# Cloud-Edge Non-Orthogonal Transmission for Fog Networks with Delayed CSI at the Cloud

Jingjing Zhang and Osvaldo Simeone

*Abstract*—In a Fog Radio Access Network (F-RAN), the cloud processor (CP) collects channel state information (CSI) from the edge nodes (ENs) over fronthaul links. As a result, the CSI at the cloud is generally affected by an error due to outdating. In this work, the problem of content delivery based on fronthaul transmission and edge caching is studied from an information-theoretic perspective in the high signal-to-noise ratio (SNR) regime. For the set-up under study, under the assumption of perfect CSI, prior work has shown the (approximate or exact) optimality of a scheme in which the ENs transmit information received from the cloud and cached contents over orthogonal resources. In this work, it is demonstrated that a non-orthogonal transmission scheme is able to substantially improve the latency performance in the presence of imperfect CSI at the cloud.

*Index Terms*—F-RAN, caching, imperfect CSI.

## I. INTRODUCTION

Consider the scenario shown in Fig. 1, in which a wireless cellular system delivers contents to a number of users by means of a centralized CP with full access to the library content, fronthaul links, and ENs endowed with caching capabilities. This set-up, referred to as a F-RAN, is motivated by current trends in the evolution of wireless cellular systems.

The F-RAN model has been recently studied from an information-theoretic perspective [1], with special cases excluding CP and fronthaul links investigated in [2, 3]. These works consider the high-SNR regime in order to focus on the impact of interference. Reference [1] proves that, in the presence of full CSI at both CP and ENs, it is approximately optimal to have the ENs transmit information received from the cloud and cached contents over orthogonal resources using time division multiplexing.

In practice, as seen in Fig. 1, the cloud processor collects CSI from the ENs over fronthaul links. The ENs, instead, can directly receive CSI from the users via feedback under a Frequency Division Multiplexing (FDD) operation. As a result, the CSI at the CP is generally affected by an additional error due to outdating associated to fronthaul transmission latency.

In [4], a novel approach is proposed that is shown to improve the high-SNR performance of multi-antenna systems in the presence of imperfect CSI. The scheme is based on rate splitting and superposition coding. Accordingly, a message of interest for multiple receivers is transmitted on the same radio resources as private messages intended for individual users in a non-orthogonal fashion with a proper power allocation.

In this work, it is demonstrated that for an F-RAN with imperfect CSI at the CP, a non-orthogonal transmission scheme
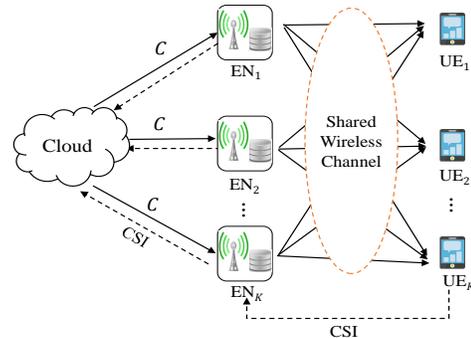


Fig. 1. Cloud and cache-based F-RAN system with solid lines representing downlink communication for content delivery and dashed lines indicating uplink CSI feedback.

is able to substantially improve the latency performance in the presence of imperfect CSI at the cloud in the high-SNR regime. Unlike [4], the signal being multiplexed non-orthogonally are the information received from the cloud on the fronthaul links and locally encoded functions of the cached contents. A similar superposition scheme, known as hybrid fronthauling, was first proposed in [5]. There, assuming full CSI, it was shown via numerical results that the scheme offers performance advantages in the finite-SNR regime, particularly for lower cache capacities.

**Notation:** For any integer $K$, we define the set $[K] \triangleq \{1, 2, \cdots, K\}$. We also define the notation $\{f_n\}_{n=1}^N \triangleq \{f_1, \cdots, f_N\}$. For a set $\mathcal{A}$, $|\mathcal{A}|$ represents the cardinality. We use the symbol $\doteq$ to denote an exponential equality in the sense that we write $f(P) \doteq P^\alpha$ if the limit $\lim_{P\to\infty} \log f(P)/\log P = \alpha$ holds.

## II. SYSTEM MODEL

### A. F-RAN Model

As shown in Fig. 1, we consider a F-RAN model in which $K$ ENs serve $K$ users through a shared wireless channel, with all nodes having a single antenna. We study downlink content delivery from a library of $N$ files $\{W_n\}_{n=1}^N$, each with size $L$ bits. Each EN is equipped with a cache, which can store $\mu NL$ bits from the library during the offline caching phase, with $\mu \in [0, 1]$ being the fractional cache capacity. A CP connects to each EN via a dedicated fronthaul link with capacity $C$ bits per symbol, and it has access to the whole library. A symbol refers to a channel use of the wireless channel.

In each dedicated time-slot, the signal received at each symbol for each user $k$ is given as

$$y_k = \boldsymbol{h}_k^T \boldsymbol{x} + n_k = \sum_{i=1}^K h_k^i x_i + n_k, \tag{1}$$

where we have defined the vectors $\boldsymbol{h}_k = [h_k^1, \cdots, h_k^K]^T$ and $\boldsymbol{x} = [x_1, \cdots, x_K]^T$ with $x_i$ being the transmitted signal of

EN $i$; $h_k^i$ is the complex channel coefficient between EN $i$ and user $k$; and $n_k$ is a zero-mean complex Gaussian noise with normalized unit power. Channels are independent, drawn from a continuous distribution, and constant within each time-slot. We impose the power constraint $\mathbb{E}[|x_i|^2] \leq P$ for each EN $i$.

As illustrated in Fig. 1, we assume an FDD operation in which the users perform channel estimation based on downlink training, and then feed back the estimated CSI to the ENs over an uplink control channel. The CSI is finally communicated from the ENs to the cloud over fronthaul links. To study the problem in its simplest instantiation, we assume here that the feedback channel from users to ENs is ideal, so that each EN has full information about all channel gains $\mathbf{H} = \{h_k^i\}$, with $i, k \in [K]$. In contrast, owing to fronthaul transmission delays and processing, the CSI available at the cloud is assumed to be delayed. The distortion of the outdated CSI $\{\hat{h}_k^i\}$ available at the cloud, i.e., the innovation $h_k^i - \hat{h}_k^i$, is characterized by its power $\mathbb{E}[|h_k^i - \hat{h}_k^i|^2]$. Following the standard model considered in, e.g., [6, 7], we assume that this mean squared error scales as $P^{-\alpha}$ with respect to the SNR $P$, for some $\alpha \geq 0$. More formally, we impose the exponential equality

$$\mathbb{E}[|h_k^i - \hat{h}_k^i|^2] \doteq P^{-\alpha}. \tag{2}$$

In the high-SNR regime, the case $\alpha = 0$, which corresponds to finite-precision CSI, is equivalent to having no CSIT; while, at the other extreme, the case $\alpha = 1$ yields a negligible CSI error. We study the general case in which the cloud has imperfect CSI in the sense of (2) with any arbitrary value $\alpha \in [0, 1]$.

### B. Caching and Delivery Policies

The communication protocol consists of caching and delivery phases.

*1) Caching phase:* In the caching phase, the cache of each EN $i$ is proactively filled with information from the content library. More precisely, each file $W_n$ is mapped to a cached content $V_n^i$ by an arbitrary function $f_n^i$ as $V_n^i = f_n^i(W_n)$. To satisfy the capacity limitation of the cache, we have the inequality $\log_2 |V_n^i| \leq \mu L$ bits, where $|V_n^i|$ represents the alphabet of variable $V_n^i$. The overall cache content of EN $i$ is hence given as $V^i = \{V_n^i\}_{n=1}^N$.

*2) Delivery phase:* In the delivery phase, for any demand vector $\boldsymbol{d} = [d_1, d_2, \cdots, d_K]$, where $W_{d_k}$ is the file requested by user $k$, the delivery of the $K$ files is completed via fronthaul and edge transmission. Fronthaul transmission occurs first with a duration of $T_F$ symbols. On each $i$th fronthaul link, the CP sends a message $U^i$ about the requested files to EN $i \in [K]$. This message is obtained as a function of the available channel estimates $\hat{\mathbf{H}} = \{\hat{h}_k^i\}_{i,k \in [K]}$, of the requested files, and of information about the cached files, as $U^i = g_f^i(\boldsymbol{d}, \hat{\mathbf{H}}, \{V^i\}_{i=1}^K)$. By the fronthaul capacity constant, we have the condition $\log_2 |U^i| \leq T_F C$ bits. Fronthaul transfer is followed by edge transmission, whereby each EN $i$ sends $T_E$ symbols obtained as the function $x_i = g_e^i(\boldsymbol{d}, \mathbf{H}, U^i, V^i)$ on channel (1).

### C. Performance Metric: NDT

A sequence of policies defined by functions $\{\{f_n^i\}, g_f^i, g_e^i\}$ is feasible if each user $k$ is able to decode the desired file $W_{d_k}$ with negligible probability of error when $L \to \infty$. For any feasible policy, we are interested in the delivery latency performance in the high-SNR regime. As in [8], we parameterize the fronthaul capacity as $C = r \log(P)$, where $r$ is the fronthaul rate. Furthermore, we normalize the delivery latency $T_F + T_E$, where $T_F$ and $T_E$ are the corresponding fronthaul and edge latencies, by the term $L/\log(P)$. This represents the downlink latency of an ideal system where each user is served without interference at the high-SNR capacity $\log P$ bit/symbol via the wireless channels [8]. As a result, the fronthaul and edge NDTs are defined as

$$\delta_F = \lim_{P \to \infty} \lim_{L \to \infty} \frac{\mathbb{E}[T_F]}{L/\log(P)} \text{ and } \delta_E = \lim_{P \to \infty} \lim_{L \to \infty} \frac{\mathbb{E}[T_E]}{L/\log(P)}.$$

The overall NDT achieved by a sequence of feasible policies is given as

$$\delta = \delta_E + \delta_E. \tag{3}$$

For given parameters $(\mu, r, \alpha)$, the minimum NDT across all feasible policies is denoted as $\delta^*(\mu, r, \alpha)$.

### III. ORTHOGONAL CLOUD-EDGE DELIVERY

In this section, we study the NDT performance of various policies based on state-of-the-art caching and delivery strategies, as described in [2, 3, 8, 9]. According to these approaches, the ENs transmit information obtained from the caches and from the cloud on the fronthaul links in orthogonal time-slots. We refer to this class of techniques as performing *orthogonal cloud-edge delivery*. As shown in [8], the mentioned orthogonal strategies yield the minimum NDT in the presence of perfect CSI at the cloud, i.e., with $\alpha = 1$, when $K = 2$, and are more generally optimal within a multiplicative factor of two. We start with edge-based and cloud-based policies, which are then orthogonally multiplexed by means of time-sharing. It is noted that cloud-based soft-transfer fronthauling, proposed in [8], is generalized here to account for imperfect CSI at the cloud (Lemma 1).

### A. Edge-based Policies

We first consider caching polices based on only edge resources, for which the NDT performance does not depend on the CSI quality $\alpha$ at the cloud. Note that here we have zero fronthaul NDT, and hence the NDT (3) is given as $\delta = \delta_E$.

*1) Edge-based ZF-beamforming.* When the fractional cache capacity is $\mu = 1$, full cooperation at the ENs is possible given that all ENs store the entire library. As a result, the ideal NDT of 1 can be achieved by means of ZF beamforming [8, Lemma 2], i.e., we have the achievable NDT

$$\delta_{EZ} = 1. \tag{4}$$

*2) Edge-based interference alignment (IA).* Consider now the case with cache capacity $\mu = 1/K$. This is the minimum cache size allowing for delivery based only on cached contents. By caching disjoint fractions of all files and using X channel IA for delivery [8, Lemma 3], the following NDT is achievable

$$\delta_{EI} = 2 - \frac{1}{K}. \tag{5}$$

## B. Cloud-based Policies

We now focus on cloud-based policies assuming that the ENs have zero cache capacity, i.e., $\mu = 0$. As discussed in [8], we can distinguish hard and soft-transfer fronthauling approaches. The former techniques send uncoded fractions of files via the fronthaul links. In contrast, the latter method implements ZF beamforming at the cloud, and sends quantized precoded signals to the ENs. Unlike hard-transfer fronthauling, soft-transfer fronthauling relies on perfect CSI at the cloud in order to perform ZF beamforming. Here, we generalize the analysis of the corresponding achievable NDT to the case of imperfect CSI.

*3) Cloud-based hard-transfer fronthauling.* With zero cache capacity, i.e., with $\mu = 0$, the NDT

$$\delta_{CH} = \min\left\{1 + \frac{K}{r}, 2 - \frac{1}{K} + \frac{1}{r}\right\} \quad (6)$$

is achievable by using one of the following hard-transfer fronthauling approaches, which attain the two NDTs in (6). In the first scheme, the cloud sends all the requested files to each EN via the respective fronthaul link, and the ENs carry out cooperative ZF beamforming. In the second scheme, the requested files are split into disjoint fragments and sent to the ENs, which performs X channel IA (see [8, Proposition 2]).

*4) Achievable NDT with cloud-based soft-transfer fronthauling.* With soft-transfer fronthauling, we have the achievable NDT described in the following lemma.

*Lemma 1:* In an $K \times K$ F-RAN with CSI quality $\alpha \in [0, 1]$ at the cloud and cache capacity $\mu = 0$, the NDT

$$\delta_{CS} = \frac{1}{r} + \frac{1}{\alpha} \quad (7)$$

is achievable by means of soft-transfer fronthauling.

*Proof:* For any request vector $\boldsymbol{d}$, the cloud precodes the $K$ requested files producing the $K \times 1$ vector

$$\bar{\boldsymbol{x}}_F = \sum_{k=1}^{K} \boldsymbol{v}_{d_k} s_{d_k}, \quad (8)$$

where $s_{d_k}$ is a symbol of the codeword encoding file $W_{d_k}$, and $\boldsymbol{v}_{d_k} \in \mathbb{C}^{K \times 1}$ is the corresponding beamforming vector. The symbols $\{s_{d_k}\}_{k=1}^{K}$ are taken from a Gaussian codebook of equal rate $R_F$ bits/symbol and equal power. The power of the symbols $s_{d_k}$ is set to be exponentially equal to $P$, i.e., $\mathbb{E}[|s_{d_k}|^2] \doteq P$. The unit-norm precoding vector $\boldsymbol{v}_{d_k}$ is selected to be orthogonal to all the estimates of the vectors $\{\boldsymbol{h}_{k'}\}_{k' \in [K], k' \neq k}$, i.e., we have the equalities $\hat{\boldsymbol{h}}_{k'}^T \boldsymbol{v}_{d_k} = 0$ for all $k' \neq k$. The cloud then quantizes each element $\bar{x}_{Fi}$ of the vector signal $\bar{\boldsymbol{x}}_F = [\bar{x}_{F1}, \cdots, \bar{x}_{FK}]^T$ as

$$x_{Fi} = \bar{x}_{Fi} + q_i, \quad (9)$$

for $i \in [K]$, where $q_i$ is the quantization noise, which is modeled as a Gaussian variable $q_i \sim \mathcal{CN}(0, \sigma^2)$ with variance $\sigma^2$, which is independent across index $i$. By standard rate distortion arguments, the variance $\sigma^2$ is related to the number $B$ of bits for each of the $L/R_F$ samples $x_{Fi}$ by the equalities $B = I(\bar{x}_{Fi}; x_{Fi}) = \log(1 + \mathbb{E}[|\bar{x}_{Fi}|^2]/\sigma^2)$, and hence we have the equality $\sigma^2 = \mathbb{E}[|\bar{x}_{Fi}|^2]/(2^B - 1)$. We choose $B = \alpha \log P$, so that we have the exponential equality $\sigma^2 \doteq$

$P/P^\alpha = P^{1-\alpha}$. Note that the power constraint $\mathbb{E}[|x_{Fi}|^2] \doteq P$ is satisfied. The cloud sends the quantized signal $x_{Fi}$ to EN $i$ for $i \in [K]$ at the fronthaul rate $C = r \log P$ on the fronthaul links. As a result, the fronthaul latency is given as $T_F = B(L/R_F)/C = \alpha L/(R_F r)$.

Each EN $i$ then forwards the quantized signal $x_{Fi}$, i.e., we set $x_i = x_{Fi}$ in (1), and hence each user $k$ receives the signal

$$y_k = \boldsymbol{h}_k^T \boldsymbol{x}_F + n_k = \boldsymbol{h}_k^T \left( \sum_{k=1}^{K} \boldsymbol{v}_{d_k} s_{d_k} + \boldsymbol{q} \right) + n_k \quad (10a)$$

$$= \boldsymbol{h}_k^T \boldsymbol{v}_{d_k} s_{d_k} + \boldsymbol{h}_k^T \left( \sum_{k'=1, k' \neq k}^{K} \boldsymbol{v}_{d_{k'}} s_{d_{k'}} + \boldsymbol{q} \right) + n_k \quad (10b)$$

$$\overset{(a)}{=} \boldsymbol{h}_k^T \boldsymbol{v}_{d_k} s_{d_k} + \underbrace{\tilde{\boldsymbol{h}}_k^T \sum_{k'=1, k' \neq k}^{K} (\boldsymbol{v}_{d_{k'}} s_{d_{k'}})}_{\triangleq z_k} + \boldsymbol{h}_k^T \boldsymbol{q} + n_k, \quad (10c)$$

where we have defined the vectors $\boldsymbol{q} = [q_1, \cdots, q_K]^T$ and $\tilde{\boldsymbol{h}}_k = \boldsymbol{h}_k - \hat{\boldsymbol{h}}_k$; and equality (a) holds due to the conditions $\boldsymbol{h}_k^T \boldsymbol{v}_{d_{k'}} = (\hat{\boldsymbol{h}}_k^T + \tilde{\boldsymbol{h}}_k^T) \boldsymbol{v}_{d_{k'}} = \tilde{\boldsymbol{h}}_k^T \boldsymbol{v}_{d_{k'}}$. Given the CSI error scaling (2), the power of the interference $z_k$ satisfies the exponential equality $\mathbb{E}[|z_k|^2] \doteq P^{1-\alpha}$. Furthermore, the power of the effective noise, namely quantization plus Gaussian noise, is given as $\mathbb{E}[|\boldsymbol{h}_k^T \boldsymbol{q}|^2] + \mathbb{E}[|n_k|^2] \doteq P^{1-\alpha}$. It follows that the file $W_{d_k}$, encoded by $s_{d_k}$, can be decoded reliably in the high-SNR regime with rate $R_F = \log(P/P^{1-\alpha}) = \alpha \log P$ by treating interference as noise. Finally, we have the fronthaul NDT $\delta_F = 1/r$ and the duration of edge transmission is given by $T_E = L/(\alpha \log P)$, yielding the edge NDT $\delta_E = 1/\alpha$. This completes the proof. ∎

## C. Cloud and Edge-based Policies

By means of time-sharing of the cloud- and edge-based solutions described above, cloud-edge orthogonal delivery achieves the following NDT.

*Proposition 1: Achievable NDT with cloud-edge orthogonal transmission.* In an $K \times K$ F-RAN with CSI quality $\alpha \in [0, 1]$ at the cloud and cache capacity $\mu > 0$, the following NDT is achievable by cloud-edge orthogonal delivery

$$\delta_O(\mu, r, \alpha) = \min\{(\delta_{EZ} - \delta_C)\mu + \delta_C, \delta_O'(\mu, r, \alpha)\} \quad (11)$$

where we have defined $\delta_C = \min\{\delta_{CH}, \delta_{CS}\}$ and

$$\delta_O'(\mu, r, \alpha) = \begin{cases} 2 - \mu, & \text{if } \mu \geq \frac{1}{K} \\ (\delta_{EI} - \delta_C)K\mu + \delta_C, & \text{if } \mu \leq \frac{1}{K}. \end{cases} \quad (12)$$

*Proof:* The result follows by the standard time sharing [8, Lemma 1]. ∎

## IV. NON-ORTHOGONAL CLOUD-EDGE DELIVERY

In this section, we introduce and analyze the proposed non-orthogonal cloud-edge delivery scheme. Unlike the orthogonal strategies studied in the previous section, cloud-based delivery, which is based on soft-transfer fronthauling, and edge-based delivery, which leverages the cached contents, are performed simultaneously at the ENs. As illustrated in Fig. 2, the technique is specifically based on the superposition at the ENs of
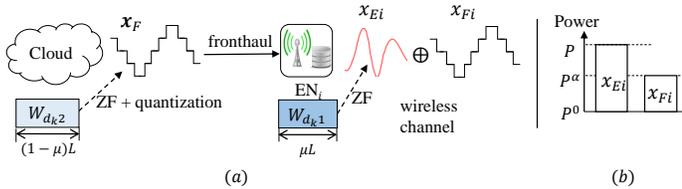
Fig. 2. Illustration of non-orthogonal cloud-edge delivery: (a) caching and delivery; and (b) powers of the signals $x_{Ei}$ and $x_{Fi}$ at each EN $i$.

the quantized fronthaul signals precoded at the cloud and of the signals encoding cached information, as well as on successive interference cancellation (SIC) at the users.

*Proposition 2: Achievable NDT with cloud-edge non-orthogonal transmission.* In an $K \times K$ F-RAN with CSI quality $\alpha \in [0, 1]$ at the cloud and cache capacity $\mu > 0$, the NDT

$$\delta_{NO}(\mu, r, \alpha) = \begin{cases} 1 + \frac{1-\mu}{r}, & \mu \geq 1 - \alpha \\ (1-\mu)(\frac{1}{\alpha} + \frac{1}{r}), & \mu \leq 1 - \alpha, \end{cases} \quad (13)$$

is achievable by means of non-orthogonal cloud-edge delivery.

A sketch of the proof is as follows and a full proof can be found in the Appendix. As illustrated in Fig. 2, in the caching phase, a fraction $\mu$ of each file is stored at all ENs. In the delivery phase, the cloud precodes the uncached $(1 - \mu)$-fraction of the requested files using the cloud-based soft-transfer fronthauling scheme detailed in the proof of Lemma 1. The quantized signals are then sent to the ENs via the fronthaul links. The ENs perform cooperative ZF precoding for the fraction $\mu$ of the cached contents using edge-based ZF beamforming as described in Section III-A. The resulting fronthaul and locally encoded signals are summed and sent to the users, with the former being transmitted with a lower power than the latter. Each user decodes the two signals by using SIC: the locally precoded signal, which has a higher power, is decoded first by treating the fronthaul precoded signal as noise. This signal is then removed from the received signals, and, finally, the fronthaul precoded signal is decoded by the user.

*Remark 1:* In a $K \times K$ F-RAN with $r \geq 1$ and $\mu \geq 1 - \alpha$, the NDT (13) coincides with the NDT derived in [8] for the case of perfect CSI. This NDT is proved in [9] to be optimal for $K = 2$ and to be generally within a multiplicative gap of 2 from the minimum NDT under perfect CSI. This suggests that imperfect CSI at the cloud may not cause any performance degradation as long as $\mu$ and $r$ are sufficiently large if non-orthogonal transmission is used.

Finally, combining orthogonal and non-orthogonal cloud-edge approaches, an improved achievable NDT can be obtained by means of time-sharing.

*Proposition 3: (Achievable NDT).* In an $K \times K$ F-RAN with CSI quality $\alpha \in [0, 1]$ at the cloud and cache capacity $\mu > 0$, the following NDT is achievable

$$\delta(\mu, r, \alpha) = \text{l.c.e.}\big(\delta_O(\mu, r, \alpha), \delta_{NO}(\mu, r, \alpha)\big), \quad (14)$$

where the lower convex envelope (l.c.e.)[1] is evaluated with respect to $\mu$.

[1] The l.c.e., is the supremum of all convex functions that lie under the given function.
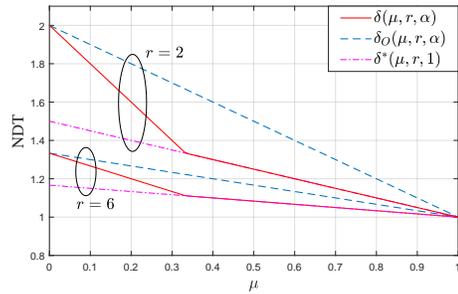


Fig. 3. NDT $\delta(\mu, r, \alpha)$ in (14) with non-orthogonal delivery and $\delta_O(\mu, r, \alpha)$ in (13) with orthogonal delivery, and minimum NDT $\delta^*(\mu, r, 1)$ with $\alpha = 1$, versus $\mu$ for the $2 \times 2$ F-RAN with different values of $r$ and $\alpha = 2/3$.
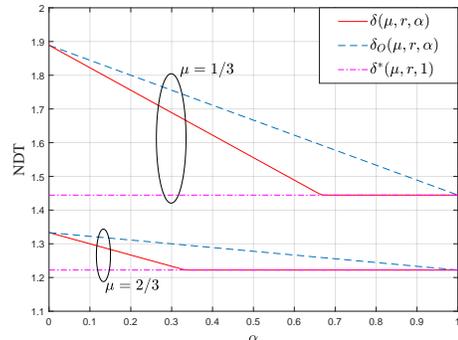


Fig. 4. NDT $\delta(\mu, r, \alpha)$ in (14) with non-orthogonal delivery and $\delta_O(\mu, r, \alpha)$ in (13) with orthogonal delivery versus $\alpha$, and minimum NDT $\delta^*(\mu, r, 1)$ with $\alpha = 1$, for the $2 \times 2$ F-RAN with $r = 1.5$ and different values of $\mu$.

## V. NUMERICAL EXAMPLES

In this section, we compare the NDT $\delta_O(\mu, r, \alpha)$ in (13) with conventional orthogonal delivery, the NDT $\delta(\mu, r, \alpha)$ in (14) with non-orthogonal delivery, and the minimum NDT $\delta^*(\mu, r, 1)$ derived in [9] for the case of perfect CSI ($\alpha = 1$) at the cloud. Throughout, we set $K = 2$.

Fig. 3 plots the mentioned NDTs with $\alpha = 2/3$ for both cases $r = 2$ and $r = 6$. It is observed that latency gains can be reaped via non-orthogonal delivery for the whole range of values $\mu \in (0, 1)$. The largest gains are obtained in the intermediate regime where both cloud and edge-caching contributions are similarly relevant. Furthermore, as observed in Remark 1, when $\mu \geq 1 - \alpha = 1/3$, the achievable NDT $\delta(\mu, r, \alpha)$ coincides with the minimum NDT $\delta^*(\mu, r, 1)$ obtained with perfect CSI. Note that this result can be achieved with orthogonal delivery only when $\mu = 1$.

Fig. 4 plots the NDTs versus the CSI quality $\alpha$. As $\alpha$ increases, non-orthogonal delivery benefits more than orthogonal delivery from the improved CSI, and it obtains the latency savings for any value of $\alpha \in (0, 1)$. Moreover, when $\alpha \geq 1 - \mu$, as observed in Remark 1, non-orthogonal delivery can obtain the optimal performance under perfect CSI.

## VI. CONCLUSIONS

In this paper, we studied the problem of content delivery in F-RAN in the presence of heterogeneous CSI availability between edge and cloud. A non-orthogonal transmission scheme superimposes signals produced at the cloud and at the edge based on cached contents was shown to reduce delivery latency as compared to conventional orthogonal methods. The approach can obtain optimal full-CSI performance for sufficiently large cache and fronthaul capacities.

## VII. APPENDIX: PROOF OF PROPOSITION 2

In the caching phase, each file is split into two disjoint subfiles, i.e., $W_n = \{W_{n1}, W_{n2}\}$, where $W_{n1}$ has size $\mu L$ bits. The first subfiles $\{W_{n1}\}_{n=1}^{N}$ are placed in each EN's cache. In the delivery phase, consider any demand vector $\boldsymbol{d}$. The cloud precodes the uncached subfiles $\{W_{d_k 2}\}_{k=1}^{K}$ producing signal $\bar{\boldsymbol{x}}_F = \sum_{k=1}^{K} \boldsymbol{v}_{d_k} s_{d_k}$. Unlike the signal in (8), here each symbol $s_{d_k}$ only encodes subfile $W_{d_k 2}$ instead of the whole file. Moreover, the power of the symbols $s_{d_k}$ is set to satisfy the exponential equality $\mathbb{E}[|s_{d_k}|^2] \doteq P^\alpha$. Similar to (9), upon fronthaul quantization, the signal $x_{Fi} = \bar{x}_{Fi} + q_i$ for each EN $i$ is produced, where $q_i \sim \mathcal{CN}(0, \sigma^2)$ is the quantization noise. We again set $B = \alpha \log P$ as in Lemma 1. As a result, we have the exponential equality $\sigma^2 \doteq P^\alpha / P^\alpha = 1$. Hence, the fronthaul latency is given as $T_F = B(L(1-\mu)/R_F)/C = BL(1-\mu)/(R_F C)$, yielding the fronthaul NDT $\delta_F = \alpha \log P(1-\alpha)/(R_F r)$.

The cached subfiles $\{W_{d_k 1}\}_{k=1}^{K}$ available at all ENs are precoded cooperatively using edge-based ZF beamforming. To elaborate, the $K \times 1$ precoded signal produced across the ENs is given as

$$\boldsymbol{x}_E = \sum_{k=1}^{K} \boldsymbol{u}_{d_k} c_{d_k}, \tag{15}$$

where symbol $c_{d_k}$ denotes the codeword encoding the cached subfile $W_{d_k 1}$, with rate $R_E$ bits/symbol; and the corresponding precoder $\boldsymbol{u}_{d_k} \in \mathbb{C}^{K \times 1}$ is designed to be orthogonal to all the channels $\{\boldsymbol{h}_{k'}\}_{k' \in [K], k' \neq k}$, i.e., $\boldsymbol{h}_{k'}^T \boldsymbol{u}_{d_k} = 0$. Recall that this is feasible since the ENs have perfect CSI. The edge encoded signals $c_{d_k}$ have power $\mathbb{E}[|c_{d_k}|^2] \doteq P$.

The ENs superimpose the edge precoded signal $\boldsymbol{x}_E$ with the cloud precoded and quantized signal $\boldsymbol{x}_F$, yielding the signal $\boldsymbol{x} = \boldsymbol{x}_E + \boldsymbol{x}_F$. Note that the power constraint $\mathbb{E}[|x_i|^2] \doteq P$ is satisfied. As a result, the received signal at user $k$ is given as

$$y_k = \boldsymbol{h}_k^T \boldsymbol{x} + n_k \tag{16a}$$
$$= \boldsymbol{h}_k^T (\boldsymbol{x}_E + \boldsymbol{x}_F) + n_k \tag{16b}$$
$$= \boldsymbol{h}_k^T \left( \sum_{k=1}^{K} \boldsymbol{u}_{d_k} c_{d_k} + \sum_{k=1}^{K} \boldsymbol{v}_{d_k} s_{d_k} + \boldsymbol{q} \right) + n_k \tag{16c}$$
$$\stackrel{(a)}{=} \underbrace{\boldsymbol{h}_k^T \boldsymbol{u}_{d_k} c_{d_k}}_{\triangleq \tilde{x}_{E,k}} + \underbrace{\boldsymbol{h}_k^T \boldsymbol{v}_{d_k} s_{d_k}}_{\triangleq \tilde{x}_{F,k}} + \underbrace{\tilde{\boldsymbol{h}}_k^T \sum_{k'=1, k' \neq k}^{K} (\boldsymbol{v}_{d_{k'}}, s_{d_{k'}})}_{\triangleq z_k}$$
$$+ \boldsymbol{h}_k^T \boldsymbol{q} + n_k, \tag{16d}$$

where equality (a) holds due to the conditions $\boldsymbol{h}_k^T \boldsymbol{u}_{d_{k'}} = 0$ and $\boldsymbol{h}_k^T \boldsymbol{v}_{d_{k'}} = (\hat{\boldsymbol{h}}_k^T + \tilde{\boldsymbol{h}}_k^T) \boldsymbol{v}_{d_{k'}} = \tilde{\boldsymbol{h}}_k^T \boldsymbol{v}_{d_{k'}}$ for any $k' \neq k$. The interference term $z_k$ in (16d) lies power $\mathbb{E}[|z_k|^2] \doteq 1$ due to the CSI error scaling (2); and for the effective noise terms, we have $\mathbb{E}[|\boldsymbol{h}_k^T \boldsymbol{q}|^2] + \mathbb{E}[|n_k|^2] \doteq 1$.

It follows that the edge-encoded symbols $c_{d_k}$ can be decoded first by user $k$ at rate $R_E = (1-\alpha) \log P$ by treating the interference-plus-noise-term, of power exponentially equal to $P^\alpha$, as noise. Having canceled the signal $\tilde{x}_{E,k}$, the user $k$ can decode the cloud-encoded signal $s_{d_k}$ at rate $R_F = \alpha \log P$. We now analyze the resulting edge NDT. To this end, we define the

time $T_{E1} = \mu L/((1-\alpha)\log P)$ required to decode reliably the edge-precoded signals and the time $T_{E2} = (1-\mu)L/(\alpha \log P)$ needed for the cloud-precoded signals to be reliably decoded.

When $\mu \leq (1-\alpha)$, we have $T_{E1} \leq T_{E2}$. In this case, due to the poor CSI at the cloud, the edge NDT is dominated by the latency $T_{E2}$ and, as a result, the total NDT is $\delta = (1-\mu)(1/\alpha + 1/r)$.

When $\mu \geq (1-\alpha)$ instead, we have $T_{E1} \geq T_{E2}$ and the delivery latency is dominated by the transmission of edge-precoded signals. After time $T_{E2} \leq T_{E1}$, the cloud-precoded signal are reliably decoded and hence we can set $\boldsymbol{x}_F = 0$ from the rest of the transmit duration $T_E - T_{E2}$. The delivery of a fraction $T_{E2}/T_{E1}$ of the edge-precoded signal is completed by time $T_{E2}$. The remaining fraction $(1 - T_{E2}/T_{E1})$ of each cached subfile can be sent via ZF beamforming at the edge at rate $R'_E = \log P$ bits without interference, and the required time is given as $T'_{E1} = \mu L (1 - T_{E2}/T_{E1})/\log P$, yielding the total edge time $T_E = T_{E2} + T'_{E1} = L/\log P$, i.e., $\delta_E = 1$. Hence, the overall NDT for this case is $\delta = (1-\mu)/r + 1$. This completes the proof.

## REFERENCES

[1] J. Zhang and O. Simeone, "Fundamental limits of cloud and cache-aided interference management with multi-antenna base stations," in *Proc. IEEE Int. Symp. Information Theory (ISIT)*, June 2018.

[2] J. Hachem, U. Niesen, and S. Diggavi, "Degrees of freedom of cache-aided wireless interference networks," *IEEE Trans. Inf. Theory*, pp. 1–1, 2018.

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

[4] B. Clerckx, H. Joudeh, C. Hao, M. Dai, and B. Rassouli, "Rate splitting for MIMO wireless networks: a promising PHY-layer strategy for LTE evolution," *IEEE Commun. Magazine*, vol. 54, no. 5, pp. 98–105, May 2016.

[5] S. H. Park, O. Simeone, and S. Shamai, "Joint optimization of cloud and edge processing for fog radio access networks," *IEEE Trans. Wireless Commun.*, vol. 15, no. 11, pp. 7621–7632, Nov 2016.

[6] A. G. Davoodi and S. A. Jafar, "Aligned image sets under channel uncertainty: Settling conjectures on the collapse of degrees of freedom under finite precision CSIT," *IEEE Trans. Inf. Theory*, vol. 62, no. 10, pp. 5603–5618, Oct 2016.

[7] T. Gou and S. Jafar, "Optimal use of current and outdated channel state information: Degrees of freedom of the MISO BC with mixed CSIT," *IEEE Communications Letters*, vol. 16, no. 7, pp. 1084 – 1087, Jul. 2012.

[8] A. Sengupta, R. Tandon, and O. Simeone, "Fog-aided wireless networks for content delivery: Fundamental latency tradeoffs," *IEEE Trans. Inf. Theory*, vol. 63, no. 10, pp. 6650–6678, Oct 2017.

[9] R. Tandon and O. Simeone, "Cloud-aided wireless networks with edge caching: Fundamental latency trade-offs in fog radio access networks," in *Proc. IEEE Int. Symp. Information Theory (ISIT)*, July 2016, pp. 2029–2033.