

Delivery Time Minimization in Edge Caching: Synergistic Benefits of Subspace Alignment and Zero Forcing

Jaber Kakar*, Alaa Alameer*, Anas Chaaban[†], Aydin Sezgin* and Arogyaswami Paulraj[‡]

*Institute of Digital Communication Systems, Ruhr-Universität Bochum, Germany

[†]Communication Theory Lab, King Abdullah University of Science and Technology, Thuwal, Saudi Arabia

[‡]Information Systems Laboratory, Stanford University, CA, USA

Email: {jaber.kakar, alaa.alameerahmad, aydin.sezgin}@rub.de, anas.chaaban@kaust.edu.sa, apaulraj@stanford.edu

Abstract—An emerging trend of next generation communication systems is to provide network edges with additional capabilities such as additional storage resources in the form of caches to reduce file delivery latency. To investigate this aspect, we study the fundamental limits of a cache-aided wireless network consisting of one central base station, M transceivers and K receivers from a latency-centric perspective. We use the normalized delivery time (NDT) to capture the per-bit latency for the worst-case file request pattern at high signal-to-noise ratios (SNR), normalized with respect to a reference interference-free system with unlimited transceiver cache capabilities. For various special cases with $M = \{1, 2\}$ and $K = \{1, 2, 3\}$ that satisfy $M+K \leq 4$, we establish the optimal tradeoff between cache storage and latency. This is facilitated through establishing a novel converse (for arbitrary M and K) and an achievability scheme on the NDT. Our achievability scheme is a synergistic combination of multicasting, zero-forcing beamforming and interference alignment.

I. INTRODUCTION

In the last decade, mobile usage in wireless networks has shifted from being connection-centric driven (e.g., phone calls) to content-centric (e.g., HD video) behaviors. In this context, integrating content caching in heterogeneous networks (HetNet) represents a viable solution for highly content-centric, next generation (5G) mobile networks. Specifically, when caching the most popular contents in HetNet *edge nodes*, e.g., eNBs and relays, alleviates backhaul traffic, reduces latency and ameliorates quality of service of mobile users. Thus, it is expected that future networks will be heterogeneous in nature, vastly deploying relay nodes (RN) (e.g., fixed RNs in LTE-A [1] or mobile RNs in form of drones [2], [3]) endowed with content cache capabilities.

A simplistic HetNet modeling this aspect is shown in Fig. 1. In this model, M RNs act as cache-aided transceivers. Thus, aspects of both transmitter and receiver caching in RNs is captured through this network model enabling a low *delivery time* of requested files by M RNs and K user equipments (UE).¹ These terms refer to the timing overhead required to satisfy all file demands of requesting nodes in the network. In this work, we are interested in *completely* characterizing the fundamental delivery time cache memory trade-off of this particular network for specific instances of M and K .

In prior work, it was shown that both receiver (Rx) and transmitter (Tx) caching can offer significant latency reduction. Rx caching was first studied in [4] for a shared link with one server and multiple cache-enabled receivers. The authors show

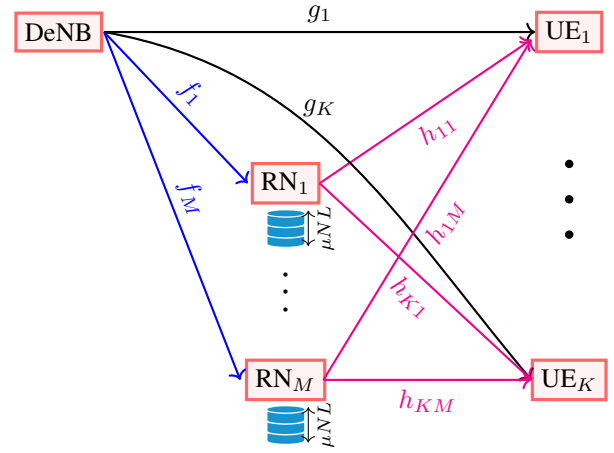


Fig. 1: A transceiver cache-aided HetNet consisting of one DeNB, M RNs and K UEs.

that appropriate caching of popular content facilitates multicast opportunities and consequently reduces latency. On the other hand, the impact of Tx caching on latency has mainly been investigated by analyzing the inverse degrees-of-freedom (DoF) metric of Gaussian interference networks [5]. To this end, the authors of [6] developed an interference alignment scheme characterizing the metric as a function of the cache storage size for a 3-user Gaussian interference network. The caches are prefetched to allow transmitter cooperation so that interference coordination techniques are applicable. The first lower bounds on the inverse DoF were developed in [7] for a network with an arbitrary number of edge nodes and users. With these bounds, the optimality of schemes presented in [6] for certain regimes of cache sizes was shown under uncoded prefetching of the cached content. Extensions of this work include the characterization of the latency-memory tradeoff in cloud and cache-assisted networks for equally and non-equally strong wireless links in [8] and [9], [10], respectively. Recently, in two new lines of research, the effect of Tx-Rx caching at *distinct* nodes [11] and transceiver caching [12] on the latency were investigated. This paper focuses on the latter.

In this paper, we study the fundamental limits on the delivery time for a *transceiver* cache-aided HetNet consisting of one donor eNBs (DeNB), M transceivers and K users. We measure the performance through a latency-centric metric known as the *normalized delivery time per bit* (NDT) (cf. formal definition of NDT in Eq. (5) in Section II). This metric, first introduced

¹We use the words *delivery time* and *latency* interchangeably.

in [7], indicates the worst-case per-bit latency incurred in the wireless network with respect to a reference interference-free system without cache capacity restrictions. Similarly to the DoF, it is a high signal-to-noise ratio (SNR) metric. The main contributions of this paper are as follows:

- We develop a novel class of information theoretic lower bounds on the NDT under the assumption of perfect channel state information (CSI) and uncoded prefetching of the cached content.
- We completely characterize the NDT-cache memory tradeoff for the settings of (a) $M = 1$ RNs and $K = \{1, 2, 3\}$ UEs and (b) $M = 2$ RNs and $K = \{1, 2\}$ UEs. To this end, we establish NDT-optimal schemes that synergistically design precoders facilitating zero-forcing (ZF) beamforming, multicasting and interference alignment. Our schemes are optimal for both time-variant and invariant channels requiring *finite* signal dimensions (time, frequency, etc.). Further, we determine the optimal schemes for the extremal cases of no caching and full caching.
- Along with our results, we discuss the relationship between (sum) DoF and NDT. To this end, we assess the results from both a rate (e.g., DoF), and latency (e.g., NDT) perspective.

Notation: For any two integers a and b with $a \leq b$, we define $[a : b] \triangleq \{a, a + 1, \dots, b\}$. When $a = 1$, we simply write $[b]$ for $\{1, \dots, b\}$. The superscript $(\cdot)^\dagger$ represents the transpose of a matrix. Furthermore, we define the function $(x)^+ \triangleq \max\{0, x\}$ and the *modified* modular operator $c = a \text{ MOD } \{b\}$ for integers a and b as $c = a$ if $a \leq b$ and $c = a \bmod b$ if $a > b$.

II. SYSTEM MODEL

We study the *downlink* of a transceiver cache-aided HetNet as shown in Fig. 1. The network consists of M causal full-duplex RNs and a donor eNB (DeNB) which serves K UEs with its desired content over a shared wireless channel. Simultaneously, each RN also requests information from the DeNB. At every transmission interval, we assume that RNs and UEs request files from the set \mathcal{W} of N popular files, whose elements are all of L bits in size. The transmission interval terminates when the requested files have been delivered. The system model, notation and main assumptions for a *single* transmission interval are summarized as follows:

- Let $\mathcal{W} = \{W_1, \dots, W_N\}$ denote the library of popular files, where each file W_n is of size L bits. Each file W_n is chosen uniformly at random from $[2^L]$. UEs and RNs request files W_{d_u} , $\forall u \in [K]$, and W_{d_r} , $\forall r \in [K + 1 : M + K]$, from the library \mathcal{W} , respectively. The demand vector $\mathbf{d} = (d_1, \dots, d_{M+K}) \in [N]^{M+K}$ denotes the request pattern of RNs and UEs.
- The RNs are endowed with a cache capable of storing μNL bits, where $\mu \in [0, 1]$ corresponds to the *fractional cache size*. It denotes how much content can be stored at each RN relatively to the entire library \mathcal{W} .
- The DeNB has access to all N popular files of \mathcal{W} .
- Global CSI at time instant t is summarized by the channel vectors $\mathbf{f}[t] = \{f_m[t]\}_{m=1}^M \in \mathbb{C}^M$ and $\mathbf{g}[t] = \{g_k[t]\}_{k=1}^K \in \mathbb{C}^K$ and the channel matrix $\mathbf{H}[t] = \{h_{km}[t]\}_{k=1, m=1}^{K, M} \in \mathbb{C}^{K \times M}$. Here, f_m and g_k represent the complex channel coefficients from DeNB to RN $_m$ and UE $_k$, respectively, while h_{km} is the channel from RN $_m$ to UE $_k$. We assume that all channel coefficients are assumed to be drawn i.i.d. from a continuous distribution.

Communication over the wireless channel occurs in two consecutive phases, (a) *placement phase* followed by (b) *delivery phase*. These are detailed next, along with the key performance metric termed as *normalized delivery time per bit* (NDT).

a) Placement phase: During this phase, every RN is given full access to the database of N files. The cached content at RN $_m$ is generated through its individual caching function.

Definition 1. (Caching function) RN $_m$, $\forall m = 1, \dots, M$, maps each file $W_n \in \mathcal{W}$ to its local *file cache content*

$$S_{m,n} = \phi_{m,n}(W_n), \quad \forall n = 1, \dots, N.$$

All $S_{m,n}$ are concatenated to form the total cache content

$$S_m = (S_{m,1}, S_{m,2}, \dots, S_{m,N})$$

at RN $_m$.

Hereby, due to the assumption of symmetry in caching, the entropy $H(S_{m,n})$ of each component $S_{m,n}$, $n = 1, \dots, N$, is upper bounded by $\mu NL/N = \mu L$. The definition of the caching function presumes that every file W_i is subjected to individual caching functions. Thus, permissible caching policies allow for intra-file coding but avoid coding across files known as inter-file coding. Moreover, the caching policy is typically kept fixed over long transmission intervals. Thus, it is indifferent to the UEs request pattern and of channel realizations.

b) Delivery phase: In this phase, a transmission policy at DeNB and all RNs is applied to satisfy the given requests \mathbf{d} under the current channel realizations \mathbf{f} , \mathbf{g} and \mathbf{H} . Throughout the remaining definitions, we denote the number of channel uses required to satisfy all file demands by T .

Definition 2. (Encoding functions) The DeNB encoding function at time instant $t \in [T]$

$$\psi_s^{[t]} : [2^{NL}] \times [N]^{M+K} \times \mathbb{C}^{Mt} \times \mathbb{C}^{Kt} \times \mathbb{C}^{Kt \times M} \rightarrow \mathbb{C}$$

determines the DeNBs transmission signal $x_s[t] = \psi_s^{[t]}(\mathcal{W}, \mathbf{d}, \mathbf{f}_{t=1}^t, \mathbf{g}_{t=1}^t, \mathbf{H}_{t=1}^t)$ subjected to an average power constraint of P . The encoding function of the causal full-duplex RN $_m$ at time instant $t \in [T]$ is defined by

$$\psi_{r,m}^{[t]} : [2^{\mu NL}] \times \mathbb{C}^{t-1} \times [N]^{M+K} \times \mathbb{C}^{Mt} \times \mathbb{C}^{Kt} \times \mathbb{C}^{Kt \times M} \rightarrow \mathbb{C},$$

which determines the codeword $x_{r,m}[t] = \psi_{r,m}^{[t]}(S_m, \mathbf{y}_{r,m}^{t-1}, \mathbf{d}, \mathbf{f}_{t=1}^t, \mathbf{g}_{t=1}^t, \mathbf{H}_{t=1}^t)$ while satisfying the average power constraint given by the parameter P .

Hereby, the codewords $x_s[t]$ and $x_{r,m}[t]$ are transmitted over $t \in [T]$ channel uses. For any time instant t , $\psi_{r,m}^{[t]}$ accounts for the simultaneous reception and transmission through incoming and outgoing wireless links at RN $_m$. To be specific, at the t -th channel use the encoding function $\psi_{r,m}^{[t]}$ maps the cached content S_m , the received signal $\mathbf{y}_{r,m}^{t-1}$ (see Eq. (2)), the demand vector \mathbf{d} and global CSI to the codeword $x_{r,m}[t]$.

After transmission, the received signals at UE $_k$ is given by

$$y_{u,k}[t] = g_k[t]x_s[t] + \sum_{m=1}^M h_{km}[t]x_{r,m}[t] + z_{u,k}[t], \forall t \in [T], \quad (1)$$

where $z_{u,k}[t]$ denotes complex i.i.d. Gaussian noise of zero mean and unit power. The received signal at RN $_m$ is given by

$$y_{r,m}[t] = f_m[t]x_s[t] + z_{r,m}[t], \forall t \in [T], \quad (2)$$

where $z_{r,m}[t]$ is additive zero mean, unit-power i.i.d. Gaussian noise. The desired files are decoded using the following functions.

Definition 3. (Decoding functions) The decoding operation at UE $_k$ follows the mapping

$$\eta_{u,k} : \mathbb{C}^T \times [N]^{M+K} \times \mathbb{C}^{MT} \times \mathbb{C}^{KT} \times \mathbb{C}^{KT \times M} \rightarrow [2^L].$$

to provide an estimate $\hat{W}_{d_k} = \eta_{u,k}(\mathbf{y}_{u,k}^T, \mathbf{d}, \mathbf{f}_{t=1}^T, \mathbf{g}_{t=1}^T, \mathbf{H}_{t=1}^T)$ of the requested file W_{d_k} . In contrast to decoding at UE $_k$, all RNs explicitly leverage their cached content according to

$$\eta_{r,m} : \mathbb{C}^T \times [2^{\mu NL}] \times [N]^{M+K} \times \mathbb{C}^{MT} \times \mathbb{C}^{KT} \times \mathbb{C}^{KT \times M} \rightarrow [2^L]$$

to generate $\hat{W}_{d_r} = \eta_{r,m}(\mathbf{y}_{r,m}^T, S_m, \mathbf{d}, \mathbf{f}_{t=1}^T, \mathbf{g}_{t=1}^T, \mathbf{H}_{t=1}^T)$ as an estimate of the requested file W_{d_r} .

A proper choice of caching, encoding and decoding functions that satisfy the reliability condition; that is, the worst-case error probability

$$P_e = \max_{\mathbf{d} \in [N]^{M+K}} \max_{j \in [M+K]} \mathbb{P}(\hat{W}_{d_j} \neq W_{d_j}) \quad (3)$$

approaches 0 as $L \rightarrow \infty$, is called a *feasible policy*. Now we are ready to define the delivery time per bit and its normalized version.

Definition 4. (Delivery time per bit [7]) The delivery time per bit (DTB) for a given request pattern \mathbf{d} and channel realization \mathbf{f}, \mathbf{g} and \mathbf{H} is defined as

$$\Delta(\mu, P) = \max_{\mathbf{d} \in [N]^{M+K}} \limsup_{L \rightarrow \infty} \frac{\mathbb{E}[T(\mathbf{d}, \mathbf{f}, \mathbf{g}, \mathbf{H})]}{L}, \quad (4)$$

where the expectation is over the channel realizations.

In the definition above, T represents the completion or delivery time [13]. The normalization of the expected delivery time by the file size L gives insight about the per bit-latency. In this context, the DTB measures the per-bit latency, i.e., the latency incurred per-bit when transmitting the requested files through the wireless channel, within a single transmission interval for the *worst-case* request pattern of RNs and UEs as $L \rightarrow \infty$. The DTB depends on the fractional cache size μ and the power level P .

In analogy to the degrees-of-freedom metric, the normalized delivery time per bit (NDT) is a high-SNR metric that relates the DTB to that of a point-to-point reference system.

Definition 5. (Normalized delivery time [7]) The NDT is defined as

$$\delta(\mu) = \lim_{P \rightarrow \infty} \frac{\Delta(\mu, P)}{1/\log(P)}. \quad (5)$$

The minimum NDT $\delta^*(\mu)$ is the infimum in NDT of all feasible policies.

The NDT compares the *delivery time per bit* achieved by the feasible coding scheme for the worst-case demand scenario to that of a baseline interference-free system in the high SNR regime. The achievable scheme, on the one hand, allows for reliable transmission of one file of L bits to a single Rx on average in $\mathbb{E}[T(\mathbf{f}, \mathbf{g}, \mathbf{H})]$ channel uses, i.e., 1 bit in $\mathbb{E}[T(\mathbf{f}, \mathbf{g}, \mathbf{H})]/L$ channel uses. The baseline system (e.g., a point-to-point channel), on the other hand, can transmit $\log(P)$ bits to a single Rx in one channel use, i.e., 1 bit in $1/\log(P)$ channel uses. Therefore, the resulting NDT $\delta(\mu)$ indicates that the worst-case delivery time for one bit of the cache-aided

network at fractional cache size μ is $\delta(\mu)$ times larger than the time needed by the baseline system.

From [6, Lemma 1], it readily follows that the NDT is a convex function in μ . This means that a cache-aided network shown in Fig. 1 operating at fractional cache size $\mu = \alpha\mu_1 + (1-\alpha)\mu_2$ for any $\alpha \in [0, 1]$ achieves at most an NDT equal to the *convex combination* $\alpha\delta(\mu_1) + (1-\alpha)\delta(\mu_2)$ through applying known feasible schemes applicable at fractional cache sizes μ_1 and μ_2 on distinct α and $1-\alpha$ -fractions of the files, respectively. This strategy is known as *memory sharing*.

III. LOWER BOUND (CONVERSE) ON NDT

For a given worst-case demand pattern \mathbf{d} ; that is all K UEs and M RNs request *distinct* files W_{d_j} ($d_j \neq d_\ell, j \neq \ell$), and given channel realizations \mathbf{f}, \mathbf{g} and \mathbf{H} , we obtain a lower bound on the delivery time $T = T(\mathbf{d}, \mathbf{f}, \mathbf{g}, \mathbf{H})$, and therefore ultimately on the NDT, of any *feasible* scheme. Note that $K + M$ distinct files W_{d_k} are available if there are at least as many files in the library, i.e., $N \geq K + M$. Without loss of generality, we assume that the requested files by the K UEs are $W_{[1:K]} = (W_1, W_2, \dots, W_K)$ and of the M RNs $W_{[K+1:K+M]} = (W_{K+1}, W_{K+2}, \dots, W_{K+M})$.

The key idea in establishing the lower bound on the NDT is that $K + \ell$ requested files, comprising of all K files $W_{[1:K]}$ requested by the UEs and ℓ files desired by a subset of ℓ RNs (out of M RNs), e.g., $W_{[K+1:K+\ell]}$, can be retrieved in the high SNR regime from

- s output signals of the UEs, e.g., $\mathbf{y}_{u,[1:s]}^T$ for $1 \leq s \leq \min\{M+1, K\}$, and
- ℓ cached contents of ℓ RNs, e.g., $S_{[1:\ell]}$, where $\bar{s} \leq \ell \leq M$ and $\bar{s} = M + 1 - s$.

We note that since $s + \ell \geq M + 1$ holds, we are able to reconstruct all $M + 1$ transmit signals ($x_s[t]$ and $x_{r,m}[t], m \in [M]$) at all T time instants of the delivery phase within bounded noise. The intuition behind the bound follows from [7], [12]. Applying standard information-theoretic bounding techniques, results in the following Lemma.

Lemma 1. For the transceiver cache-aided network with one DeNB, M RNs each endowed with a cache of fractional cache size $\mu \in [0, 1]$, K UEs and a file library of $N \geq M + K$ files, the NDT is lower bounded under perfect CSI at all nodes by

$$\delta^*(\mu) \geq \max \left\{ 1, \max_{\substack{\ell \in [\bar{s}: M], \\ s \in [\min\{M+1, K\}]}} \delta_{\text{LB}}(\mu, \ell, s) \right\}, \quad (6)$$

where $\bar{s} = M + 1 - s$ and

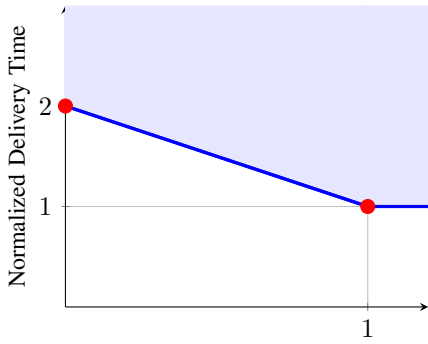
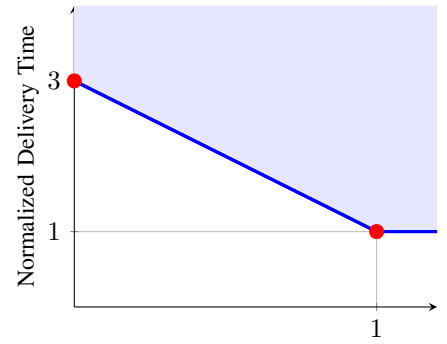
$$\delta_{\text{LB}}(\mu, \ell, s) = \frac{K + \ell - \mu(\bar{s}(K - s + \frac{\bar{s}-1}{2}) + \frac{\ell}{2}(\ell + 1))}{s}. \quad (7)$$

Proof. The key idea behind the proof is provided in the previous paragraph. Details are omitted for the sake of brevity. ■

IV. ACHIEVABILITY FOR SOME SPECIAL CASES

First let us consider two special corner points at fractional cache sizes $\mu = 0$ and $\mu = 1$. These are the cases where the RN has either *zero-cache* ($\mu = 0$) or *full-cache* ($\mu = 1$) capabilities. We now expound the optimal NDT for these two points.

Lemma 2. For the transceiver cache-aided network with one DeNB, M RNs each endowed with a cache of fractional cache

(a) NDT for $M = K = 1$ (b) NDT for $M = 1, K = 2$ Fig. 2: Optimal NDT as a function of μ for $M = 1$ and $K \leq 2$.

size μ , K UEs and a file library of $N \geq M + K$ files, the optimal NDT is

$$\delta^*(\mu) = K + M \quad \text{for } \mu = 0, \quad (8)$$

achievable via DeNB broadcasting to M RNs and K UEs, and

$$\delta^*(\mu) = \max \left\{ \frac{K}{M+1}, 1 \right\} \quad \text{for } \mu = 1, \quad (9)$$

achievable via zero-forcing beamforming for an $(M + 1, K)$ MISO² broadcast channel.

Proof. For the proof, it suffices to find a cache transmission policy that matches the lower bound in Lemma 1 for $\mu = 0$ and $\mu = 1$, respectively. On the one hand, if $\mu = 0$, we note that $\delta_{\text{LB}}(0, M, 1) = K + M$. On the other hand, if $\mu = 1$, we observe that $\delta_{\text{LB}}(1, 0, M + 1) = K/(M+1)$ if $M + 1 \leq K$ and $\delta_{\text{LB}}(1, \ell, s) < 1$ if $M + 1 > K$. Next, we consider the achievability at $\mu = 0$ and $\mu = 1$. For these two fractional cache sizes, the network in Fig. 1 reduces to a SISO broadcast channel (BC) with $K + M$ users for $\mu = 0$ and an $(M + 1, K)$ MISO broadcast channel for $\mu = 1$. The approximate *per-user* rate (neglecting $o(\log(P))$ bits) for these two channels are known to be $\frac{1}{(K+M)} \log(P)$ (achievable through unicasting each user's message) and $\frac{1}{K} \min\{M + 1, K\} \log(P)$ (achievable through zero-forcing beamforming), respectively. Equivalently, each user needs the reciprocal *per-user* rate of signaling dimensions (e.g., channel uses in time or frequency) to retrieve one desired bit. Thus, the approximate DTB becomes, respectively, $\frac{(K+M)}{\log(P)}$ and $\frac{K}{\min\{M+1, K\} \log(P)}$. Normalizing the delivery time per bit by the point-to-point reference DTB $\frac{1}{\log(P)}$ generates the NDTs $K+M$ and $\max\{K/(M+1), 1\}$. This establishes the NDT-optimality at these fractional cache sizes. ■

Remark 1. From Lemma 2, we infer that the caching problem for the system illustrated in Fig. 1 establishes the behavior of the network in terms of delivery time between the two extremes – SISO BC with $K + M$ users and an $(M + 1, K)$ MISO BC. This analysis will reveal what kind of schemes other than simple unicasting and zero-forcing will be optimal for $0 < \mu < 1$.

Now, we move to special cases of the system where $M = \{1, 2\}$ to provide a *complete* characterization of the NDT-memory trade-off. We will primarily focus on $M = 1$ for the sake of brevity.

²In MISO broadcast channels, we use the notation, (a, b) for integers a and b to denote a broadcast channel with a transmit antennas and b single antenna receivers.

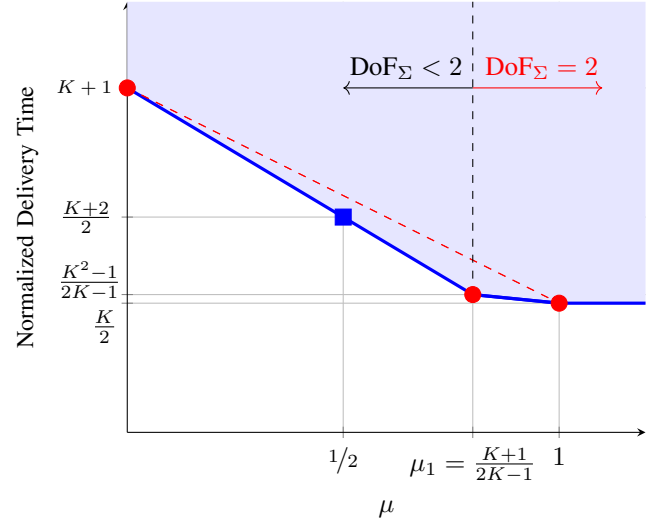


Fig. 3: NDT lower bound for $M = 1$ and $K \geq 3$. For $K = 3$, this line is in fact achievable. The dashed line shows the achievable NDT of a *suboptimal* time-sharing based unicasting-zero-forcing scheme, which is optimal for $M = 1$ and $K \leq 2$.

A. Achievability for $M = 1$

The lower bounds on the NDT are obtained from Lemma 1 by setting $M = 1$, yielding

$$\delta^*(\mu) \geq \begin{cases} \delta_{\text{LB}}(\mu, 1, 1) = K + 1 - \mu K & \text{for } K \geq 1 \\ \delta_{\text{LB}}(\mu, 1, 2) = \frac{K+1-\mu}{2} & \text{for } K \geq 2 \\ \delta_{\text{LB}}(\mu, 0, 2) = \frac{K}{2} & \text{for } K \geq 2 \end{cases} \quad (10)$$

For $K \leq 2$, the optimal NDT-memory curves are shown in Fig. 2 for both $K = 1$ (cf. Fig. 2a) and $K = 2$ (cf. Fig. 2b). The achievability at the corner points (marked by circles in Fig. 2) at $\mu = 0$ and $\mu = 1$ readily follow from Lemma 2. Intermediary points are achievable through memory sharing. Thus, the optimal NDT for $K \leq 2$ and $M = 1$ becomes

$$\delta^*(\mu) = K + 1 - \mu K. \quad (11)$$

This result is in agreement with our prior work [12]. For $K > 2$, on the other hand, the lower bound in (10) simplifies to

$$\delta^*(\mu) \geq \begin{cases} K + 1 - \mu K & \text{for } 0 \leq \mu \leq \mu_1 \\ \frac{K+1-\mu}{2} & \text{for } \mu_1 \leq \mu \leq 1 \end{cases}, \quad (12)$$

where $\mu_1 = \frac{K+1}{2K-1}$ as shown in Fig. 3. In order to show the tightness of the lower bound, we have to focus on the corner point $(\frac{K+1}{2K-1}, \frac{K^2-1}{2K-1})$. Interestingly, if this point is achievable we make two observations.

	RN & DeNB				DeNB only
File W_1	$\eta_{1,1}$	$\eta_{1,2}$	$\eta_{1,3}$	$\eta_{1,4}$	$\eta_{1,5}$
File W_2	$\eta_{2,1}$	$\eta_{2,2}$	$\eta_{2,3}$	$\eta_{2,4}$	$\eta_{2,5}$
File W_3	$\eta_{3,1}$	$\eta_{3,2}$	$\eta_{3,3}$	$\eta_{3,4}$	$\eta_{3,5}$
File W_4	$\eta_{4,1}$	$\eta_{4,2}$	$\eta_{4,3}$	$\eta_{4,4}$	$\eta_{4,5}$

ZF symbols

Fig. 4: Requested files by $K = 3$ users and $M = 1$ RNs and the availability illustrated by the symbols transmitted from the DeNB only or from both at the DeNB and the RN.

	ZF map		
UE ₁	$\eta_{2,1}$	$\eta_{2,2}$	$\eta_{3,3}$
UE ₂	$\eta_{3,1}$	$\eta_{3,2}$	$\eta_{1,3}$
UE ₃	$\eta_{1,1}$	$\eta_{1,2}$	$\eta_{2,3}$

Fig. 5: Map that assigns which symbol to zero-force at which receiver.

Remark 2. For increasing K , the point converges to $(\frac{1}{2}, \frac{K}{2})$. This fact shows that as K increases, a fractional memory size of $\mu_1 \rightarrow 1/2$ suffices in already attaining the lowest attainable NDT of $K/2$. The point $(\frac{1}{2}, \frac{K+2}{2})$ (marked by a square in Fig. 3) is achievable by leveraging interference alignment techniques for a $2 \times K$ X-channel [14] and unicasting uncached information about the RNs desired file from the DeNB.

Remark 3. If (12) is achievable, one can see that the per-user DoF of the K UEs is $\frac{2K-1}{K^2-1}$ and that of the RN $\frac{K-2}{K^2-1}$. The resulting sum DoF thus becomes $\frac{K(2K-1)}{K^2-1} + \frac{K-2}{K^2-1} = 2$. Thus, at fractional cache sizes greater than μ_1 , the sum DoF remains 2. This is shown in Fig. 3.

So far, we were able to establish the achievability for this corner point for $K = 3$ UEs only. The *generalization* to arbitrary numbers of UEs is still an *open problem*. In the sequel, we will illustrate the achievability of the corner point $(\frac{K+1}{2K-1}, \frac{K^2-1}{2K-1}) = (\frac{4}{5}, \frac{8}{5})$ for $K = 3$.

Assume without loss of generality $N = 4$ and that the UEs request files W_1, W_2 and W_3 while the RN is interested in file W_4 . According to Fig. 4, all files are comprised of 5 symbols, i.e., the i -th file is composed of symbols $\eta_{i,1}, \eta_{i,2}, \eta_{i,3}, \eta_{i,4}$ and $\eta_{i,5}$. All these symbols are available at the DeNB. However, as far as the RN is concerned, only the first four symbols of all files are locally available in its cache. Since, the RN is interested in file W_4 and it knows $\eta_{4,1}, \eta_{4,2}, \eta_{4,3}$ and $\eta_{4,4}$, the only missing symbol it desires is $\eta_{4,5}$. Thus, the transmission policy has to be designed such that DeNB and RN are involved in sending *all* symbols of files W_1, W_2 and W_3 as well as $\eta_{4,5}$. These are in total 16 information symbols. The transmission strategy will exploit the correlation that arises between the availability of shared symbols at RN and DeNB by leveraging zero-forcing (ZF) opportunities while *simultaneously* facilitating (subspace) interference alignment (IA) at the UEs. This is why our scheme (as shown in Fig. 4) only zero-forces symbols $\eta_{1,1}, \eta_{1,2}, \eta_{1,3}, \eta_{2,1}, \eta_{2,2}, \eta_{2,3}$ and $\eta_{3,1}, \eta_{3,2}, \eta_{3,3}$. Symbols $\eta_{1,4}, \eta_{2,4}$ and $\eta_{3,4}$ are *not* zero-forced but are instead used to enable alignment⁴

³This is due to the fact that $L(1 - \mu_1)$ symbols are uncached and are conveyed in $T = K^2 - 1$ channel uses, with $L = 2K - 1$. This constitutes the aforementioned DoF value.

⁴IA is facilitated by the fact that the DeNB does *not* transmit these symbols (even though it knows them). Thus, effectively, the DeNB does not need to be aware of $\eta_{i,4}, i \in [N]$.

amongst others with $\eta_{4,5}$ at the UEs. The map that assigns which symbol is zero-forced at which UE is given in Figure 5. To this end, DeNB and RN form their transmit signals according to

$$x_s[t] = \sum_{i=1}^3 \sum_{j=1, j \neq 4}^5 \nu_{\eta_{i,j}}[t] \eta_{i,j} + \nu_{\eta_{4,5}}[t] \eta_{4,5}, \quad (13)$$

$$x_r[t] = \sum_{i=1}^3 \sum_{j=1}^4 \beta_{\eta_{i,j}}[t] \eta_{i,j}, \quad (14)$$

$\forall t \in [T]$ for $T = 8$, respectively. The complex scalars $\nu_{\eta_{i,j}}[t]$ and $\beta_{\eta_{i,j}}[t]$ are precoders for symbol $\eta_{i,j}$ originating from DeNB and RN at time instant t , respectively. They are chosen such that both ZF and IA at the UEs become feasible. According to the ZF map of Fig. 5, the ZF conditions at UE _{k} , $k \in [K] = [3]$, become

$$\nu_{\eta_{(k+1) \bmod K,1}}[t] g_k[t] + \beta_{\eta_{(k+1) \bmod K,1}}[t] h_{k1}[t] = 0, \quad (15a)$$

$$\nu_{\eta_{(k+1) \bmod K,2}}[t] g_k[t] + \beta_{\eta_{(k+1) \bmod K,2}}[t] h_{k1}[t] = 0, \quad (15b)$$

$$\nu_{\eta_{(k+2) \bmod K,3}}[t] g_k[t] + \beta_{\eta_{(k+2) \bmod K,3}}[t] h_{k1}[t] = 0. \quad (15c)$$

Simultaneously, we design the precoding scalars such that the interference at each UE is aligned into a three-dimensional signal space. (The remaining 5 dimensions are reserved for the 5 symbols of the desired file.) The interference graph in Fig. 6 shows which symbols align with each other at which UE. This graph consists of 3 layers. In the first layer, two symbols, namely $\eta_{4,5}$ and $\eta_{1,4}, \eta_{2,4}$ or $\eta_{3,4}$ align at the three UEs. At layers two and three, on the other hand, three symbols align per UE. Symbols $\eta_{1,4}, \eta_{2,4}$ and $\eta_{3,4}$ link layers 1 and 2, while $\eta_{1,5}, \eta_{2,5}$ and $\eta_{3,5}$ connect layers 2 and 3. In analogy to the graph in Fig. 6, the alignment conditions at UE _{k} can be written as

$$\nu_{\eta_{4,5}}[t] g_k[t] = \beta_{\eta_{(k+1) \bmod K,4}}[t] h_{k1}[t] \quad (16)$$

for Layer 1,

$$\begin{aligned} & \beta_{\eta_{(k+2) \bmod K,4}}[t] h_{k1}[t] \\ &= \beta_{\eta_{(k+2) \bmod K,2}}[t] h_{k1}[t] + \nu_{\eta_{(k+2) \bmod K,2}}[t] g_k[t] \\ &= \nu_{\eta_{(k+1) \bmod K,5}}[t] g_k[t] \end{aligned} \quad (17)$$

for Layer 2, and

$$\begin{aligned} & \nu_{\eta_{(k+2) \bmod K,5}}[t] g_k[t] \\ &= \beta_{\eta_{(k+2) \bmod K,1}}[t] h_{k1}[t] + \nu_{\eta_{(k+2) \bmod K,1}}[t] g_k[t] \\ &= \beta_{\eta_{(k+1) \bmod K,3}}[t] h_{k1}[t] + \nu_{\eta_{(k+1) \bmod K,3}}[t] g_k[t] \end{aligned} \quad (18)$$

for Layer 3. Under the given ZF and IA conditions (cf. (15) and (16)–(18)), the precoders are functions of the channels $\mathbf{g}[t]$ and $\mathbf{H}[t]$. We fix the precoder for symbol $\eta_{4,5}$ to

$$\nu_{\eta_{4,5}}[t] = j_{13}[t] j_{23}[t] j_{33}[t] g_1[t] g_2[t] g_3[t] h_{11}[t] h_{21}[t] h_{31}[t], \quad (19)$$

where

$$j_{13}[t] = g_2[t] h_{31}[t] - g_3[t] h_{21}[t], \quad (20a)$$

$$j_{23}[t] = g_3[t] h_{11}[t] - g_1[t] h_{31}[t], \quad (20b)$$

$$j_{33}[t] = g_1[t] h_{21}[t] - g_2[t] h_{11}[t]. \quad (20c)$$

We omit the solutions of the remaining precoders for the sake of brevity. However, these scalars can be computed by using (19) in (15) and (16)–(18). Note that our approach also works

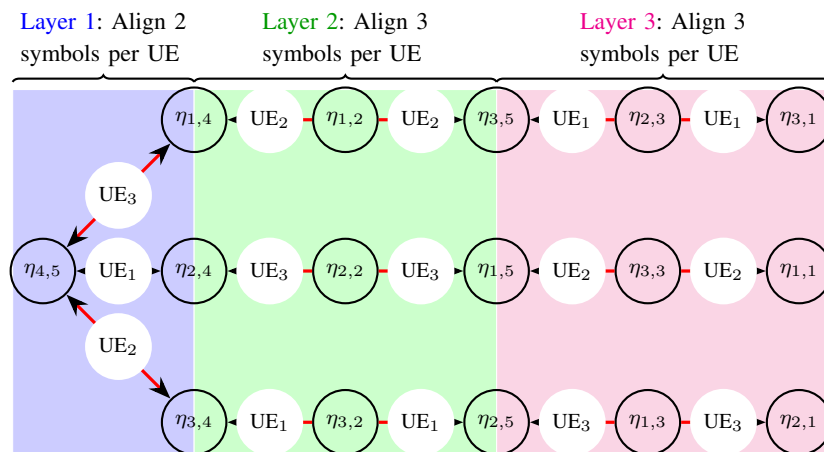


Fig. 6: Interference Alignment Graph for the achievability at corner point $(\frac{4}{5}, \frac{8}{5})$ for $M = 1$ and $K = 3$. The graph consists of three (subspace) alignment chains.

with constant channels under the umbrella of *real interference alignment* [15], [16]. Whether a two-phase precoding design with constant channels (similar to previous work on relay-aided X-channels [17]) attains close-to-optimal performance is an interesting extension to work on. However, it is beyond the scope of this paper. Now we will go through the decoding from the perspective of both the RN and the UEs.

At the receiver side, each UE observes interference aligned to 3 independent signal dimensions and 5 desired symbols occupying 5 independent dimensions. In total, we require $T = 8$ channel uses to allow for reliable decoding. Thus, the achievable DoF of each UE becomes $\frac{5}{8}$ and that of the RN (it only requires $\eta_{4,5}$) $\frac{1}{8}$. The NDT, on the other hand, corresponds to $\frac{8}{5}$. This establishes the achievability for $M = 1$ and $K = 3$.

B. Achievability for $M = 2$

As we increase the number of RNs from $M = 1$ to $M = 2$, the concept of *cooperative interference neutralization* becomes relevant. This enables the exploitation of side information at the RNs to allow them to receive their desired symbols while zero-forcing their contribution at the UEs. With this approach, we are able to show the achievability (and as such the optimality) for all corner points when $M = 2$ and $K = \{1, 2\}$. Details on the achievability is left out due to page limitations. In short, the optimal NDTs, are, respectively, given by

$$\delta^*(\mu) = \max \left\{ 3 - 4\mu, 1 \right\} \text{ for } M = 2, K = 1, \quad (21a)$$

$$\delta^*(\mu) = \max \left\{ 4 - 6\mu, \frac{4 - 3\mu}{2}, \frac{3 - \mu}{2} \right\} \text{ for } M = 2, K = 2. \quad (21b)$$

V. CONCLUSION

In this paper, we studied the fundamental limits on the delivery time for cache-aided wireless networks where relay nodes act as cache-equipped transceivers. We utilized the normalized delivery time (NDT) as a delivery time metric which captures the worst-case latency of the requested file retrieval. To this end, we developed, on the one hand, a novel lower bound for a cache-aided network with M relay nodes and K users. On the other hand, we determined NDT-optimal schemes with which we were able to completely characterize the trade-off between delivery time and cache memory for specific instances where $M = \{1, 2\}$ and $K = \{1, 2, 3\}$ and $M + K \leq 4$. Our achievability schemes determine optimal precoders such that

zero-forcing, interference alignment and cooperative interference neutralization are synergistically combined. The presented schemes are applicable to both time-variant and time-invariant channels. In future work, we would like to generalize our scheme. Specifically, we first aim to determine whether our scheme can be generalized for $M = 1$ and $K > 3$.

REFERENCES

- [1] Network, Evolved Universal Terrestrial Radio Access, "M2 Application Protocol (M2ap) (Release 10)," 2011.
- [2] J. Kakar, "UAV Communications: Spectral Requirements, MAV and SUAV Channel Modeling, OFDM Waveform Parameters, Performance and Spectrum Management," Master's thesis, Virginia Tech, Blacksburg, USA, 2015.
- [3] J. Kakar and V. Marojevic, "Waveform and Spectrum Management for Unmanned Aerial Systems Beyond 2025," 2017. [Online]. Available: <http://arxiv.org/abs/1708.01664>
- [4] M. A. Maddah-Ali and U. Niesen, "Fundamental limits of caching," *Trans. on Info. Theory*, vol. 60, no. 5, pp. 2856–2867, May 2014.
- [5] S. Gharekhloo and A. Sezgin, "Latency-Limited Broadcast Channel with Cache-Equipped Helpers," *IEEE Trans. on Wireless Communications*, vol. 16, no. 7, pp. 4192–4203, July 2017.
- [6] M. A. Maddah-Ali and U. Niesen, "Cache-aided interference channels," in *ISIT*, June 2015, pp. 809–813.
- [7] A. Sengupta, R. Tandon, and O. Simeone, "Cache aided wireless networks: Tradeoffs between storage and latency," in *CISS*, March 2016, pp. 320–325.
- [8] R. Tandon and O. Simeone, "Cloud-aided wireless networks with edge caching: Fundamental latency trade-offs in fog radio access networks," in *ISIT*, July 2016, pp. 2029–2033.
- [9] J. Kakar, S. Gharekhloo, and A. Sezgin, "Fundamental Limits on Delivery Time in Cloud- and Cache-Aided Heterogeneous Networks," 2017. [Online]. Available: <http://arxiv.org/abs/1706.07627>
- [10] J. Kakar, S. Gharekhloo, Z. H. Awan, and A. Sezgin, "Fundamental limits on latency in cloud- and cache-aided HetNets," in *ICC*, May 2017, pp. 1–6.
- [11] N. Naderializadeh, M. A. Maddah-Ali, and A. S. Avestimehr, "Fundamental limits of cache-aided interference management," in *ISIT*, July 2016, pp. 2044–2048.
- [12] J. Kakar, S. Gharekhloo, and A. Sezgin, "Fundamental limits on latency in transceiver cache-aided HetNets," in *ISIT*, June 2017, pp. 2955–2959.
- [13] Y. Liu and E. Erkip, "Completion time in broadcast channel and interference channel," in *Allerton*, Sept 2011, pp. 1694–1701.
- [14] V. R. Cadambe and S. A. Jafar, "Interference Alignment and the Degrees of Freedom of Wireless X Networks," *IEEE Trans. on Info. Theory*, vol. 55, no. 9, pp. 3893–3908, Sept 2009.
- [15] A. S. Motahari, S. Oveis-Gharan, M. A. Maddah-Ali, and A. K. Khandani, "Real Interference Alignment: Exploiting the Potential of Single Antenna Systems," *IEEE Trans. on Info. Theory*, vol. 60, no. 8, pp. 4799–4810, Aug 2014.
- [16] M. A. Maddah-Ali, "On the degrees of freedom of the compound MISO broadcast channels with finite states," in *ISIT*, June 2010, pp. 2273–2277.
- [17] D. Frank, K. Ochs, and A. Sezgin, "A systematic approach for interference alignment in CSIT-less relay-aided X-networks," in *WCNC*, April 2014, pp. 1126–1131.