)| = \frac{(l!)^m \; q^{ml^2\left(\frac{m-1}{2}+n\right)}}{(ml)! \; (m!)} \; \frac{\prod_{i=0}^{ml-1} \begin{bmatrix} k-nl-i \\ 1 \end{bmatrix}_q}{\left( \prod_{i=0}^{l-1} \begin{bmatrix} l-i \\ 1 \end{bmatrix}_q \right)^m} \; y_2.$$

where the last equation hold for any $X \in \mathbb{X}$, and $y_1 = \prod_{i=0}^{n-1} \binom{(n-i)l}{l}$, $\quad y_2 = \prod_{i=0}^{m-1} \binom{(m-i)l}{l}$.

*Proof:* See Appendix E. ∎

Note that, by Lemma 12, $B$ is a $D$-left regular bipartite graph with $K$ left vertices and $F$ right vertices. Therefore, by Definition 1, $B(K, F, D)$ is a valid bipartite caching graph. We remark that, similar to scheme B, the graph $B(K, F, D)$ is a bi-regular bipartite graph. We now construct the induced matching cover similar to that of Scheme B.

**Induced matching cover:** $\mathcal{C}_Z \triangleq \left\{ \{X, Z \setminus X\} : X \in \binom{Z}{n} \right\}$, where $Z \in \mathbb{Z}$ is a valid induced matching and $\{\mathcal{C}_Z : Z \in \mathbb{Z}\}$ is an induced matching cover of $B$. The proof of these statements is similar to Lemma 10 and Lemma 11. We now present the coded caching scheme corresponding to the bipartite caching graph constructed above.

---

**Theorem 4.** *Let $k, m, n, l, q \in \mathbb{Z}^+$ such that $nl + ml \leq k$ and $q$ is some prime power. The bipartite caching graph $B$ constructed in Section V-F is a $B(K, F, D)$-bipartite caching graph with an induced matching cover having*

$$S = \frac{(l!)^{(n+m)} \; q^{\frac{(n+m)(n+m-1)l^2}{2}}}{((n+m)l)! \; ((n+m)!)} \; \frac{\prod_{i=0}^{(n+m)l-1} \begin{bmatrix} k-i \\ 1 \end{bmatrix}_q}{\left( \prod_{i=0}^{l-1} \begin{bmatrix} l-i \\ 1 \end{bmatrix}_q \right)^n} \times$$

$$\times \prod_{i=0}^{n+m-1} \binom{(n+m-i)l}{l}$$

*induced matchings, each having $g = \binom{n+m}{n}$ edges and defines a coded caching scheme with, number of users $K$ and subpacketization $F$ (where $K, F, D$ are as given in Lemma 12) and*

$$\frac{M}{N} = 1 - q^{mnl^2} \prod_{i=0}^{ml-1} \frac{\begin{bmatrix} k-nl-i \\ 1 \end{bmatrix}_q}{\begin{bmatrix} k-i \\ 1 \end{bmatrix}_q},$$

$$R = \frac{S}{F}, \qquad \text{Global caching gain } \gamma = g.$$

---

*Proof:* Similar to the proof of Theorem 3. ∎

**Remark 5.** *For the special case of $l = 1$, the scheme in Theorem 4, reduces to Scheme B (Theorem 3). The asymptotic analysis for the scheme in Theorem 4 is similar to that of Scheme B (Appendix D). When $\frac{M}{N}$ is upper bounded by a constant $\left( \frac{M}{N} \leq \frac{nl}{q^{k-ml-nl+1}} \right)$ and as $K$ grows large, we can show (by following the similar techniques as that of Appendix D) that the scheme in Theorem 4 achieves subexponential subpacketization i.e., $F = q^{O\left((\log_q K)^2\right)}$ and rate $R = \Theta\left( \frac{K}{(\log_q K)^n} \right)$. Thus, this generalized scheme has more flexibility than Scheme B, but unfortunately does not improve upon Scheme B asymptotically.*

## VI. Extension to Other Settings

In this section, we extend our main scheme, Scheme B, to the distributed computing setting [27], and the wireless interference channel setting [28]. Our main motivation for this section is to present a low subpacketization scheme for each of these two settings which improves upon the currently existing schemes. Most of the existing schemes for these settings are adapted from existing broadcast coded caching schemes, and hence inherit the large subpacketization issue from the same.

The connection from broadcast coded caching schemes to distributed computing schemes was established in [32], through the concept of placement delivery array (PDA). In [32], a method was shown to derive distributed computing schemes from a special class of PDAs known as $g$-PDAs. After recollecting a result from [15], which establishes the connection between bipartite caching graphs and PDA's, we use the result from [32] to adapt our Scheme B to the distributed computing setting. Towards that end, we recall the definition of PDA presented in [13].

**Definition 3** (Placement Delivery Array [13]). *For positive integers $K, F, Z$ and $S$, an $F \times K$ array $\boldsymbol{A} = [a_{j,k}], j \in [F], k \in [K]$, composed of a specific symbol "$*$" and $S$ integers $1, \cdots, S$, is called a $(K, F, Z, S)$ placement delivery array (PDA), if it satisfies the following conditions:*

*C1. The symbol "$*$" appears $Z$ times in each column.*

*C2. Each integer occurs at least once in the array.*

*C3. For any two distinct entries $a_{j_1,k_1}$ and $a_{j_2,k_2}$ we have $a_{j_1,k_1} = a_{j_2,k_2} = s$, an integer, only if*

　　*1. $j_1 \neq j_2, k_1 \neq k_2$, i.e., they lie in distinct rows and distinct columns; and*

　　*2. $a_{j_1,k_2} = a_{j_2,k_1} = *$.*

*If each integer $s \in [S]$ occurs exactly $g$ times, then $\boldsymbol{A}$ is called a regular $g - (K, F, Z, S)$ PDA, or $g$-PDA for short.*

Most known coded caching schemes in the literature correspond to PDAs. The following result shows that any bipartite caching graph $B(K, F, D)$ with an induced matching cover with the size of each induced matching being at least 2 is equivalent to a $g$-PDA. This is equivalent to the result in [15] which relates a PDA with a strong edge coloring of the bipartite caching graph $B$, which as we have discussed in Section III-B is equivalent to an induced matching cover of $B$.

**Theorem 5** (equivalent to Theorem 1 in [15]). *$B$ is a $(K, F, D)$ bipartite caching graph with an induced matching cover consisting of $S$ induced matchings, each of size $g \geq 2$ if and only if there exist a $g - (K, F, Z = F - D, S)$ regular PDA.*

By Theorem 5, it is clear that the bipartite caching graph developed in Section V-C, which resulted in Scheme B, corresponds to a $g - (K, F, F - D, S)$ regular PDA, where the expressions of $K, F, S, g$ are given in Theorem 3 and the expression for $D$ is given in Lemma 9.

### A. Extension to Distributed Computing Systems

Using coded transmissions to reduce the communication load is a general technique that can potentially be used in any distributed communication system. One such system is the distributed computing framework called the Map-Reduce framework [33]. In this framework, one or more functions have to be evaluated on some given large amount of data. To do this in a distributed manner, subsets of the data are assigned to a number of distributed computing nodes, each of which compute the functions on their own data subset. After this, the intermediate values are accumulated at the nodes to give the total computed function values. In [27], for this framework presented in [33], a technique to reduce the communication load was presented, which is inspired from the coded caching framework.

We now recall the model from [27]. Consider $N^{\mathcal{C}}$ files, $K^{\mathcal{C}}$ computing nodes and each node is assigned one or more functions to be computed finally (we use superscript $\mathcal{C}$ to represent the parameters of the distributing computing system). The total number of functions to be computed on the files is $Q$. In the *map* phase, the set of $N^{\mathcal{C}}$ files are split into $F^{\mathcal{C}}$ batches, each containing $\frac{N^{\mathcal{C}}}{F^{\mathcal{C}}}$ files. Each node stores some batches of files and computes all the $Q$ functions on the files present in each batch. We denote the files present at node $k \in [K^{\mathcal{C}}]$ as $\mathcal{M}_k$. These computed function values are referred as intermediate values, and we assume they are represented as $T$-length bit-vectors. Thus, there are $Q|\mathcal{M}_k|$ intermediate values computed at each node $k$ at the end of the map phase. In the *shuffle* phase, each node $k \in [K^{\mathcal{C}}]$ computes a signal denoted by $X_k$ (of

length $l_k$ bits) from the intermediate values, it computed in the map phase, and communicates to other nodes. This is known as *data shuffling*. In the final *reduce* phase, each node is assigned to reduce a set of $\frac{Q}{K^\mathcal{C}}$ functions (we assume $\frac{Q}{K^\mathcal{C}}$ is an integer). Using the signals received from the other nodes via data shuffling, and the intermediate values computed at itself during the map phase, each node decodes all the intermediate values of its assigned output functions, to finally compute the complete value of the output functions.

The *computation load*, denoted by $r^\mathcal{C}$, is the number of map functions computed at all the nodes, normalized by $N$. The communication load, denoted by $L^\mathcal{C}$, is the ratio of the total number of bits transmitted to the total number of bits in all the intermediary values $QNT$. Therefore,

$$r^\mathcal{C} = \frac{\sum\limits_{k \in [K^\mathcal{C}]} |\mathcal{M}_k|}{N}, \qquad L^\mathcal{C} = \frac{\sum\limits_{i=1}^{K^\mathcal{C}} l_i}{QNT}.$$

Similar to the coded caching problem, it is a practical necessity to obtain distributed computing schemes for the above model with smaller number of batches $F^\mathcal{C}$ for low computation and communication loads. We adapt our Scheme B to this setting in order to obtain a distributed computing scheme which has lower number of batches for similar computation load requirements compared to other major schemes in the literature, at the cost of having larger communication loads. For this purpose, we rely on the literature, which relates the PDAs to the distributed computing schemes [32], [34], [35]. In particular, we are interested in the following result from [32].

**Lemma 13.** *(Corollary 3 in [32]) . For a given $g - (K, F, Z, S)$ regular PDA with $g \geq 2$, there exists a scheme for distributed computing system with $K^\mathcal{C} = K$ nodes, achieving the computation load $r^\mathcal{C} = \frac{KZ}{F}$ and communication load $L^\mathcal{C} = \frac{g}{g-1}\frac{S}{KF}$, which can be implemented with the minimum number of batches requirement $F^\mathcal{C} = F$.*

Now, by applying Lemma 13, we can get the corresponding distributed computing system which is presented in Theorem 6 (the proof follows from Theorem 5 and Lemma 13).

---

**Theorem 6.** *Let $k, n, m, q$ be positive integers such that $n + m \leq k$ and $q$ be some prime power. The bipartite caching graph given in Section V-C corresponds to a distributed computing scheme with,*

$$K^\mathcal{C} = \frac{1}{n!} \, q^{\frac{n(n-1)}{2}} \prod_{i=0}^{n-1} \begin{bmatrix} k - i \\ 1 \end{bmatrix}_q,$$

$$r^\mathcal{C} = K^\mathcal{C} \left( 1 - q^{nm} \prod_{i=0}^{m-1} \frac{\begin{bmatrix} k - n - i \\ 1 \end{bmatrix}_q}{\begin{bmatrix} k - i \\ 1 \end{bmatrix}_q} \right),$$

$$F^\mathcal{C} = \frac{1}{m!} \, q^{\frac{m(m-1)}{2}} \prod_{i=0}^{m-1} \begin{bmatrix} k - i \\ 1 \end{bmatrix}_q,$$

$$L^\mathcal{C} = \frac{q^{nm}}{\binom{n+m}{n} - 1} \prod_{i=0}^{n-1} \frac{\begin{bmatrix} k - m - i \\ 1 \end{bmatrix}_q}{\begin{bmatrix} k - i \\ 1 \end{bmatrix}_q}.$$

---

In Table X, we compare our distributed computing system presented in Theorem 6 with that of [27]. The scheme in [27] is the original scheme in the distributed computing literature for the given setting, and can be viewed as an adaptation of the fundamental coded caching scheme [6] into a distributed computing scheme. This scheme is optimal in terms of the communication load for a given computation load, but uses a large batch size, as shown in Table X. The parameters of the

scheme presented in [27] are $K_2^{\mathcal{C}} = K, F_2^{\mathcal{C}} = \binom{K}{r}, r_2^{\mathcal{C}} = r, L_2^{\mathcal{C}} = \frac{1}{r}\left(1 - \frac{r}{K}\right)$. where $K$ is a positive integer and $r \in [K]$. The parameters corresponding to Theorem 6 are labelled as $K_1^{\mathcal{C}}, F_1^{\mathcal{C}}, r_1^{\mathcal{C}}, L_1^{\mathcal{C}}$. We choose matching values for $K^{\mathcal{C}}$ and $r^{\mathcal{C}}$ for the purpose of comparison.

TABLE X: Comparison of the distributed computing scheme in Theorem 6 with the scheme in [27]. (inf represents $> 10^{307}$).

| Number of computing nodes | | Computation load | | Number of batches | | Communication load | |
|---|---|---|---|---|---|---|---|
| $(k, n, m, q)$ $K_1^{\mathcal{C}}$ (Theorem 6) | $(K, r)$ $K_2^{\mathcal{C}}$ [27] | $r_1^{\mathcal{C}}$ (Theorem 6) | $r_2^{\mathcal{C}}$ [27] | $F_1^{\mathcal{C}}$ (Theorem 6) | $F_2^{\mathcal{C}}$ [27] | $L_1^{\mathcal{C}}$ (Theorem 6) | $L_2^{\mathcal{C}}$ [27] |
| $(6, 2, 2, 2)$ 1953 | $(1953, 273)$ 1953 | 273 | 273 | 1953 | inf | 0.1720 | 0.0032 |
| $(4, 2, 1, 3)$ 780 | $(780, 78)$ 780 | 78 | 78 | 40 | $10^{108}$ | 0.45 | 0.0115 |
| $(5, 2, 2, 2)$ 465 | $(465, 129)$ 465 | 129 | 129 | 465 | $10^{117}$ | 0.1445 | 0.0056 |
| $(5, 1, 2, 3)$ 121 | $(121, 4)$ 121 | 4 | 4 | 7260 | $10^{6}$ | 0.4835 | 0.2417 |
| $(4, 2, 1, 2)$ 105 | $(105, 21)$ 105 | 21 | 21 | 15 | $10^{21}$ | 0.40 | 0.0381 |
| $(6, 1, 2, 2)$ 63 | $(63, 3)$ 63 | 3 | 3 | 1953 | 39711 | 0.4762 | 0.3175 |

We see from Table X that for a given number of computing nodes and computation load, our scheme presented in Theorem 6 performs much better than the scheme of [27] in terms of the number of file batches required (which makes our scheme more practical for distributed systems with tens of nodes or more) for the same computation load. However, the communication load is higher in general. However, for some examples, for instance the entries in Table X with number of nodes equal to 63 or 121, for the communication load at most twice of what the scheme of [27] achieves, we achieve order-of-magnitude gains in the batch-number $F^{\mathcal{C}}$.

**Remark 6.** *The coded caching technique has also been utilized in the Device-to-Device (D2D) setting in recent work, for instance, [26]. The underlying communication model is very similar to the distributed computing model. Hence, similar to Theorem 6, we can also extend the proposed Scheme B to the D2D coded caching system [26], using Theorem 5 and [32] (Corollary 1 in [32]) to get low subpacketization schemes for D2D systems. For more details, the reader is referred to [4].*

### B. Extension to Coded Caching in a Wireless Interference Channel

We now present another setting in which we can apply our low subpacketization scheme, namely the interference channel. Coded caching for the interference channel setting with caches at both transmitters and receivers was considered in [28], where the fundamental coded caching scheme of Ali-Niesen [6] was adapted to this setting. In [36], this scheme was further refined to present a scheme with lower subpacketization requirement (particularly, no subpacketization at the transmitter end). We leave the details of these schemes to the reader and refer them to the respective papers. In this subsection, we present the adaptation of our low-subpacketization Scheme B for the setting presented in [28], also motivated by the techniques from [36]. We then compare via numerical examples our scheme's performance with that of [36] which itself is an improvement over [28] in terms of the subpacketization, with all other parameters being the same.

We first review the model given in [28] as shown in Fig. 4. Consider a wireless channel with $K_T$ transmitters and $K_R$ receivers with $N$ files (denoted as $W_i : i \in [N]$) such that each transmitter can store $M_T$ files and each receiver can store
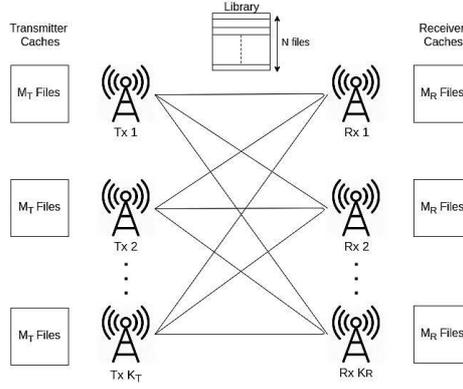
Fig. 4: Channel model for coded caching in the wireless interference channel.

$M_R$ files. We assume $L \triangleq \frac{K_T M_T}{N}$ is an integer, and also that $L$ divides $K_R$. Each transmitter and receiver has one antenna, and the channel coefficient between any particular transmitter-to-receiver pair is a complex number that is a realization of an independent continuous probability distribution, and assumed to be constant over the time of communication. The coded caching scheme in this setting involves two phases as before, the caching phase where the transmitter and receiver caches are populated, and then the delivery phase in which the receiver demands (each receiver demands a particular file as before) are served collaboratively by the transmitters. In every round of transmission, some subset of $K_T$ transmitters send signals to the receivers. In a valid scheme, at the end of a finite number of rounds, the decoding of demands at the respective receivers should be complete. The sum-DoF in this setting refers intuitively to the number of receivers served per transmission. For a precise definition of this parameter, we refer to the reader to [28].

For this setting, under the condition that $L + \frac{K_R M_R}{N} \leq K_R$, the scheme in [28] achieves a sum-DoF of $L + \frac{K_R M_R}{N}$, with a subpacketization level of $\binom{K_T}{L}\binom{K_R}{\frac{K_R M_R}{N}}$ which arises as a result of subpacketization for caching at both the transmitter-side and the receiver-side. The work [36] achieves the same sum-DoF with a subpacketization level of $\binom{K_R/L}{K_R M_R/(LN)}$, which avoids subpacketization at the transmitter-side completely. For more details, we refer the reader to [28] and [36]. In the rest of this subsection, we shall adapt our low subpacketization scheme, Scheme B, to this interference channel setting. For the sake of notational convenience, we focus on adapting a particular special class of Scheme B (which has parameters $(k, m, n, q)$) with parameter $n = 1$. However, the more general scheme can also be adapted, which gives more flexibility in terms of the parameters.

The principle we use is similar to the grouping scheme for the multi-transmitter coded caching scheme given in [36]. The idea behind this is to divide the receivers into a number of $K \triangleq \frac{K_R}{L}$ groups, each containing $L$ receivers. Then the broadcast-channel coded caching scheme (in our case, Scheme B) is applied to the groups considering them to be super-users, while also utilizing the presence of transmitter caches to zero-force some non-demanded and un-cached subfiles. We now give the detailed scheme.

Let $k, m, q$ be positive integers such that $m \leq k$ and $q$ is some prime power. We assume that $K_R = KL$, where $K = \begin{bmatrix} k \\ 1 \end{bmatrix}_q$ for some $k$. Let $\mathbb{V} \triangleq \{V \in PG_q(k-1, 0)\}$ (Note that in Section V-C, we used the notation $\mathbb{T}$ for the same). Consider $\frac{M_R}{N} = 1 - \frac{q^m \begin{bmatrix} k-m \\ 1 \end{bmatrix}_q}{\begin{bmatrix} k \\ 1 \end{bmatrix}_q}$ (note that this is equal to the $M/N$ value we get from Scheme $B$ with $n = 1$ in Theorem 3). We group the $K_R$ users into $K$ groups, such that in each group there are $L$ users. The groups are indexed by distinct 1-dim subspaces of $\mathbb{F}_q^k$, i.e., by the elements of $\mathbb{V}$. Let the users in a group $V \in \mathbb{V}$ be denoted as $V(1), ..., V(L)$. We now describe the placement and delivery phase.

*1) Placement Phase:* *Transmitter Caching Strategy:* We follow the caching strategy described in [36] (Appendix A in [36]) for the transmitter side. We require each subfile of each file to be cached at precisely $L$ transmitters.

- For $i \in [K_T]$, the cache content $C_{T_i}$ of the $i^{th}$ transmitter denoted by $T_i$ is given by

$$C_{T_i} = \{W_{1+(p-1)(mod\ N)} : p \in \{1 + (i-1)M_T, ..., M_T i\}\}.$$

Thus, the transmitters cache successive $M_T$ files into their caches. Because of this caching strategy, it should also be clear that each file cached at $L = \frac{K_T M_T}{N}$ transmitters. This property of the transmitter-side caching will be used to obtain a sum-DoF $= L + \frac{K_R M_R}{N}$ via the zero-forcing technique.

_Receiver Caching Strategy:_ On the receiver side, the caching strategy is based on the projective geometry as described in Section V-C (Scheme B, with the parameter $n = 1$). Consider $m \in \mathbb{Z}^+$ such that $m \leq k$. Let

$$\mathbb{S} \triangleq \{S \in PG_q(k-1, m-1)\}.$$

$$\mathbb{Y} \triangleq \left\{ \{V_1, V_2, ..., V_m\} : \forall V_i \in \mathbb{V}, \sum_{i=1}^{m} V_i \in \mathbb{S} \right\}.$$

Following Scheme B, the subfiles are indexed by the elements of $\mathbb{Y}$. The subfiles of file $W_i$ are denoted as $W_i = \{W_i^Y : Y \in \mathbb{Y}\}$. The subpacketization $F$ is thus equal to $|\mathbb{Y}|$, which was shown in Lemma 9 in Section V to be $\frac{1}{m!} q^{\frac{m(m-1)}{2}} \prod_{i=0}^{m-1} \begin{bmatrix} k-i \\ 1 \end{bmatrix}_q$.

- The caches of the receivers $V(r) : r \in [L]$ in the group $V \in \mathbb{V}$ are populated with the same content, given as

$$Z_V \triangleq \left\{ W_i^Y : Y \in \mathbb{Y}, V \not\subseteq \sum_{V_j \in Y} V_j \right\}_{i=1}^{N}.$$

Thus we have the $K = \begin{bmatrix} k \\ 1 \end{bmatrix}_q$ groups of our current setting in the place of $K$ users in the original Scheme B. We now describe the delivery phase corresponding to this caching strategy.

_2) **Delivery Phase:**_ In the delivery phase, each receiver demands for a file. We use a projective geometry based delivery scheme developed in Section V-C over the groups of users. Let the demand of receiver $V(i)$ be denoted as $W_{d_{V(i)}}$. The delivery scheme serves $m + 1$ groups of receivers in each transmission. For this purpose, we develop some notations following the description of Scheme B.

Let

$$\mathbb{Z} = \left\{ \{V_1, V_2, ..., V_{m+1}\} \subseteq \mathbb{V} : \sum_{j=1}^{m+1} V_j \in PG_q(k-1, m) \right\}.$$

Consider $Z = \{V_1, V_2, ..., V_{m+1}\} \in \mathbb{Z}$, denoting a set of $(m+1)$ groups of receivers, the $(m+1)L$ receivers which will be served in one round of transmission. Let $Y_j = Z \backslash V_j : j \in [m+1]$. In this round, the subfiles to be transmitted to the receivers in the groups given in $Z$ are $\{W_{d_{V_j(i)}}^{Y_j} : j \in [m+1], i \in [L]\}$. Clearly, the subfile $W_{d_{V_j(i)}}^{Y_j}$ which is desired at the user $V_j(i)$, is not cached at any receiver in the group $V_j$ but available at all the receivers in the groups $V_{j'} : j' \in [m+1] \backslash j$. For the purpose of transmitting on a wireless channel, let $\widetilde{W}_{d_{V_j(i)}}^{Y_j}$ denote the signal (from some complex constellation) denoting the mapping of the subfile $W_{d_{V_j(i)}}^{Y_j}$; this mapping is known to all receivers and transmitters. We now construct the idea behind the delivery scheme, by showing the round corresponding to $Z \in \mathbb{Z}$.

We denote by $s_i \in \mathbb{C}^{1 \times 1} : i \in [K_T]$ the signal transmitted by $i^{th}$ transmitter during this round. We have to design our transmission signals such that the mapped subfile $\widetilde{W}_{d_{V_j(i)}}^{Y_j}$ can be obtained at the user $V_j(i)$, while (a) the same subfile can be zero-forced at the receivers in $V_j \backslash V_j(i)$ by utilizing the presence of the $L$ transmitters in which the subfile $W_{d_{V_j(i)}}^{Y_j}$ is available, and (b) the same subfile can be cancelled using the cache content at each receiver in the group $V_{j'} : j' \in [m+1] \backslash j$.

Let $\overline{w}_j = [\widetilde{W}_{d_{V_j(1)}}^{Y_j}, ..., \widetilde{W}_{d_{V_j(L)}}^{Y_j}]^T$ denote the mapped subfiles involved in this round of transmissions corresponding to the group $V_j$. The round of transmissions is described as follows:

$$\begin{bmatrix} s_1 \\ \vdots \\ s_{K_T} \end{bmatrix} = \begin{bmatrix} \overline{A}_1 & \cdots & \overline{A}_{m+1} \end{bmatrix} \begin{bmatrix} \overline{w}_1 \\ \vdots \\ \overline{w}_{m+1} \end{bmatrix}, \tag{15}$$

where $\overline{A}_i \in \mathbb{C}^{K_T \times L} : i \in [m+1]$ contains as its columns the $L$ precoding vectors of length $K_T$ corresponding to the mapped subfiles $\widetilde{W}^{Y_j}_{d_{V_j(i)}} : i \in [L]$ generated by the $K_T$ transmitters. Note that any such precoding vector can have only $L$ non-zeros, as the subfile $W^{Y_j}_{d_{V_j(i)}}$ is available only at some $L$ transmitters. Thus, every column of the matrices $\overline{A}_i : i \in [m+1]$ contains $K_T - L$ zeros, with the remaining entries to be chosen to satisfy the zero-forcing requirements. In the remainder of this subsection, we show that such precoding vectors can indeed be chosen, and show that decoding of the mapped subfiles (and thereby, via the inverse mapping, the uncached subfiles) can be decoded at the respective receivers. The delivery scheme can therefore be completed by constructing transmissions as in (15) for every $Z \in \mathbb{Z}$, with appropriately chosen precoding vectors to effect successful decoding.

We now show that the precoding vectors of (15) in the matrices $\overline{A}_i : i \in [m+1]$ can be chosen to ensure decoding of the desired subfiles at the respective users. The received signals at the receivers of groups $Z = (V_1, ..., V_{m+1})$ is given as,

$$\begin{bmatrix} y_1 \\ \vdots \\ y_{m+1} \end{bmatrix} = \begin{bmatrix} H_1 \\ \vdots \\ H_{m+1} \end{bmatrix} \begin{bmatrix} s_1 \\ \vdots \\ s_{K_T} \end{bmatrix} + z, \tag{16}$$

where $y_j = [y_{j,1}, ..., y_{j,L}]^T : j \in [m+1]$, with $y_{j,i} : i \in [L]$ being the received signal at the user $V_j(i)$, the matrix $H_j : j \in [m+1]$ is the $L \times K_T$ channel matrix from the $K_T$ transmitters to the $L$ receivers of the group $V_j$ and $z = (z_1, \ldots, z_{m+1})^T$ denotes the additive white Gaussian noise, each component of which is distributed as a circular symmetric complex Gaussian.

We now have from (15) and (16), the received signal at the group $V_j$ as

$$y_j = \begin{bmatrix} H_j A_1 \ldots H_j A_{m+1} \end{bmatrix} \begin{bmatrix} \overline{w}_1 \\ \vdots \\ \overline{w}_{m+1} \end{bmatrix} + z_j \tag{17}$$

Because the receiver $V_j(i)$ has the subfiles in (17) except $W^{Z \backslash V_j}_{d_{V_j(i)}} : i \in [L]$, from (17), the receivers of the group $V_j$ can obtain,

$$\tilde{y}_j = H_j A_j \overline{w}_j + z_j.$$

Because of the fact that the entries of $H_j$ are picked from a continuous distribution, any $L$ columns of $H_j$ are linearly independent almost surely, there exists a vector, with the condition that it is non-zero only in the positions corresponding to the transmitters caching $W_{d_{V_j(i)}}$ and $0$ everywhere else, which can be fixed as the $i^{th}$ column of $A_j$, such that $H_j A_j = I$ (the identity matrix of size $L$). Hence, we have found the desired solution for the precoding vectors. This ensures decoding (as SNR grows large, in the DoF sense) of the mapped subfiles $\{\widetilde{W}^{Y_j}_{d_{V_j(i)}} : j \in [m+1], i \in [L]\}$ at the respective receivers, and thus by inverse mapping the original desired subfiles can be decoded. The delivery scheme, which constructs the transmissions as in (15) for every $Z \in \mathbb{Z}$ thus ensures decoding of all uncached and desired subfiles at the respective receivers. As in each round of transmissions, the number of users served is $L(m+1)$, which is the sum-DoF for the presented scheme. The results of the coded caching scheme constructed above is presented in the following theorem.

**Theorem 7.** *In a wireless channel consisting of $N$ files, $K_T$ transmitters each with cache that can contain $M_T$ files such that $L = \frac{K_T M_T}{N}$ is a positive integer, and consisting of $K_R = L \begin{bmatrix} k \\ 1 \end{bmatrix}_q$ receivers each of cache size $M_R$ files where $\frac{M_R}{N} = 1 - \dfrac{q^m \begin{bmatrix} k-m \\ 1 \end{bmatrix}_q}{\begin{bmatrix} k \\ 1 \end{bmatrix}_q}$, we can achieve with high SNR, a sum-$\mathsf{DoF} = L(m+1)$ with subpacketization $F = \frac{1}{m!} \, q^{\frac{m(m-1)}{2}} \prod_{i=0}^{m-1} \begin{bmatrix} k-i \\ 1 \end{bmatrix}_q$, where $k, m \in \mathbb{Z}^+$ with $m \leq k$ and $q$ is some prime power.*

TABLE XI: Comparison of the coded caching scheme for the interference channel presented in Theorem 7 with the scheme in [36], for some specific values of $K_R, L, \frac{M_R}{N}$.

| $(\mathbf{k, m, q})$ $\mathbf{K_R}$ | $\mathbf{L}$ | $\frac{\mathbf{M_R}}{\mathbf{N}}$ | $\mathbf{F_1}$ (Theorem 7) | $\mathbf{F_2}$ [36] | sum-$\mathsf{DoF_1}$ (Theorem 7) | sum-$\mathsf{DoF_2}$ [36] |
|---|---|---|---|---|---|---|
| $(4,3,2)$ 30 | 2 | 0.4667 | 420 | 6435 | 8 | 16 |
| $(5,3,2)$ 62 | 2 | 0.2258 | 4340 | $7.3 \times 10^5$ | 8 | 15 |
| $(4,3,3)$ 80 | 2 | 0.3250 | 9360 | $1.2 \times 10^{10}$ | 8 | 28 |
| $(5,4,2)$ 124 | 4 | 0.4839 | $2.6 \times 10^4$ | $3 \times 10^8$ | 20 | 64 |
| $(6,4,2)$ 252 | 4 | 0.2381 | $5.4 \times 10^5$ | $1.22 \times 10^{14}$ | 20 | 64 |

Table XI gives a numerical comparison of our scheme's parameters, subpacketization and sum-DoF, with that of [36] by choosing matching values for $K_R, L$ and $\frac{M_R}{N}$. The comparison parameters corresponding to Theorem 7 are labelled as $F_1$ and sum-$\mathsf{DoF_1}$. The comparison parameters corresponding to the scheme in [36] are $F_2 = \binom{K_R/L}{\lfloor (K_R M_R)/(NL) \rfloor}$ and sum-$\mathsf{DoF_2} = L + \lfloor \frac{K_R M_R}{N} \rfloor$. As can be seen, our scheme performs better in terms of the subpacketization, while having lesser sum-DoF. The quantities are rounded off to a few decimal places wherever applicable.

## VII. LOWER BOUNDS ON RATE OF DELIVERY SCHEME FOR SYMMETRIC CACHING

In this section, we present two information theoretic lower bounds on the rate of the transmission scheme associated with a $(K, F, D = F(1 - \frac{M}{N}))$ bipartite caching scheme (which is associated with the bipartite caching graph $B(K, F, D)$) and numerically compare with the existing lower bounds from [11], [18]. We obtain two bounds for the rate of coded caching with the fixed parameters $K, F, M$ and $N$. The first bound holds for all symmetric caching schemes, but is quite loose. The second one holds for a special class of symmetric caching schemes, in which each subfile is stored in the same number of users (equivalently the bipartite caching scheme associated with a bi-regular bipartite graph). These bi-regular schemes include all the symmetric caching schemes in the literature to the best of our knowledge, including those in Table I. In this special class of caching schemes, the second bound is shown to be better, for a number of parameter choices, than the existing bounds from the prior work via numerical comparisons.

We first give some preliminary ideas and definitions before we present our bounds. As the subfile $W_{i,f}$ of the file $W_i$ takes values from the finite set $\mathcal{A}$ with uniform distribution, taking the base of logarithm as $|\mathcal{A}|$, we have the Shannon entropy of $W_{i,f}$ as $H(W_{i,f}) = 1, \forall i, f$. Thus $H(W_i) = F, \forall i \in [N]$.

**Definition 4.** *For the given parameters $K, \frac{M}{N}$, and $F$ such that $\frac{FM}{N} \in \mathbb{Z}^+$, a rate $R$ is said to be achievable if there exists some $(K, F, D = F(1 - \frac{M}{N}))$ bipartite caching scheme, with a delivery scheme with the rate $R$ that satisfies all the client demands. For the given parameters $K, F, D$, we define the optimal rate $R^*(K, F, D)$ as follows:*

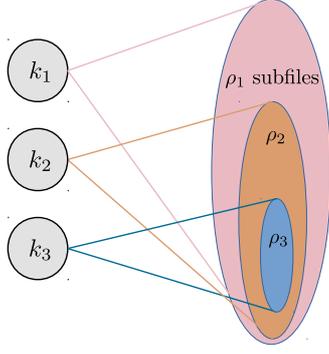$$R^*(K, F, D) \triangleq \inf\{R : R \text{ is achievable}\}.$$

Fig. 5: Intuition behind the lower bound presented in Theorem 8. The nesting of the $\rho_3$ neighbours of the vertex $k_3$ within the $\rho_2$ neighbours of $k_2$, and these further within the $\rho_1$ neighbours of $k_1$, results in the number of transmissions $R^*F \geq \rho_1 + \rho_2 + \rho_3$.

**Remark 7.** *We abbreviate $R^*(K, F, D)$ as simply $R^*$, as the parameters involved will be clear from the context.*

It is known from [11] that $R^* \geq \frac{K(1-\frac{M}{N})}{1+\frac{MK}{N}}$ for any value of $F$, and this is achieved by the scheme in [6] with $F = \binom{K}{\frac{MK}{N}}$. Further, for the coded caching schemes with given parameters derived from the PDAs (which correspond to the bipartite caching schemes with the induced matching based delivery schemes), it was shown in [18] that

$$R^*F \geq \left\lceil \frac{DK}{F} \right\rceil + \left\lceil \frac{D-1}{F-1}\left\lceil \frac{DK}{F} \right\rceil \right\rceil + \cdots$$
$$\cdots + \left\lceil \frac{1}{\frac{FM}{N}+1}\left\lceil \frac{2}{\frac{FM}{N}+2}\left\lceil \cdots \left\lceil \frac{D-1}{F-1}\left\lceil \frac{DK}{F} \right\rceil \right\rceil \cdots \right\rceil \right\rceil \right\rceil. \tag{18}$$

Note that in (18), the notation $R^*$ is abused to correspond to the optimal rate only among the PDA based delivery schemes (which corresponds to the induced matching based delivery schemes), which is a restricted class of delivery schemes among all the delivery schemes obtainable for any symmetric caching scheme. We compare our new lower bounds with these two bounds.

We will first prove a generic lower bound on the rate, in Theorem 8, for the coded caching schemes whose caching scheme is based on the bipartite caching model, about which we briefly give some intuition. Consider the induced-matching based delivery scheme in Section III-B of the bipartite caching graph. The cardinality of any subset of edges of the bipartite caching graph such that no two among the subset can appear in any single induced matching, gives us a lower bound on the number of transmissions in any induced-matching based delivery scheme. This is illustrated in Fig. 5. Consider a bipartite caching graph $B$, of which a subgraph is shown in Fig. 5. Let $k_1$ be an arbitrary user vertex. Let $\rho_1 = |\mathcal{N}(k_1)|$. By Definition 2, it is clear that no two of these $\rho_1$ edges (from $k_1$ to $\mathcal{N}(k_1)$) must lie in the same induced matching. Now consider an arbitrary user vertex $k_2 (\neq k_1)$. Let $\rho_2 = |\mathcal{N}(k_1) \cap \mathcal{N}(k_2)|$. We can see that no two of the $\rho_1 + \rho_2$ edges (from $k_1$ to $\mathcal{N}(k_1)$ and from $k_2$ to $\mathcal{N}(k_1) \cap \mathcal{N}(k_2)$) lie in the same induced matching, by Definition 2. Following this idea, consider $N'$ user vertices $\{k_1, k_2, \cdots, k_{N'}\}$. Let $\rho_j = \left| \bigcap_{i=1}^{j} \mathcal{N}(k_i) \right|$ where $j \in [N']$. Now consider the set of edges of $B$ from $k_j$ to $\bigcap_{i=1}^{j} \mathcal{N}(k_i), \forall j \in [N']$. There are exactly $\sum_{j=1}^{N'} \rho_j$ edges in this set and no two edges lie in the same induced matching, by Definition 2. Therefore, $\sum_{j=1}^{N'} \rho_j$ is a lower bound on the number of transmissions of any induced-matching based delivery scheme.

Interestingly, we show in Theorem 8 that the bound $\sum_{j=1}^{N'} \rho_j$ derived in the above discussion applies in the information theoretic sense also, for a given bipartite caching scheme. We then use this result to show our two lower bounds on the rate of the coded caching schemes with the fixed parameters $K, F$ and $D = F(1 - M/N)$ in Corollary 3 and Theorem 9. Thus, these bounds hold for non-linear schemes as well, in contrast with (18), which was shown to be true for induced matching based delivery schemes. The lower bound in Theorem 8 can be viewed as the maximum acyclic induced subgraph (MAIS) lower bound [37] for the index coding problem induced by the coded caching setup. The proof of this bound uses a technique from

[38] (Appendix A in [38]).

> **Theorem 8.** *Let $B$ be a bipartite caching graph representing a caching scheme on a broadcast network with $N$ files at the server. For some $N' \leq N$, let $U = \{k_j : j \in [N']\}$ be an arbitrary subset of $N'$ user vertices of $B$. For $j \in [N']$, let $\rho_j$ be the number of right vertices (subfiles) in $B$ which are adjacent to each vertex in $\{k_i : i \in [j]\}$. Let $\tilde{R}^*$ be the infimum of all achievable rates for the bipartite caching scheme defined by $B$. Then,*
> $$\tilde{R}^* F \geq \sum_{j=1}^{N} \rho_j.$$

*Proof:* We are given a valid caching scheme associated with $B$. As $N' \leq N$, we can assume a demand scenario in which the $N'$ users all demand different files. Let $\boldsymbol{Y}$ denote the set of all transmissions in a valid transmission scheme. Let $W_{d_j}$ be the demand and $Z_j$ be the cache content of the user $k_j \in U$. Let $S_j$ denote the set of subfiles of $W_{d_j}$ not cached in any of the users $k_i : i \in [j]$. This corresponds to the subfile vertices adjacent to all the users $k_i : i \in [j]$. In our notation, $|S_j| = \rho_j$. Since $W_{d_j}$s are distinct, thus each subfile in $S_j : j \in [N']$ is distinct. We then follow an idea similar to [38]. We construct a virtual receiver which contains an empty cache at first. In the $j^{th}$ step, the cache of this virtual user is populated with all the cache contents of the user $k_j$ except those pertaining to the files demanded by $k_i : i \in [j-1]$. Let $\tilde{Z}_j = Z_j \backslash \{W_{d_i,f} : i < j, \forall f\}$. Then $\{\tilde{Z}_j : j \in [N']\}$ is the final cache content of this virtual user. By the given transmission scheme, the receivers can decode their demands. Hence, we must have

$$H\left(\{W_{d_j} : j \in [N']\} \mid \{\tilde{Z}_j : j \in [N']\}, \boldsymbol{Y}\right) = 0, \tag{19}$$

as the virtual user must be successively able to decode all the demands of the $N'$ users. Since $RF$ denotes the size of transmissions (assuming each subfile to be of unit size) in a code of rate $R$, we must have the following inequalities.

$$\begin{aligned}
\tilde{R}^* F &\geq H(\boldsymbol{Y}) \\
&\geq I\left(\boldsymbol{Y}; \{W_{d_j} : j \in [N']\} \mid \{\tilde{Z}_j : j \in [N']\}\right) \\
&= H\left(\{W_{d_j} : j \in [N']\} \mid \{\tilde{Z}_j : j \in [N']\}\right) \quad \text{(by (19))} \\
&\geq H\left(\{S_j : j \in [N']\}\right) = \sum_{j=1}^{N'} \rho_j,
\end{aligned}$$

where $I(;)$ denotes the mutual information, and the last inequality is obtained by noting the missing subfiles in $\{\tilde{Z}_j : j \in [N']\}$. This completes the proof. ∎

Once again, we note that as the bound in Theorem 8 applies to the rate of any delivery scheme for the given bipartite caching scheme, not just the induced-matching based scheme. The following example illustrates the lower bound presented in Theorem 8.

**Example 6.** *(Continuation of Example 1, Example 2) Label the user vertices of the bipartite caching graph $B(4,5,3)$ presented in Fig. 2 as $k_1 = 1, k_2 = 3, k_3 = 2, k_4 = 4$. Assume that the number of files $N \geq 4$. Following Theorem 8, we get $\rho_1 = |\mathcal{N}(k_1)| = |\{f_1, f_2, f_3\}| = 3$, $\rho_2 = \left|\bigcap_{i=1}^{2} \mathcal{N}(k_i)\right| = |\{f_2, f_3\}| = 2$, $\rho_3 = \left|\bigcap_{i=1}^{3} \mathcal{N}(k_i)\right| = |\{f_2\}| = 1$, $\rho_4 = \left|\bigcap_{i=1}^{4} \mathcal{N}(k_i)\right| = |\phi| = 0$. We then get the lower bound on the rate as, $R^* \geq \frac{\sum_{j=1}^{4} \rho_j}{5} = \frac{6}{5}$. The induced matching based delivery scheme presented in Example 2 achieves the rate $R = \frac{6}{5}$. Therefore, for the symmetric caching scheme defined by $B(4,5,3)$, the rate $R = \frac{6}{5}$ is information theoretically optimal.*

Using the result in Theorem 8, we obtain our first lower bound in Corollary 3 for any symmetric caching scheme (left-regular bipartite caching schemes), and then the second bound in Theorem 9 for all bi-regular bipartite caching schemes.

**Corollary 3.** *For any symmetric caching scheme with $K$ users, cache fraction $\frac{M}{N}$, subpacketization $F$, and the number of files $N \geq K(1 - \frac{M}{N})$; the optimal rate $R^*$ satisfies*

$$R^* F \geq (K + F)\left(1 - \frac{M}{N}\right) - 1. \tag{20}$$

*Proof:* A symmetric caching scheme, with the given parameters, is equivalently represented by a bipartite caching graph $B(K, F, D)$ as in Section III-B. In any such graph, by the pigeon-holing argument, it is easy to see that there is a subfile vertex $f \in \mathcal{F}$ in $B$ having at least $K\left(1 - \frac{M}{N}\right)$ adjacent user vertices, which we refer to as $U = \{k_i : i \in [K\left(1 - \frac{M}{N}\right)]\}$. Clearly, by definition of $B$, the vertex $k_1$ has $D = F(1 - \frac{M}{N})$ neighbours. Thus, following the notations in Theorem 8, we have $\rho_1 = D$. Further, for $2 \leq i \leq K(1 - \frac{M}{N})$, we have $\rho_i \geq 1$, as $f$ is adjacent to each vertex in $U$. As the bounds on $\rho_i : \forall i \in \left[K(1 - \frac{M}{N})\right]$ are independent of the bipartite graph chosen and depend only on the given parameters $K, F, \frac{M}{N}$, by applying Theorem 8 for the bipartite graph corresponding to the caching scheme which gives the delivery scheme with optimal rate $R^*$, we get (20). ∎

In the following bound, we abuse the notation $R^*$ to denote the optimal rate (given $K, F, \frac{M}{N}$) among all delivery (including non-linear) schemes defined for a *bi-regular* bipartite caching graph $B(K, F, D)$. It is clear that for such a graph the right degree is $d \triangleq K\left(1 - \frac{M}{N}\right)$ which means that, each subfile is stored in a constant $\frac{KM}{N}$ number of clients. All known symmetric caching schemes belong to this class of graphs, to the best of our knowledge, and hence this bound is applicable to each of them. This bound also applies to the schemes which we have developed in this work, as they are obtained from bi-regular bipartite graphs (see Remarks 3 and 4). We show in Table XII that this second lower bound outperforms prior bounds for many chosen values of parameters.

**Theorem 9.** *Let $R^*$ be the infimum of all achievable rates for the coded caching problem defined by a bi-regular bipartite caching graph $B(K, F, D)$, with the right degree $d = K\left(1 - \frac{M}{N}\right)$ and assume that the number of files $N \geq K\left(1 - \frac{M}{N}\right)$. Then $R^*$ satisfies*

$$R^* F \geq D + \left\lceil \frac{(d-1)D}{K-1} \right\rceil + \cdots$$
$$\cdots + \left\lceil \frac{1}{\frac{KM}{N} + 1} \left\lceil \frac{2}{\frac{KM}{N} + 2} \left\lceil \cdots \left\lceil \frac{d-2}{K-2} \left\lceil \frac{(d-1)D}{K-1} \right\rceil \right\rceil \cdots \right\rceil \right\rceil \right\rceil$$

*Proof:* Consider a bi-regular bipartite caching graph $B(K, F, D)$. We know that $D = F\left(1 - \frac{M}{N}\right)$ and $d = K\left(1 - \frac{M}{N}\right)$. Consider an arbitrary user vertex and call it $k_1$. We know that $|\mathcal{N}(k_1)| = D$. From the notations of Theorem 8, we have $\rho_1 = D$. WLOG, let $\mathcal{N}(k_1) = \{f_1, f_2, \cdots, f_{\rho_1}\}$. Now, consider the graph induced by $\mathcal{K} \cup \mathcal{N}(k_1)$ of $B$ and call it $B'$.

**Finding lower bound on $\rho_2$:** Since the degree of each $f \in \mathcal{N}(k_1)$ is $d$, there are exactly $d\rho_1$ edges in $B'$. It is easy to see that the number of edges in $B'$ between $\mathcal{K} \setminus \{k_1\}$ and $\mathcal{N}(k_1)$ is $(d-1)\rho_1$. By the pigeon-holing argument, there exists a user vertex in $\mathcal{K} \setminus \{k_1\}$ with degree at least $\left\lceil \frac{(d-1)\rho_1}{K-1} \right\rceil$, in $B'$. Consider such a user vertex and call it $k_2$. Therefore,

$$\rho_2 \geq \left\lceil \frac{(d-1)\rho_1}{K-1} \right\rceil = \left\lceil \frac{(d-1)D}{K-1} \right\rceil.$$

WLOG, let $\mathcal{N}(k_2) = \{f_1, f_2, \cdots, f_{\rho_2}\}$ in $B'$. It is clear that $\mathcal{N}(k_2) \subseteq \mathcal{N}(k_1)$. Therefore, from each subfile vertex in $\mathcal{N}(k_2)$ two edges will go to $k_1, k_2$ each.

**Finding lower bound on $\rho_3$:** Since the degree of each $f \in \mathcal{N}(k_2)$ is $d$, there are exactly $d\rho_2$ edges incident at the edges in $\mathcal{N}(k_2)$. Now, it should be clear that the number of edges in $B'$ between $\mathcal{K} \setminus \{k_1, k_2\}$ and $\mathcal{N}(k_2)$ is $(d-2)\rho_2$. Again by the pigeon-holing argument, there exists a user vertex in $\mathcal{K} \setminus \{k_1, k_2\}$ with degree at least $\left\lceil \frac{(d-2)\rho_2}{K-2} \right\rceil$, in $B'$. Consider such a user vertex and call it $k_3$. Therefore,

$$\rho_3 \geq \left\lceil \frac{(d-2)\rho_2}{K-2} \right\rceil.$$

TABLE XII: Comparison of the two proposed information theoretic lower bounds on the rate of the transmission schemes associated with the special class of caching schemes defined by the bi-regular bipartite caching graphs with parameters $K, F, D$, with that of [18], [11]. The last column gives the number of transmissions in the scheme constructed in Section V-C for whatever values are applicable. (NA represents not applicable).

| $K$ | $F$ | $D$ | $\frac{M}{N}$ $=1-\frac{D}{F}$ | (Corollary 3) $R^* \geq$ | (Theorem 9) $R^* \geq$ | [18] $R^* \geq$ | [11] $R^* \geq$ | Scheme B [Section V-C] $R$ |
|---|---|---|---|---|---|---|---|---|
| 15 | 50 | 30 | 0.4 | 0.76 | 1.42 | 1.08 | 1.2857 | NA |
| 24 | 54 | 36 | 0.3333 | 0.9444 | 2.0185 | 1.6667 | 1.7778 | NA |
| 15 | 20 | 12 | 0.4 | 1 | 1.5 | 1.55 | 1.2857 | NA |
| 7 | 21 | 12 | 0.4286 | 0.7143 | 1.0476 | 0.8571 | 1 | 1.33 |
| 13 | 78 | 54 | 0.3077 | 0.7949 | 1.8333 | 1.3205 | 1.8 | 3 |
| 15 | 105 | 84 | 0.2 | 0.9048 | 3.0952 | 2.2286 | 3 | 4 |
| 21 | 210 | 160 | 0.2381 | 0.8333 | 2.7381 | 1.7952 | 2.6667 | 105 |
| 31 | 465 | 420 | 0.0968 | 0.9613 | 7.0839 | 4.9247 | 7 | 9.3333 |
| 40 | 780 | 702 | 0.1 | 0.9449 | 7.2936 | 4.7397 | 7.2 | 12 |
| 105 | 105 | 48 | 0.5429 | 0.9048 | 1.2571 | 1.2571 | 0.8276 | 8 |
| 465 | 4340 | 1792 | 0.5871 | 0.4569 | 0.7435 | 0.4880 | 0.7007 | 19.2 |
| 4340 | 465 | 192 | 0.5871 | 4.2645 | 4.5548 | 6.9398 | 0.7030 | 179.2 |
| 465 | 465 | 335 | 0.28 | 1.4387 | 3.8215 | 3.8215 | 2.5573 | 56 |
| 8001 | 9,921,240 | 6,666,081 | 0.3281 | 0.6724 | 2.0479 | 1.0330 | 2.0471 | 358.4 |

By using the lower bound on $\rho_2$ we can write,

$$\rho_3 \geq \left\lceil \frac{d-2}{K-2} \left\lceil \frac{(d-1)\rho_1}{K-1} \right\rceil \right\rceil.$$

Note that these $\rho_3$ edges are incident on $\rho_3$ subfile vertices in $B'$, which we denote as $\mathcal{N}(k_3)$. By construction of $k_3$ we note that $\mathcal{N}(k_3) \subseteq \mathcal{N}(k_2)$. Iterating this procedure for $j = 4, \ldots, d$, we can identify vertices $k_j$, with neighbouring vertices $\mathcal{N}(k_j)$ in $B'$ numbering $\rho_j$ respectively, such that $\mathcal{N}(k_j) \subseteq \mathcal{N}(k_{j-1}), \forall j \leq d$, satisfying the following inequality for each $j$ by the pigeon-hole principle

$$\rho_j \geq \left\lceil \frac{(d-(j-1))\,\rho_{j-1}}{K-(j-1)} \right\rceil.$$

Note that $d = K\left(1 - \frac{M}{N}\right)$. As $N \geq K\left(1 - \frac{M}{N}\right)$, we have in Theorem 8, $N' = K\left(1 - \frac{M}{N}\right)$. Now, applying Theorem 8 upon noticing that the bounds on $\rho_j$s depend only on the parameters $K, F, \frac{M}{N}$, completes the proof. ∎

We now present the numerical comparisons (in Table XII) between the information theoretic bounds we have obtained in this section, with earlier results. Throughout this numerical comparison, we assume that the caching schemes are defined by the bi-regular bipartite caching graphs. For a number of choices of parameters $\left(K, F \text{ and } \frac{M}{N}\right)$, in Table XII, we compare numerically the new subpacketization-dependent lower bounds based on Corollary 3 and Theorem 9 on the optimal rate (column 5 and column 6 of Table XII), with the lower bound (18) of [18] (given in column 7), as well as the lower bound of [11] $\left(R^* \geq \frac{K(1-\frac{M}{N})}{1+\frac{MK}{N}}\right.$ , calculated in column 8$\Big)$. Note that the bound in [18] holds for PDA based delivery schemes with given

subpacketization level, while the bound in [11] holds for non-linear schemes as well and is subpacketization-independent.

It can be seen that for many of the chosen parameters, our bound of Theorem 9 is better than those in [18], [11] (when applied to the special class of bi-regular bipartite caching graphs). However, the bound in Corollary 3 is quite loose. Further, the last column of Table XII denotes the rate achieved (for whichever parameters are applicable) by our new coded caching scheme, titled Scheme B, whose construction we have presented in Section V-C (from Remark 4, recall that Scheme B is defined on a bi-regular bipartite caching graph). Also, seeing the table, we remark that there is in general a wide gap between the lower bounds and achievable rates for small subpacketization levels, which indicate scope for future work.

## VIII. Conclusion and Discussion

In this work, we have presented the coded caching schemes which achieve low subpacketization compared to a number of existing schemes in the literature. Our coded caching schemes are constructed over the foundations of the bipartite graphs and the projective geometry. The literature in this area is now extensive, therefore we have presented comparisons with only a few important existing coded caching schemes, which are considered state-of-the-art to the best of our knowledge. We have also extended our scheme to other channel settings, thereby showing similar gains in the subpacketization in those settings also. There are a number of questions and open problems, which are yet to be answered in terms of the subpacketization-rate trade-off, some of which are listed here.

- Is there a closed-form expression for the optimal rate achievable for a given subpacketization level? Can we obtain impossibility results, for certain regimes of (asymptotic) coded caching gains given a certain level of subpacketization, or vice-versa? (For instance, the work [14] gives an impossibility of the subpacketization being linear in $K$, for constant rate and cache-fraction).

- The current paper can be said to subsume the construction of [6] in the following sense. Suppose we substitute $n = 1$ in Theorem 3, and let $q \to 0$, then it is easy to see that we obtain the scheme from [6]. Further, retaining $n$ as a parameter, just letting $q \to 0$ gives us a scheme from [15] (which is based on the set-theoretic principles). The likely common theory which underlies these types of combinatorial constructions, which are based on set-containment and subspace-containment principles, is the notion of geometric lattices. Exploring such connections may lead to further interesting and useful constructions.

- One drawback of our constructions in this work is the number of parameters and the lack of flexibility in designing for all possible values for the number of users, and for cache sizes, without introducing dummy users or wasting existing user cache memory. It is possible that using a similar grouping scheme as in [12], we could achieve some flexibility in the number of users. Exploring this is a direction for future work.

- The single server with multi-antenna and multi-receiver wireless scenario can be looked at as a multi-transmitter scenario (which we discuss in this current work) with all the transmitters having access to the entire library. Recently, the works [39], [40] carry forward the discussion of subpacketization in the multi-antenna wireless communication. In particular, in [40], a simple linear subpacketization scheme is presented for multi-antenna scenario provided some parameter conditions are satisfied. This raises the question of whether there are schemes for other scenarios, which are of low theoretical complexity, low subpacketization, and offer good caching gain.

## Appendix A
### Proof of Lemma 4

Construct a bipartite graph with left vertices as $\left\{ V_i, i = 1, \ldots, \begin{bmatrix} m+t \\ t \end{bmatrix}_q \right\}$ (the $t$-dim subspaces of $X$) and right vertices as $\{T : T \text{ is a } m\text{-dim subspace of } X\}$. By the lemma statement, we know that the number of right vertices is $\begin{bmatrix} m+t \\ m \end{bmatrix}_q$. By A1 of Lemma 1, we have $\begin{bmatrix} m+t \\ m \end{bmatrix}_q = \begin{bmatrix} m+t \\ t \end{bmatrix}_q$. Therefore, the number of right vertices is equal to the number of left vertices.

We now define the edges of the bipartite graph. For a left vertex $V$, let the adjacent right-vertices in the bipartite graph be $\{T : V \cap T = \{\mathbf{0}\}\}$.

Thus, the left-degree is $\begin{bmatrix} m+t \\ t \end{bmatrix}_q - |\{T : V \cap T \neq \{\mathbf{0}\}\}|$. Now, $V \cap T$ is a subspace. By A3 of Lemma 1, the number of $m$-dim subspaces of $X$ intersecting a fixed $t$-dim subspace in some $i$-dim subspace ($1 \leq i \leq min(t,m)$) is

$$q^{(m-i)(t-i)} \begin{bmatrix} (m+t) - (t) \\ m-i \end{bmatrix}_q \begin{bmatrix} t \\ i \end{bmatrix}_q = q^{(m-i)(t-i)} \begin{bmatrix} m \\ i \end{bmatrix}_q \begin{bmatrix} t \\ i \end{bmatrix}_q .$$

Thus the left-degree in this bipartite graph is

$$\begin{bmatrix} m+t \\ t \end{bmatrix}_q - \sum_{i=1}^{min(m,t)} q^{(m-i)(t-i)} \begin{bmatrix} m \\ i \end{bmatrix}_q \begin{bmatrix} t \\ i \end{bmatrix}_q ,$$

where the second term above is precisely $|\{T : V \cap T \neq \{\mathbf{0}\}\}|$.

Similarly, by Lemma 1, the number of $t$-dim subspaces of $X$ intersecting a fixed $m$-dim subspace in some $i$-dim subspace is

$$q^{(t-i)(m-i)} \begin{bmatrix} (m+t) - (m) \\ t-i \end{bmatrix}_q \begin{bmatrix} m \\ i \end{bmatrix}_q = q^{(t-i)(m-i)} \begin{bmatrix} t \\ i \end{bmatrix}_q \begin{bmatrix} m \\ i \end{bmatrix}_q .$$

And hence the right degree is equal to the left-degree. Hence, the bipartite graph we have constructed is regular.

A perfect matching of a graph $G$ is a matching of $G$ such that every vertex of $G$ is incident on some edge of the matching. It should be clear that what we are looking for is a perfect matching of the regular bipartite graph we have constructed. The reason is as follows. Define $T_{V_i,X}$ as the $m$-dim subspace adjacent to $V_i$ in the perfect matching. Since for given $V_i$, any $T$ adjacent to $V_i$ in our bipartite graph is such that $T \oplus V_i = X$, thus we have $T_{V_i,X} \oplus V_i = X$. Thus, a perfect matching gives us the collection of $T_{V_i,X}, \forall V_i$ as we desire.

Now, for a regular bipartite graph with $n$ left-vertices, algorithms are known to find a perfect matching with complexity as small as $O(n \log n)$ [41]. This completes the proof.

APPENDIX B

ASYMPTOTIC ANALYSIS OF SCHEME A (THEOREM 2)

We analyse the asymptotic behaviour of our scheme as $K$ grows large in two cases. In the first case, we bound $\frac{M}{N}$ from above by a constant and analyse the asymptotic behaviour of $F$ and $R$. In the second case, we bound $R$ from above by a constant and analyse the asymptotic behaviour of $F$ and $\frac{M}{N}$.

We first give the equivalent expressions of $\frac{M}{N}$ and $R$ (which can be easily verified using the definition of the Gaussian binomial coefficient). We do this because we can apply the bounds in Lemma 2 conveniently to these expressions in order to obtain our asymptotics.

$$\frac{M}{N} = 1 - \frac{\begin{bmatrix} k-t \\ m \end{bmatrix}_q}{\begin{bmatrix} k \\ m+t \end{bmatrix}_q} = 1 - \frac{\begin{bmatrix} m+t \\ t \end{bmatrix}_q}{\begin{bmatrix} k \\ t \end{bmatrix}_q} \quad \text{and}$$

$$R = \frac{\begin{bmatrix} k \\ m \end{bmatrix}_q}{\begin{bmatrix} k \\ m+t \end{bmatrix}_q} = \frac{\begin{bmatrix} m+t \\ t \end{bmatrix}_q}{\begin{bmatrix} k-m \\ t \end{bmatrix}_q} .$$

Throughout we assume $q$ is constant. We have $K = \begin{bmatrix} k \\ t \end{bmatrix}_q$. We analyse our scheme as $k$ grows large.

By using (4) we have

$$q^{(k-t)t} \leq K \leq q^{(k-t+1)t}. \tag{21}$$

We can write this as,

$$\frac{1}{\sqrt{K}} q^{\frac{-t^2}{2}} \leq q^{\frac{-kt}{2}} \leq \frac{1}{\sqrt{K}} q^{\frac{-t^2+t}{2}}. \tag{22}$$

**Case 1:** If $\frac{M}{N}$ is upper bounded by a constant

From Theorem 2 we have,

$$1 - \frac{M}{N} = \frac{\begin{bmatrix} m+t \\ t \end{bmatrix}_q}{\begin{bmatrix} k \\ t \end{bmatrix}_q} \overset{(5)}{\geq} q^{(m+t-k-1)t}. \tag{23}$$

To lower bound $1 - \frac{M}{N}$ (or upper bound $\frac{M}{N}$) by a constant, assume $t$ and $k-m$ as constants. Note that $m+t \leq k$.

**Asymptotics of $F$:** We now analyse the asymptotics for $F$. Consider,

$$\frac{F}{K} = \frac{\begin{bmatrix} k \\ m+t \end{bmatrix}_q}{\begin{bmatrix} k \\ t \end{bmatrix}_q} \overset{(6)}{\leq} q^{(k-t-m-t+1)m} \leq q^{(k-2t-m+1)(k-t)}$$

$$\overset{(21)}{\leq} K^{\frac{k-2t-m+1}{t}}$$

$$\text{Hence} \quad F \leq K^{\frac{k-t-m+1}{t}}.$$

Therefore $F = O(poly(K))$. (since $t$ and $k-m$ are constants)

**Asymptotics of $R$:**

We now analyse the asymptotics for $R = \frac{\begin{bmatrix} m+t \\ t \end{bmatrix}_q}{\begin{bmatrix} k-m \\ t \end{bmatrix}_q}$.

By using (5) we get,

$$q^{(m+t-k+m-1)t} \leq R \leq q^{(m+t-k+m+1)t}.$$

Now by using (21) we can write,

$$q^{2(t+m-k-1)t} \leq \frac{R}{K} \leq q^{2(t+m-k)t+t}.$$

Therefore $R = \Theta(K)$. (since $t$ and $k-m$ are constants)

**Case 2:** If $R$ is upper bounded by a constant

We have, $R \overset{(5)}{\leq} q^{(2m-k+t+1)t}$.

To bound $R$ from above by a constant, assume $t$ and $k-2m$ as constants. Note that $m+t \leq k$.

**Asymptotics of $F$:** We now analyse the asymptotics for $F$. Consider,

$$\frac{F}{K} \overset{(6)}{\leq} q^{(k-t-m-t+1)m} \leq q^{(k-2t-m+1)(k-t)}$$

$$F \leq K \; q^{\left(\frac{k}{2}+\frac{k-2m}{2}-2t+1\right)(k-t)}$$

$$\leq q^{\log_q K} \; q^{\frac{k^2-kt}{2}+\alpha_1(k-t)}, \tag{24}$$

where $\alpha_1 = \frac{k-2m}{2} - 2t + 1$ is a constant. From (21) we have $k \leq \frac{1}{t} \log_q K + t$.

By using this in the inequality in (24), it is easy to see that $F = q^{O\left((\log_q K)^2\right)}$. (since $t$ and $k - 2m$ are constants)

**Asymptotics of $\frac{M}{N}$:** We now analyse the asymptotics for $\frac{M}{N}$. By using (5) we have,

$$q^{(m+t-k-1)t} \leq 1 - \frac{M}{N} \leq q^{(m+t-k+1)t}$$

$$q^{\left(\frac{2m-k}{2}+t-1\right)t}q^{\frac{-kt}{2}} \leq 1 - \frac{M}{N} \leq q^{\left(\frac{2m-k}{2}+t+1\right)t}q^{\frac{-kt}{2}}.$$

By using (22) we get,

$$\frac{1}{\sqrt{K}}q^{\left(\frac{2m-k}{2}+t-1\right)t-\frac{t^2}{2}} \leq 1 - \frac{M}{N} \leq \frac{1}{\sqrt{K}}q^{\left(\frac{2m-k}{2}+t+1\right)t-\frac{t^2+t}{2}}.$$

Therefore $\frac{M}{N} = 1 - \Theta\left(\frac{1}{\sqrt{K}}\right)$. (since $t$ and $k - 2m$ are constants)

## APPENDIX C
### PROOF OF LEMMA 9

**Finding the number of user vertices $K\,(= |\mathbb{X}|)$:**

Finding $K$ means finding the number of distinct sets $\{T_1, T_2, \cdots, T_n\}$ such that $T_i \in \mathbb{T}, \; \forall i \in [n]$ and $\sum\limits_{i=1}^{n} T_i \in \mathbb{R}$. By invoking Lemma 8 (with $a = 0$), we have,

$$K = \frac{1}{n!} \prod_{i=0}^{n-1} (\theta(k) - \theta(i)) \tag{25}$$

$$= \frac{1}{n!} \prod_{i=0}^{n-1} \left(\frac{q^k - 1}{q-1} - \frac{q^i - 1}{q-1}\right)$$

$$= \frac{1}{n!} \prod_{i=0}^{n-1} \frac{q^k - q^i}{q-1} = \frac{1}{n!} \left(\prod_{i=0}^{n-1} q^i\right) \left(\prod_{i=0}^{n-1} \frac{q^{k-i} - 1}{q-1}\right)$$

$$= \frac{1}{n!} \; q^{\frac{n(n-1)}{2}} \prod_{i=0}^{n-1} \begin{bmatrix} k-i \\ 1 \end{bmatrix}_q.$$

**Finding the number subfile vertices $F\,(= |\mathbb{Y}|)$:**

Proof is similar to that of $K$ (replace $n$ with $m$).

**Finding the degree of user vertex $D\,(= |\mathcal{N}(X)|)$:**

Consider an arbitrary $X = \{T_1, T_2, \cdots, T_n\} \in \mathbb{X}$. We have $\sum\limits_{i=1}^{n} T_i = R$, for some $R \in \mathbb{R}$. We know that $dim(R) = n$. Now, finding $|\mathcal{C}_X|$ is equivalent to counting the number of distinct sets $Y = \{T'_1, T'_2, \cdots, T'_m\} \in \mathbb{Y}$ such that $X \cup Y \in \mathbb{Z}$, that is $dim\left(\sum\limits_{i=1}^{n} T_i + \sum\limits_{i=1}^{m} T'_i\right) = n + m$. By Lemma 8 we have,

$$|\mathcal{N}(X)| = \frac{1}{m!} \prod_{i=0}^{m-1} (\theta(k) - \theta(n+i))$$

$$= \frac{1}{m!} \prod_{i=0}^{m-1} \frac{q^k - q^{n+i}}{q-1}$$

$$= \frac{1}{m!} \left( \prod_{i=0}^{m-1} q^{n+i} \right) \left( \prod_{i=0}^{m-1} \frac{q^{k-n-i} - 1}{q-1} \right)$$

$$= \frac{q^{nm}}{m!} \, q^{\frac{m(m-1)}{2}} \prod_{i=0}^{m-1} \begin{bmatrix} k-n-i \\ 1 \end{bmatrix}_q.$$

This completes the proof.

<center>APPENDIX D</center>

<center>ASYMPTOTIC ANALYSIS OF SCHEME B (THEOREM 3)</center>

In this appendix, we analyse the asymptotic behaviour of $F, R$ for our coded caching scheme presented in Theorem 3 (Scheme B) as $\frac{M}{N}$ is upper bounded by a constant and $K \to \infty$. We show that $F = q^{O\left((\log_q K)^2\right)}$, while $R = \Theta\left( \frac{K}{(\log_q K)^n} \right)$. Throughout our analysis, we assume $q$ is a constant and some prime power, and $n$ is some constant. We now upper bound $\frac{M}{N}$ by a constant. From (11) we have,

$$1 - \frac{M}{N} = \frac{\prod_{i=0}^{m-1} (\theta(k) - \theta(n+i))}{\prod_{i=0}^{m-1} (\theta(k) - \theta(i))} = \frac{\prod_{i=n}^{n+m-1} (\theta(k) - \theta(i))}{\prod_{i=0}^{m-1} (\theta(k) - \theta(i))}$$

$$= \frac{\prod_{i=0}^{n+m-1} (\theta(k) - \theta(i))}{\left( \prod_{i=0}^{n-1} (\theta(k) - \theta(i)) \right) \left( \prod_{i=0}^{m-1} (\theta(k) - \theta(i)) \right)}$$

$$= \frac{\prod_{i=m}^{n+m-1} (\theta(k) - \theta(i))}{\prod_{i=0}^{n-1} (\theta(k) - \theta(i))} = \frac{\prod_{i=0}^{n-1} (\theta(k) - \theta(m+i))}{\prod_{i=0}^{n-1} (\theta(k) - \theta(i))}$$

$$= \prod_{i=0}^{n-1} \frac{q^k - q^{m+i}}{q^k - q^i} = \prod_{i=0}^{n-1} \frac{q^{k-i} - q^m}{q^{k-i} - 1}$$

$$\geq \prod_{i=0}^{n-1} \frac{q^{k-i} - q^m}{q^{k-i}}$$

$$1 - \frac{M}{N} \geq \prod_{i=0}^{n-1} \left( 1 - \frac{q^i}{q^{k-m}} \right).$$

Let $\alpha = k - m$. ($\alpha \geq n$, since $k \geq n + m$)

$$1 - \frac{M}{N} \geq \prod_{i=0}^{n-1} \left( 1 - \frac{q^i}{q^\alpha} \right) \geq \prod_{i=0}^{n-1} \left( 1 - \frac{q^{n-1}}{q^\alpha} \right)$$

$$\geq \left( 1 - \frac{1}{q^{\alpha-n+1}} \right)^n \geq 1 - \frac{n}{q^{\alpha-n+1}}. \tag{26}$$

Therefore, the upper bound on $\frac{M}{N}$ is given as $\frac{M}{N} \leq \frac{n}{q^{\alpha-n+1}}$, where $\alpha = k - m$ and $n$ are constants.

We have $K = \frac{1}{n!} \, q^{\frac{n(n-1)}{2}} \prod_{i=0}^{n-1} \begin{bmatrix} k-i \\ 1 \end{bmatrix}_q$ . We analyse our scheme as $k$ grows large (thus $K$ grows large). By Lemma 2 we have,

$$\frac{1}{n!} \, q^{\frac{n(n-1)}{2}} \prod_{i=0}^{n-1} q^{k-i-1} \leq K \leq \frac{1}{n!} \, q^{\frac{n(n-1)}{2}} \prod_{i=0}^{n-1} q^{k-i}$$

$$\frac{\prod_{i=0}^{n-1} q^i}{n!} \, q^{(k-1)n} \prod_{i=0}^{n-1} q^{-i} \leq K \leq \frac{\prod_{i=0}^{n-1} q^i}{n!} \, q^{kn} \prod_{i=0}^{n-1} q^{-i}$$

$$\frac{1}{n!} \, q^{(k-1)n} \leq K \leq \frac{1}{n!} \, q^{kn} \tag{27}$$

$$(k-1)n \leq \log_q (n!K) \leq kn.$$

Hence, we have,

$$\frac{1}{n} \log_q (n!K) \leq k \leq \frac{1}{n} \log_q (n!K) + 1. \tag{28}$$

**Asymptotics of $R$:**

We now get the asymptotics for the rate. We have, $R = \dfrac{K(1 - \frac{M}{N})}{\gamma}$. From Theorem 3, we have $\gamma = \binom{n+m}{n}$. Since $\alpha = k-m$ we can write, $\gamma = \binom{k-\alpha+n}{n}$. We have the following well known bounds on the binomial coefficient ($e$ being the base of the natural logarithm),

$$\left(\frac{a}{b}\right)^b \leq \binom{a}{b} \leq e^b \left(\frac{a}{b}\right)^b.$$

By using this result, the bounds on $\gamma$ can be written as,

$$\left(\frac{k-\alpha+n}{n}\right)^n \leq \gamma \leq \left(\frac{e(k-\alpha+n)}{n}\right)^n.$$

By using (28) the lower bound on $\gamma$ can be written as,

$$\left(\frac{\frac{1}{n}\log_q (n!K) - \alpha + n}{n}\right)^n \leq \gamma \quad,$$

and the upper bound on $\gamma$ can be written as,

$$\gamma \leq \left(\frac{e\left(\frac{1}{n}\log_q (n!K) - \alpha + n + 1\right)}{n}\right)^n.$$

After some simple manipulations we get $\gamma = \Theta\left((\log_q n!K)^n\right) = \Theta\left((\log_q K)^n\right)$. (since $n$ is a constant). Therefore, we get $R = \Theta\left(\frac{K}{(\log_q K)^n}\right)$.

**Asymptotics of $F$:**

We now obtain the asymptotics for the subpacketization $F$. By using $K, F$ expressions in Theorem 3 we get,

$$\frac{F}{K} = \frac{n!}{m!} \frac{q^{\frac{m(m-1)}{2}} \prod_{i=0}^{m-1} \begin{bmatrix} k-i \\ 1 \end{bmatrix}_q}{q^{\frac{n(n-1)}{2}} \prod_{i=0}^{n-1} \begin{bmatrix} k-i \\ 1 \end{bmatrix}_q}.$$

By Lemma 2 we have,

$$\frac{F}{K} \leq \frac{n!}{m!} \frac{q^{\frac{m(m-1)}{2}}}{q^{\frac{n(n-1)}{2}}} \frac{\prod_{i=0}^{m-1} q^{k-i}}{\prod_{i=0}^{n-1} q^{k-i-1}}$$

$$= \frac{n!}{m!} \frac{q^{\frac{m(m-1)}{2}}}{q^{\frac{n(n-1)}{2}}} \frac{\prod_{i=0}^{m} q^{-i}}{\prod_{i=0}^{n} q^{-i}} \frac{q^{km}}{q^{(k-1)n}}$$

$$F \leq \frac{n! \ K}{m!} \ q^{km-kn+n}.$$

By using $m = k - \alpha$ we get,

$$F \leq \frac{q^{\log_q (n!K)} \ q^{\left(k^2 + (\alpha+n)(-k)+n\right)}}{(k-\alpha)!}.$$

By (28) we have,

$$k^2 + (\alpha + n)(-k) + n$$

$$\leq \left(\frac{1}{n} \log_q (n!K) + 1\right)^2 + (\alpha + n)\left(\frac{-1}{n} \log_q (n!K)\right) + n$$

$$= \left(\frac{1}{n} \log_q (n!K)\right)^2 + \left(\frac{2-\alpha-n}{n} \log_q (n!K)\right) + n + 1.$$

By the lower bound of (28) we have,

$$\frac{1}{(k-\alpha)!} \leq \frac{1}{\lfloor \frac{1}{n} \log_q (n!K) - \alpha \rfloor!}.$$

By using these bounds, we get,

$$F \leq \frac{q^{\left(\frac{1}{n} \log_q (n!K)\right)^2 + \left(\frac{2-\alpha}{n} \log_q (n!K)\right) + n + 1}}{\lfloor \frac{1}{n} \log_q (n!K) - \alpha \rfloor!}.$$

Using Stirling's approximation for $x!$ as $\sqrt{2\pi x} \left(\frac{x}{e}\right)^x$ for large $x$, and after some simple manipulations, we see that $F = q^{O\left((\log_q (n!K))^2\right)} = q^{O\left((\log_q K)^2\right)}$ (Since $n$ is a constant).

## APPENDIX E
## PROOF OF LEMMA 12

**Finding the number of user vertices $K (= |\mathbb{X}|)$:**

Let $\mathbb{T} \triangleq PG_q(k-1, 0)$ (set of all 1-dim subspaces). Finding $K$ means finding the number of distinct sets $\{L_1, L_2, \cdots, L_n\}$ such that $L_i \in \mathbb{L}, \ \forall i \in [n]$ and $\sum_{i=1}^{n} L_i \in \mathbb{R}$. Now to find $K$, we prove the following smaller claims.

*Claim 1:* The number (say $x_1$) of distinct $nl$ sized sets $\{T_1, T_2, \cdots, T_{nl}\}$ such that $T_i \in \mathbb{T}, \forall i \in [nl]$ and $\sum_{i=1}^{n} T_i \in \mathbb{R}$ is $x_1 = \frac{1}{(nl)!} \prod_{i=0}^{nl-1} (\theta(k) - \theta(i))$.

*Proof of Claim 1:* By invoking Lemma 8 with $a = 0$ and $b = nl$ we see the result.

*Claim 2:* Consider a set $\mathcal{T} = \{T_1, T_2, \cdots, T_{nl}\}$ such that $T_i \in \mathbb{T}, \forall i \in [nl]$ and $\sum_{i=1}^{n} T_i \in \mathbb{R}$. The number (say $x_2$) of distinct $X \in \mathbb{X}$ generated from $\mathcal{T}$ is $x_2 = \frac{1}{n!} \prod_{i=0}^{n-1} \left(\binom{(n-i)l}{l}\right)$.

*Proof of Claim 2:* It is clear that $\mathcal{T}$ contains $nl$ number of linearly independent 1-dim subspaces. So the addition of any $l$ subspaces from $\mathcal{T}$ will generate a unique $l$-dim subspace. Hence, finding $x_2$ is equivalent to counting the number of distinct $n$-sized sets of $l$-sized sets that can be formed from $\mathcal{T}$. It is clear that each such $n$-sized set will generate a unique $X$. From the basic combinatorics, the expression for $x_2$ can be inferred.

Therefore, the total number of $X \in \mathbb{X}$ which can be generated from $\mathbb{F}_q^k$ is $x_1 x_2$. But there are some repetitions in $x_1 x_2$. We identify them in the following claim.

<u>*Claim 3:*</u> Consider an arbitrary $X = \{L_1, L_2, \cdots, L_n\} \in \mathbb{X}$. The number (say $x_3$) of distinct $\{T_1, T_2, \cdots, T_{nl}\}$ (such that $T_i \in \mathbb{T}, \forall i \in [nl]$ and $\sum_{i=1}^{nl} T_i \in \mathbb{R}$) which generate $X$ is $x_3 = \left( \frac{1}{(l!)} \prod_{i=0}^{l-1} (\theta(l) - \theta(i)) \right)^n$.

*Proof of Claim 3:* Consider an arbitrary $L_j \in X$. By Lemma 8, the number of distinct sets $\{T_1^{'}, T_2^{'}, \cdots, T_l^{'}\}$ ($\forall T_i^{'} \in \mathbb{T}, i \in [l]$) such that $\sum_{i=1}^{l} T_i^{'} = L_j$ is (substitute $a = 0, b = l, k = l$ in Lemma 8) $\frac{1}{(l!)} \prod_{i=0}^{l-1} (\theta(l) - \theta(i))$. There are $n$ such $L_i$'s in $X$ (all are linearly independent subspaces). Therefore $x_3 = \left( \frac{1}{(l!)} \prod_{i=0}^{l-1} (\theta(l) - \theta(i)) \right)^n$.

Hence $K = \frac{x_1 x_2}{x_3}$. Now by using $\theta(k) = \frac{q^k - 1}{q - 1}$ and by doing some simple manipulations we see the expression of $K$ as per the lemma statement. The proofs of $F$ and $D$ follows similarly.

## REFERENCES

[1] P. Krishnan, "Coded caching via line graphs of bipartite graphs," in *2018 IEEE Information Theory Workshop (ITW)*, 2018, pp. 1–5.

[2] H. H. Suthan Chittoor, B. M., and P. Krishnan, "Coded caching via projective geometry: A new low subpacketization scheme," in *2019 IEEE International Symposium on Information Theory (ISIT)*, 2019, pp. 682–686.

[3] H. H. Suthan Chittoor and P. Krishnan, "Projective geometry based coded caching schemes with subexponential and linear subpacketizations," in *2019 19th International Symposium on Communications and Information Technologies (ISCIT)*, 2019, pp. 537–542.

[4] H. H. S. Chittoor and P. Krishnan, "Low subpacketization coded caching via projective geometry for broadcast and d2d networks," in *2019 IEEE Global Communications Conference (GLOBECOM)*, 2019, pp. 1–6.

[5] *Cisco Visual Networking Index: Forecast and Trends, 2017-2022*. [Online]. Available: www.cisco.com

[6] M. A. Maddah-Ali and U. Niesen, "Fundamental limits of caching," *IEEE Transactions on Information Theory*, vol. 60, no. 5, pp. 2856–2867, May 2014.

[7] M. A. Maddah-Ali and U. Niesen, "Decentralized coded caching attains order-optimal memory-rate tradeoff," *IEEE/ACM Transactions on Networking*, vol. 23, no. 4, pp. 1029–1040, Aug 2015.

[8] J. Zhang, X. Lin, and X. Wang, "Coded caching under arbitrary popularity distributions," *IEEE Transactions on Information Theory*, vol. 64, no. 1, pp. 349–366, 2018.

[9] R. Pedarsani, M. A. Maddah-Ali, and U. Niesen, "Online coded caching," *IEEE/ACM Transactions on Networking*, vol. 24, no. 2, pp. 836–845, April 2016.

[10] N. Karamchandani, U. Niesen, M. A. Maddah-Ali, and S. N. Diggavi, "Hierarchical coded caching," *IEEE Transactions on Information Theory*, vol. 62, no. 6, pp. 3212–3229, June 2016.

[11] K. Wan, D. Tuninetti, and P. Piantanida, "An index coding approach to caching with uncoded cache placement," *IEEE Transactions on Information Theory*, vol. 66, no. 3, pp. 1318–1332, 2020.

[12] K. Shanmugam, M. Ji, A. M. Tulino, J. Llorca, and A. G. Dimakis., "Finite-length analysis of caching-aided coded multicasting," *IEEE Transactions on Information Theory*, vol. 62, no. 10, pp. 5524–5537, Oct 2016.

[13] Q. Yan, M. Cheng, X. Tang, and Q. Chen, "On the placement delivery array design for centralized coded caching scheme," *IEEE Transactions on Information Theory*, vol. 63, no. 9, pp. 5821–5833, Sep. 2017.

[14] C. Shangguan, Y. Zhang, and G. Ge, "Centralized coded caching schemes: A hypergraph theoretical approach," *IEEE Transactions on Information Theory*, vol. 64, no. 8, pp. 5755–5766, Aug 2018.

[15] Q. Yan, X. Tang, Q. Chen, and M. Cheng, "Placement delivery array design through strong edge coloring of bipartite graphs," *IEEE Communications Letters*, vol. 22, no. 2, pp. 236–239, Feb 2018.

[16] L. Tang and A. Ramamoorthy, "Coded caching schemes with reduced subpacketization from linear block codes," *IEEE Transactions on Information Theory*, vol. 64, no. 4, pp. 3099–3120, April 2018.

[17] K. Shanmugam, A. M. Tulino, and A. G. Dimakis, "Coded caching with linear subpacketization is possible using ruzsa-szeméredi graphs," in *2017 IEEE International Symposium on Information Theory (ISIT)*, 2017, pp. 1237–1241.

[18] M. Cheng, Q. Yan, X. Tang, and J. Jiang, "Coded caching schemes with low rate and subpacketizations," *arXiv preprint arXiv:1703.01548*, 2017.

[19] M. Cheng, J. Jiang, Q. Yan, and X. Tang, "Constructions of coded caching schemes with flexible memory size," *IEEE Transactions on Communications*, vol. 67, no. 6, pp. 4166–4176, 2019.

[20] S. A. Saberali, L. Lampe, and I. F. Blake, "Decentralized coded caching without file splitting," *IEEE Transactions on Wireless Communications*, vol. 18, no. 2, pp. 1289–1303, Feb 2019.

[21] M. Cheng, J. Li, X. Tang, and R. Wei, "Linear coded caching scheme for centralized networks," *IEEE Transactions on Information Theory*, vol. 67, no. 3, pp. 1732–1742, 2021.

[22] J. Michel and Q. Wang, "Placement delivery arrays from combinations of strong edge colorings," *IEEE Transactions on Communications*, vol. 68, no. 10, pp. 5953–5964, 2020.

[23] W. Song, K. Cai, and L. Shi, "Some new constructions of coded caching schemes with reduced subpacketization," *ArXiv*, vol. abs/1908.06570, 2019.

[24] M. Cheng, J. Wang, and X. Zhong, "A unified framework for constructing centralized coded caching schemes," *ArXiv*, vol. abs/1908.05865, 2019.

[25] S. Agrawal, K. V. Sushena Sree, and P. Krishnan, "Coded caching based on combinatorial designs," in *2019 IEEE International Symposium on Information Theory (ISIT)*, 2019, pp. 1227–1231.

[26] M. Ji, G. Caire, and A. F. Molisch, "Fundamental limits of caching in wireless d2d networks," *IEEE Transactions on Information Theory*, vol. 62, no. 2, pp. 849–869, Feb 2016.

[27] S. Li, M. A. Maddah-Ali, Q. Yu, and A. S. Avestimehr, "A fundamental tradeoff between computation and communication in distributed computing," *IEEE Transactions on Information Theory*, vol. 64, no. 1, pp. 109–128, Jan 2018.

[28] N. Naderializadeh, M. A. Maddah-Ali, and A. S. Avestimehr, "Fundamental limits of cache-aided interference management," *IEEE Transactions on Information Theory*, vol. 63, no. 5, pp. 3092–3107, May 2017.

[29] K. Wan, D. Tuninetti, and P. Piantanida, "On caching with more users than files," in *2016 IEEE International Symposium on Information Theory (ISIT)*, 2016, pp. 135–139.

[30] J. Hirschfeld, *Projective Geometries Over Finite Fields. Oxford Mathematical Monographs.* Oxford University Press New York, 1998.

[31] L. Takács, "Some asymptotic formulas for lattice paths," *Journal of Statistical Planning and Inference*, vol. 14, pp. 123–142, 05 1986.

[32] Q. Yan, X. Tang, and Q. Chen, "Placement delivery array and its applications," in *2018 IEEE Information Theory Workshop (ITW)*, 2018, pp. 1–5.

[33] J. Dean and S. Ghemawat, "Mapreduce: Simplified data processing on large clusters," in *OSDI'04: Sixth Symposium on Operating System Design and Implementation*, San Francisco, CA, 2004, pp. 137–150.

[34] Q. Yan, S. Yang, and M. Wigger, "Storage, computation, and communication: A fundamental tradeoff in distributed computing," in *2018 IEEE Information Theory Workshop (ITW)*, 2018, pp. 1–5.

[35] Q. Yan, M. Wigger, S. Yang, and X. Tang, "A fundamental storage-communication tradeoff for distributed computing with straggling nodes," *IEEE Transactions on Communications*, vol. 68, no. 12, pp. 7311–7327, 2020.

[36] E. Lampiris and P. Elia, "Adding transmitters dramatically boosts coded-caching gains for finite file sizes," *IEEE Journal on Selected Areas in Communications*, vol. 36, no. 6, pp. 1176–1188, June 2018.

[37] Z. Bar-Yossef, Y. Birk, T. S. Jayram, and T. Kol, "Index coding with side information," *IEEE Transactions on Information Theory*, vol. 57, no. 3, pp. 1479–1494, March 2011.

[38] Q. Yu, M. A. Maddah-Ali, and A. S. Avestimehr, "The exact rate-memory tradeoff for caching with uncoded prefetching," *IEEE Transactions on Information Theory*, vol. 64, no. 2, pp. 1281–1296, Feb 2018.

[39] M. Salehi, A. Tolli, S. P. Shariatpanahi, and J. Kaleva, "Subpacketization-rate trade-off in multi-antenna coded caching," in *2019 IEEE Global Communications Conference (GLOBECOM)*, 2019, pp. 1–6.

[40] M. Salehi, A. Tölli, and S. P. Shariatpanahi, "A multi-antenna coded caching scheme with linear subpacketization," in *ICC 2020 - 2020 IEEE International Conference on Communications (ICC)*, 2020, pp. 1–6.

[41] A. Goel, M. Kapralov, and S. Khanna, "Perfect matchings in o(n log n) time in regular bipartite graphs," in *Proceedings of the Forty-second ACM Symposium on Theory of Computing*, ser. STOC '10. New York, NY, USA: ACM, 2010, pp. 39–46.