

# Fronthaul Data Reduction in Massive MIMO Aided C-RAN via Two-timescale Hybrid Compression

An Liu<sup>1</sup>, Senior Member, IEEE, Xihan Chen<sup>1</sup>, Wei Yu<sup>2</sup>, Fellow, IEEE, Vincent Lau<sup>3</sup>, Fellow, IEEE and Min-Jian Zhao<sup>1</sup>

<sup>1</sup>College of Information Science and Electronic Engineering, Zhejiang University

<sup>2</sup>The Edward S. Rogers Sr. Department of Electrical and Computer Engineering, University of Toronto

<sup>3</sup>Department of Electronic and Computer Engineering, Hong Kong University of Science and Technology

**Abstract**—In massive MIMO aided cloud radio access network (C-RAN), plenty of remote radio heads (RRHs), each equipped with a massive MIMO array, are distributed within a specific geographical area and are connected to a centralized baseband unit (BBU) pool through fronthaul links. One major performance bottleneck in the uplink of massive MIMO aided C-RAN is that, the RRHs need to transport a huge amount of data to the BBU for baseband processings. Existing fronthaul compression methods that rely on fully-digital processing are not suitable for the massive MIMO regime due to their high implementation cost. To overcome this challenge, we propose a two-timescale hybrid analog-and-digital spatial compression scheme at RRHs to reduce the fronthaul data, where the analog filter is updated at a slow timescale according to the channel statistics to achieve massive MIMO array gain, and the digital filter is updated at a fast timescale according to the instantaneous effective channel state information (CSI) to achieve spatial multiplexing gain. Such a design can alleviate the performance bottleneck of limited fronthaul with reduced hardware cost and power consumption, and is more robust to the CSI delay. We propose an online algorithm for the two-timescale non-convex optimization of analog and digital filters. Simulations verify the advantages of the proposed scheme over state-of-the-art baseline schemes.

**Index Terms**—Cloud radio access network, Massive MIMO, Hybrid compression and forward

## I. INTRODUCTION

Recently, massive MIMO aided C-RAN has been proposed to improve the spectral efficiency of wireless systems [1]. However, such an architecture requires a huge amount of digital sampled data to be transported over the fronthaul link. Therefore, it is necessary to compress the uplink data at each RRH to satisfy the limited fronthaul capacity constraint. Various fully-digital fronthaul compression techniques have been proposed for C-RAN with small-scale multi-antenna RRHs [2], [3]. In particular, the spatial compression and forward scheme proposed in [3] combines fully-digital spatial filtering and uniform scalar quantization to alleviate the performance bottleneck caused by the limited fronthaul capacity. Unfortunately, fully-digital spatial filtering requires a larger number of analog-to-digital converter (ADCs) and radio frequency (RF) chains at each massive MIMO RRH. In [4], a fully-analog linear spatial filtering is used at each RRH to achieve the fronthaul compression with reduced hardware cost and power consumption. However, fully-analog processing is known to be less efficient than hybrid analog and digital processing.

In this paper, we propose a two-timescale hybrid (analog and digital) compression and forward (THCF) scheme for the uplink transmission of massive MIMO aided C-RAN, to

alleviate the performance bottleneck of the limited fronthaul, with reduced hardware cost and power consumption. In this scheme, each RRH first performs a two-timescale hybrid analog and digital spatial filtering to reduce the dimension of its received signal. Specifically, the analog filtering matrix is adapted to the long-term channel statistics to achieve massive MIMO array gain, and the digital filtering matrix is adapted to the instantaneous effective CSI to achieve spatial multiplexing gain. Then, each RRH applies the uniform scalar quantization over each of these dimensions. Finally, the quantized signals at the RRHs are sent to the BBU for joint decoding.

The power allocation at users, analog/digital filtering matrices and quantization bits allocation at RRHs, and the receive beamforming matrix at the BBU are jointly optimized to maximize a general utility function. We propose an online *block-coordinate stochastic successive convex approximation* (BC-SSCA) algorithm to solve this joint optimization problem. Simulations show that the proposed two-timescale hybrid scheme achieves better tradeoff performance than the baselines.

## II. SYSTEM MODEL

### A. Network Architecture and Channel Model

Consider the uplink of a massive MIMO aided C-RAN, where  $N$  RRHs, each equipped with a massive MIMO array of  $M \gg 1$  antennas and  $S < M$  Rx RF chains, are distributed within a specific geographical area to serve  $K$  single-antenna users. Each RRH  $n$  serves as a relay between the BBU and users, and is connected to the BBU via a fronthaul link of capacity  $C_n$  bits per second (bps). We assume that the number of users  $K$  is fixed and  $NS \gg K$  so that there are enough spatial degrees of freedom to serve all the  $K$  users. In this case, the received signal at RRH  $n$  is given by

$$\mathbf{y}_n = \sum_{k=1}^K \mathbf{h}_{n,k} \sqrt{p_k} s_k + \mathbf{z}_n = \mathbf{H}_n \mathbf{P}^{1/2} \mathbf{s} + \mathbf{z}_n,$$

where  $\mathbf{h}_{n,k} \in \mathbb{C}^M$  is the channel vector from user  $k$  to RRH  $n$ ,  $\mathbf{H}_n = [\mathbf{h}_{n,1}, \dots, \mathbf{h}_{n,K}] \in \mathbb{C}^{M \times K}$ ,  $s_k \sim \mathcal{CN}(0, 1)$  is the data symbol of user  $k$ ,  $\mathbf{s} = [s_1, \dots, s_K]^T$ ,  $p_k$  is the transmit power of user  $k$ ,  $\mathbf{P} = \text{diag}(p_1, \dots, p_K)$ , and  $\mathbf{z}_n \sim \mathcal{CN}(0, \mathbf{I})$  is the noise vector.

### B. Two-timescale Hybrid Compression and Forward at RRHs

Each RRH  $n$  applies the THCF scheme to make sure that the compressed received signal  $\tilde{\mathbf{y}}_n$  satisfies the fronthaul capacity

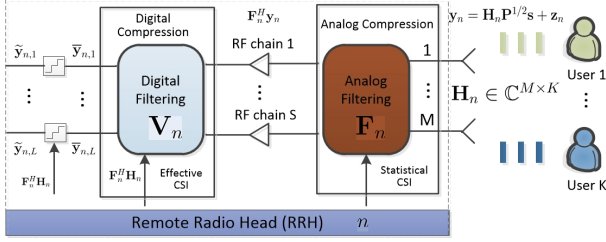


Figure 1: An illustration of the THCF scheme

constraint of  $C_n$  bps, as illustrated in Fig. 1. Specifically, a hybrid filtering matrix  $F_n V_n \in \mathbb{C}^{M \times L}$  is applied at RRH  $n$  to compress the received signal  $\mathbf{y}_n$  into a low-dimensional signal  $\bar{\mathbf{y}}_n = \mathbf{V}_n^H \mathbf{F}_n^H \mathbf{y}_n = [\bar{y}_{n,1}, \dots, \bar{y}_{n,L}]^T \in \mathbb{C}^L$ , where  $\mathbf{F}_n \in \mathbb{C}^{M \times S}$  and  $\mathbf{V}_n = [\mathbf{v}_{n,1}, \dots, \mathbf{v}_{n,L}] \in \mathbb{C}^{S \times L}$  are the analog and digital filtering matrices, respectively, and we set  $L = \min(K, S)$  such that there is no information loss due to digital filtering [3]. The analog filtering matrix  $\mathbf{F}_n$  is usually implemented using an RF phase shifting network [5]. Hence,  $\mathbf{F}_n$  can be represented by a phase vector  $\boldsymbol{\theta}_n \in [0, 2\pi]^{MS}$ , whose  $((j-1)M+i)$ -th element  $\theta_{n,i,j}$  is the phase of the  $(i,j)$ -th element of  $\mathbf{F}_n$ . Then, a simple uniform scalar quantization [3] is applied to each element of  $\bar{\mathbf{y}}_n$  at RRH  $n$ .

After the uniform scalar quantization, the compressed received signal  $\tilde{\mathbf{y}}_n = [\tilde{y}_{n,1}, \dots, \tilde{y}_{n,L}]^T$  is modeled by

$$\tilde{\mathbf{y}}_n = \bar{\mathbf{y}}_n + \mathbf{e}_n = \mathbf{V}_n^H \mathbf{F}_n^H (\mathbf{H}_n \mathbf{P}^{1/2} \mathbf{s} + \mathbf{z}_n) + \mathbf{e}_n,$$

where  $\mathbf{e}_n = [e_{n,1}, \dots, e_{n,L}] \in \mathbb{C}^L$  with  $e_{n,l}$  denoting the quantization error for  $\bar{y}_{n,l}$ . Let  $d_{n,l}$  denote the number of bits that RRH  $n$  uses to quantize the real or imaginary part of  $\bar{y}_{n,l}$ . With uniform scalar quantization, the covariance matrix of  $\mathbf{e}_n$  is given by a function of  $\mathbf{p} = [p_1, \dots, p_K]^T$ ,  $\mathbf{F}_n \mathbf{V}_n$  and  $\mathbf{d}_n = [d_{n,1}, \dots, d_{n,L}]^T$  as [3]

$$\mathbf{Q}_n(\mathbf{p}, \mathbf{F}_n \mathbf{V}_n, \mathbf{d}_n) = \text{diag}(q_{n,1}, \dots, q_{n,L}),$$

$$q_{n,l} = \begin{cases} \frac{3}{4^{d_{n,l}}} (\sum_{k=1}^K p_k |\mathbf{h}_{n,k}^H \tilde{\mathbf{v}}_{n,l}|^2 + \|\tilde{\mathbf{v}}_{n,l}\|^2) & \text{if } d_{n,l} > 0, \\ \infty & \text{if } d_{n,l} = 0, \end{cases} \quad (1)$$

where  $\tilde{\mathbf{v}}_{n,l} = \mathbf{F}_n \mathbf{v}_{n,l}$ . Finally, each RRH forwards the quantized bits to the BBU via the fronthaul link.

### C. Joint Rx Beamforming at the BBU

The received signal  $\tilde{\mathbf{y}} = [\tilde{\mathbf{y}}_1^T, \dots, \tilde{\mathbf{y}}_N^T]^T$  at the BBU from all RRHs can be expressed as

$$\tilde{\mathbf{y}} = \tilde{\mathbf{V}}^H \mathbf{H} \mathbf{P}^{1/2} \mathbf{s} + \tilde{\mathbf{V}}^H \mathbf{z} + \mathbf{e},$$

where  $\tilde{\mathbf{V}} = \text{diag}(\mathbf{F}_1 \mathbf{V}_1, \dots, \mathbf{F}_N \mathbf{V}_N) \in \mathbb{C}^{MN \times LN}$ ,  $\mathbf{H} = [\mathbf{h}_1, \dots, \mathbf{h}_K] \in \mathbb{C}^{MN \times K}$  with  $\mathbf{h}_k = [\mathbf{h}_{1,k}^T, \dots, \mathbf{h}_{N,k}^T]^T$  denoting the composite channel vector of user  $k$ ,  $\mathbf{z} = [\mathbf{z}_1^T, \dots, \mathbf{z}_N^T]^T$ , and  $\mathbf{e} = [\mathbf{e}_1^T, \dots, \mathbf{e}_N^T]^T$ . A joint Rx beamforming vector  $\mathbf{u}_k \in \mathbb{C}^{NL \times 1}$  is applied at the BBU to obtain the estimated data symbol for each user  $k$  as

$$\hat{s}_k = \mathbf{u}_k^H \tilde{\mathbf{V}}^H \mathbf{H} \mathbf{P}^{1/2} \mathbf{s} + \mathbf{u}_k^H \tilde{\mathbf{V}}^H \mathbf{z} + \mathbf{u}_k^H \mathbf{e}, \forall k.$$

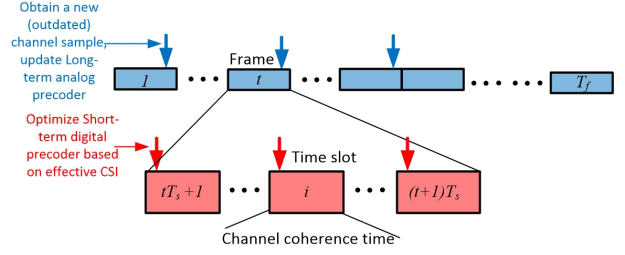


Figure 2: An illustration of two-timescale frame structure.

### D. Frame Structure and Achievable Data Rate

We focus on a coherence time interval of channel statistics, which is divided into  $T_f$  frames with each frame consisting of  $T_s$  time slots (channel coherence intervals), as illustrated in Fig. 2. We assume that the BBU can obtain the real-time effective CSI  $\mathbf{F}_n^H \mathbf{H}_n \in \mathbb{C}^{S \times K}$ ,  $\forall n$  at each time slot, and one (possibly outdated) channel sample  $\mathbf{H}$  at each frame. The long-term analog filtering matrices  $\mathbf{F}_n, \forall n$  are only updated once per frame based on a channel sample to achieve massive MIMO array gain. The short-term control variables  $\{\mathbf{p}, \mathbf{V}_n, \mathbf{d}_n, \mathbf{u}_k\}$  are adaptive to the real-time effective CSI  $\mathbf{F}_n^H \mathbf{H}_n, \forall n$  to achieve the spatial multiplexing gain. For convenience, we let  $\mathbf{v} = \text{Vec}([\mathbf{V}_1, \dots, \mathbf{V}_N])$ ,  $\mathbf{d} = [\mathbf{d}_1^T, \dots, \mathbf{d}_N^T]^T$  and  $\mathbf{u} = [\mathbf{u}_1^T, \dots, \mathbf{u}_K^T]^T$ .

For given long-term control variables  $\boldsymbol{\theta} = [\boldsymbol{\theta}_1^T, \dots, \boldsymbol{\theta}_N^T]^T$  (phase vectors of analog filtering matrices), short-term control variables  $\mathbf{x} \triangleq [\mathbf{p}^T, \mathbf{v}^T, \mathbf{d}^T, \mathbf{u}^T]^T$  and channel realization  $\mathbf{H}$ , the achievable data rate of user  $k$  is given by

$$r_k^o(\boldsymbol{\theta}, \mathbf{x}, \mathbf{H}) = \log(1 + \text{SINR}_k(\boldsymbol{\theta}, \mathbf{x}; \mathbf{H})),$$

where  $\text{SINR}_k(\boldsymbol{\theta}, \mathbf{x}; \mathbf{H})$  is the SINR of user  $k$  given by

$$\text{SINR}_k(\boldsymbol{\theta}, \mathbf{x}; \mathbf{H}) = \frac{p_k |\mathbf{u}_k^H \tilde{\mathbf{V}}^H \mathbf{h}_k|^2}{\sum_{l \neq k} p_l |\mathbf{u}_k^H \tilde{\mathbf{V}}^H \mathbf{h}_l|^2 + \|\mathbf{u}_k^H \tilde{\mathbf{V}}^H\|^2 + \mathbf{u}_k^H \mathbf{Q}(\boldsymbol{\theta}, \mathbf{p}, \mathbf{v}, \mathbf{d}) \mathbf{u}_k},$$

$$\mathbf{Q}(\boldsymbol{\theta}, \mathbf{p}, \mathbf{v}, \mathbf{d}) =$$

$$\text{diag}(\mathbf{Q}_1(\mathbf{p}, \mathbf{F}_1 \mathbf{V}_1, \mathbf{d}_1), \dots, \mathbf{Q}_N(\mathbf{p}, \mathbf{F}_N \mathbf{V}_N, \mathbf{d}_N)).$$

Note that  $\mathbf{F}_n$  is a function of  $\boldsymbol{\theta}_n$ .

Let  $\mathbf{x}(\mathbf{H})$  denote the short-term control variable under channel state  $\mathbf{H}$  and  $\bar{\Omega} \triangleq \{\mathbf{x}(\mathbf{H}) \in \bar{\mathcal{X}}, \forall \mathbf{H}\}$  denote the collection of the short-term control variables for all possible channel states, with  $\bar{\mathcal{X}}$  denoting the feasible set of the short-term control variables. To be more specific,  $\bar{\mathcal{X}}$  is the set of all short-term control variables  $\mathbf{x} = [\mathbf{p}^T, \mathbf{v}^T, \mathbf{d}^T, \mathbf{u}^T]^T$  that satisfy the following constraints:

$$p_k \in [0, P_k], \forall k, \quad (2)$$

$$2B_W \sum_{l=1}^L d_{n,l} \leq C_n, \forall n, \quad (3)$$

$$d_{n,l} \geq 0 \text{ is an integer}, \forall n, l, \quad (4)$$

where  $P_k$  is the individual power constraint at user  $k$ ,  $B_W$  is the system bandwidth, and (3) is the fronthaul capacity constraint. Then the average data rate of user  $k$  is

$$\bar{r}_k^\circ(\boldsymbol{\theta}, \bar{\Omega}) = \mathbb{E}[r_k^\circ(\boldsymbol{\theta}, \mathbf{x}(\mathbf{H}); \mathbf{H})],$$

For convenience, define  $\bar{\mathbf{r}}^\circ(\boldsymbol{\theta}, \bar{\Omega}) \triangleq [\bar{r}_1^\circ(\boldsymbol{\theta}, \bar{\Omega}), \dots, \bar{r}_K^\circ(\boldsymbol{\theta}, \bar{\Omega})]^T$  as the average data rate vector.

### III. TWO-TIMESCALE JOINT OPTIMIZATION AT BBU

Note that  $r_k^\circ(\boldsymbol{\theta}, \mathbf{x}, \mathbf{H})$  is not a smooth function of  $d_{n,l}$ ,  $\forall n, l$  because  $d_{n,l}$  is integer. To make the problem tractable, we relax the integer constraint on  $d_{n,l}$  and approximate the quantization noise power  $q_{n,l}$ ,  $\forall n, l$  with the following smooth function of a real variable  $d_{n,l} \geq 0$  as [3]

$$\hat{q}_{n,l} = \frac{3}{4d_{n,l}} \left( \sum_{k=1}^K p_k |h_{n,k}^H \tilde{\mathbf{v}}_{n,l}|^2 + \|\tilde{\mathbf{v}}_{n,l}\|^2 \right). \quad (5)$$

We use  $r_k(\boldsymbol{\theta}, \mathbf{x}; \mathbf{H})$  to denote the approximate data rate of user  $k$  obtained by replacing  $q_{n,l}$  in (1) with  $\hat{q}_{n,l}$  in (5) and the integer constraint in (4) with constraint  $d_{n,l} \geq 0$ . Moreover, define  $\bar{\mathbf{r}}(\boldsymbol{\theta}, \Omega) = [\bar{r}_1(\boldsymbol{\theta}, \Omega), \dots, \bar{r}_K(\boldsymbol{\theta}, \Omega)]^T$  as the approximate average data rate vector, where  $\bar{r}_k(\boldsymbol{\theta}, \Omega) = \mathbb{E}[r_k(\boldsymbol{\theta}, \mathbf{x}(\mathbf{H}); \mathbf{H})]$ . Then, the two-timescale joint optimization of long-term and short-term control variables can be formulated as the following utility maximization problem

$$\mathcal{P} : \max_{\boldsymbol{\theta} \in \Theta, \Omega} g(\bar{\mathbf{r}}(\boldsymbol{\theta}, \Omega)), \quad (6)$$

where  $\Omega \triangleq \{\mathbf{x}(\mathbf{H}) \in \mathcal{X}, \forall \mathbf{H}\}$  with  $\mathcal{X}$  denoting the set of all short-term control variables that satisfy constraint (2), (3) and  $d_{n,l} \geq 0$ , the utility function  $g(\bar{\mathbf{r}})$  is continuously differentiable function of  $\bar{\mathbf{r}}$ ,  $\Theta \triangleq [0, 2\pi]^{NMS}$  is the feasible set of  $\boldsymbol{\theta}$ . Moreover,  $g(\bar{\mathbf{r}})$  is non-decreasing w.r.t.  $\bar{r}_k$ ,  $\forall k$  and its derivative  $\nabla_{\bar{\mathbf{r}}} g(\bar{\mathbf{r}})$  w.r.t.  $\bar{\mathbf{r}}$  is Lipschitz continuous. This general utility function  $g(\bar{\mathbf{r}})$  includes many important network utilities as special cases, such as average sum rate ( $g(\bar{\mathbf{r}}) = \sum_{k=1}^K \bar{r}_k$ ) and proportional fairness (PFS) utility ( $g(\bar{\mathbf{r}}) = \sum_{k=1}^K \log(\bar{r}_k + \varepsilon)$ , where  $\varepsilon > 0$  is a small number to avoid the singularity at  $\bar{r}_k = 0$ ).

Since Problem  $\mathcal{P}$  is a two-timescale stochastic non-convex problem, we focus on designing an efficient algorithm to find stationary solutions of Problem  $\mathcal{P}$ , as defined below.

**Definition 1** (Stationary solution of  $\mathcal{P}$ ). A solution  $(\boldsymbol{\theta}^*, \Omega^* = \{\mathbf{x}^*(\mathbf{H}) \in \mathcal{X}, \forall \mathbf{H}\})$  is called a stationary solution of Problem  $\mathcal{P}$  if it satisfies the following conditions:

- 1) For every  $\mathbf{H}$  outside a set of probability zero,

$$(\mathbf{x} - \mathbf{x}^*(\mathbf{H}))^T \mathbf{J}_x(\boldsymbol{\theta}^*, \mathbf{x}^*(\mathbf{H}); \mathbf{H}) \nabla_{\bar{\mathbf{r}}} g(\bar{\mathbf{r}}^*) \leq 0, \quad (7)$$

$\forall \mathbf{x} \in \mathcal{X}$ , where  $\mathbf{J}_x(\boldsymbol{\theta}^*, \mathbf{x}^*(\mathbf{H}); \mathbf{H})$  is the Jacobian matrix of the (approximate) rate vector  $\mathbf{r}(\boldsymbol{\theta}, \mathbf{x}; \mathbf{H}) \triangleq [r_1(\boldsymbol{\theta}, \mathbf{x}; \mathbf{H}), \dots, r_K(\boldsymbol{\theta}, \mathbf{x}; \mathbf{H})]^T$  w.r.t.  $\mathbf{x}$  at  $\boldsymbol{\theta} = \boldsymbol{\theta}^*$  and  $\mathbf{x} = \mathbf{x}^*(\mathbf{H})$ , and  $\nabla_{\bar{\mathbf{r}}} g(\bar{\mathbf{r}}^*)$  is the derivative of  $g(\bar{\mathbf{r}})$  at  $\bar{\mathbf{r}} = \bar{\mathbf{r}}^* \triangleq \bar{\mathbf{r}}(\boldsymbol{\theta}^*, \Omega^*)$ .

- 2)

$$(\boldsymbol{\theta} - \boldsymbol{\theta}^*)^T \nabla_{\boldsymbol{\theta}} g(\bar{\mathbf{r}}(\boldsymbol{\theta}^*, \Omega^*)) \leq 0, \forall \boldsymbol{\theta} \in \Theta, \quad (8)$$

where  $\nabla_{\boldsymbol{\theta}} g(\bar{\mathbf{r}}(\boldsymbol{\theta}^*, \Omega^*)) \triangleq \mathbb{E}[\mathbf{J}_{\boldsymbol{\theta}}(\boldsymbol{\theta}^*, \mathbf{x}^*(\mathbf{H}); \mathbf{H})] \nabla_{\bar{\mathbf{r}}} g(\bar{\mathbf{r}}^*)$  is the partial derivative of  $g(\bar{\mathbf{r}}(\boldsymbol{\theta}^*, \Omega^*))$  w.r.t.  $\boldsymbol{\theta}$  at  $\boldsymbol{\theta} = \boldsymbol{\theta}^*$  and  $\Omega = \Omega^*$ ,  $\mathbf{J}_{\boldsymbol{\theta}}(\boldsymbol{\theta}^*, \mathbf{x}^*(\mathbf{H}); \mathbf{H})$  is the Jacobian matrix of the (approximate) rate vector  $\mathbf{r}(\boldsymbol{\theta}, \mathbf{x}; \mathbf{H})$  w.r.t.  $\boldsymbol{\theta}$ .

Note that a stationary solution  $(\boldsymbol{\theta}^*, \Omega^*)$  of  $\mathcal{P}$  may not satisfy all the integer constraints in (4). To obtain an integer solution for the quantization bits allocation, we use the same method as in [3] to round each  $d_{n,l}^*$  to its nearby integer.

### IV. ONLINE OPTIMIZATION ALGORITHM

#### A. Summary of the BC-SSCA Algorithm

The proposed online BC-SSCA algorithm is summarized in Algorithm 1 and its time line is illustrated in Fig. 2. In BC-SSCA, an auxiliary weight vector  $\boldsymbol{\mu} = [\mu_1, \dots, \mu_K]^T$  is introduced to approximate the derivative  $\nabla_{\bar{\mathbf{r}}} g(\bar{\mathbf{r}}(\boldsymbol{\theta}, \Omega))$ . At the beginning of each coherence time of channel statistics, the BBU resets the BC-SSCA algorithm with an initial analog filter phase vector  $\boldsymbol{\theta}^0$  and a weight vector  $\boldsymbol{\mu}^0$ . Then  $\boldsymbol{\theta}$  and  $\boldsymbol{\mu}$  are updated once at the end of each frame, where  $\boldsymbol{\theta}$  is updated by maximizing a concave surrogate function  $\hat{f}^t(\boldsymbol{\theta})$  of  $g(\bar{\mathbf{r}}(\boldsymbol{\theta}, \bar{\Omega}))$  w.r.t.  $\boldsymbol{\theta}$ . Specifically, let  $\boldsymbol{\theta}^t$  and  $\boldsymbol{\mu}^t$  denote the analog filter phase vector and weight vector used during the  $t$ -th frame. The  $t$ -th iteration ( $t$ -th frame) of the BC-SSCA algorithm is described as follows.

*Step 1 (Short-term optimization at each time slot):* At time slot  $i \in [tT_s + 1, (t+1)T_s]$  in the  $t$ -th frame, the BBU first acquires the effective channel  $(\mathbf{F}_n^t)^H \mathbf{H}_n(i)$ ,  $\forall n$ , where  $\mathbf{H}_n(i)$  is the channel state of RRH  $n$  at time slot  $i$ , and  $\mathbf{F}_n^t$  is the analog filtering matrix at RRH  $n$  corresponding to  $\boldsymbol{\theta}^t$ . Then it calculates the short-term variables  $\mathbf{x}^{J_t}(\boldsymbol{\mu}^t, \boldsymbol{\theta}^t, \mathbf{H}(i))$  from  $(\mathbf{F}_n^t)^H \mathbf{H}_n(i)$ ,  $\forall n$  by running a *short-term block-coordinate (BC) algorithm* with input  $J_t$ ,  $\boldsymbol{\mu}^t$ ,  $\boldsymbol{\theta}^t$  and  $\mathbf{H}_n(i)$ , where  $J_t$  determines the total number of iterations for the short-term BC algorithm at frame  $t$ . For any finite iteration  $t < \infty$ ,  $J_t$  is finite, and we let  $J_t \rightarrow \infty$  as  $t \rightarrow \infty$ . Note that  $\mathbf{x}^{J_t}(\boldsymbol{\mu}^t, \boldsymbol{\theta}^t, \mathbf{H}(i))$  depends on  $\boldsymbol{\theta}^t, \mathbf{H}(i)$  only through the effective channels  $(\mathbf{F}_n^t)^H \mathbf{H}_n(i)$ 's. Specifically, for given input  $J$ ,  $\boldsymbol{\mu}, \boldsymbol{\theta}$  and  $\mathbf{H}$ , the short-term BC runs  $J$  iterations to find a stationary point (up to certain accuracy)  $\mathbf{x}^J(\boldsymbol{\mu}, \boldsymbol{\theta}, \mathbf{H})$  of the following weighted sum-rate maximization problem (WSRMP):

$$\mathcal{P}_S(\boldsymbol{\mu}, \boldsymbol{\theta}, \mathbf{H}) : \max_{\mathbf{x}=[\mathbf{p}^T, \mathbf{v}^T, \mathbf{d}^T, \mathbf{u}^T]^T} \sum_{k=1}^K \mu_k r_k(\boldsymbol{\theta}, \mathbf{x}; \mathbf{H}).$$

The details will be postponed to Section IV-B.

*Step 2 (Long-term optimization at the end of frame  $t$ ):* In Step 2a, the BBU obtains a full channel sample  $\mathbf{H}^t \triangleq \mathbf{H}(tT_s + 1)$  before the end of  $t$ -th frame. Then, in Step 2b (at the end of the  $t$ -th frame), the BBU updates the surrogate function  $\hat{f}^t(\boldsymbol{\theta})$  based on  $\mathbf{H}^t$ , the current iterate  $\boldsymbol{\theta}^t$ , and the short-term control variables  $\mathbf{x}(i) \triangleq \mathbf{x}^{J_t}(\boldsymbol{\mu}^t, \boldsymbol{\theta}^t, \mathbf{H}(i))$ ,  $\forall i \in \mathcal{T}_t \triangleq [tT_s + 1, (t+1)T_s]$  as

$$\hat{f}^t(\boldsymbol{\theta}) = g(\hat{\mathbf{r}}^t) + (\mathbf{f}^t)^T (\boldsymbol{\theta} - \boldsymbol{\theta}^t) - \tau \|\boldsymbol{\theta} - \boldsymbol{\theta}^t\|^2, \quad (9)$$

**Algorithm 1** BC-SSCA Algorithm

---

**Initialize:**  $\theta^0 \in \Theta$ ;  $\mu^0 = [1, \dots, 1]^T$ ,  $t = 0$ .  
**Step 1** (Short-term optimization at time slot  $i \in \mathcal{T}_t$ ):  
 Apply the short-term BC algorithm with input  $J_i$ ,  $\mu^t$ ,  $\theta^t$  and  $\mathbf{H}_n(i)$ , to obtain the short-term variable  $\mathbf{x}^{J_i}(\mu^t, \theta^t, \mathbf{H}(i))$ .  
**Step 2** (Long-term optimization at the end of frame  $t$ ):  
**2a:** Obtain a full channel sample  $\mathbf{H}^t \triangleq \mathbf{H}(tT_s + 1)$ .  
**2b:** Update the surrogate function  $\bar{f}^t(\theta)$  according to (9). Calculate  $\bar{\mu}^t = \nabla_{\bar{r}} g(\hat{\mathbf{r}}^t)$  and update  $\mu^{t+1}$  according to (11).  
**2c:** Solve (12) to obtain  $\bar{\theta}^t$ . Update  $\theta^{t+1}$  according to (13).  
 Let  $t = t + 1$  and return to Step 1.

---

where  $\tau > 0$  is a constant;  $\hat{\mathbf{r}}^t = [\hat{r}_1^t, \dots, \hat{r}_K^t]^T$  is approximate average data rate updated recursively as

$$\hat{r}_k^t = (1 - \rho_t) \hat{r}_k^{t-1} + \rho_t \sum_{i \in \mathcal{T}_t} \frac{r_k(\theta^t, \mathbf{x}(i); \mathbf{H}(i))}{|\mathcal{T}_t|}, \forall k, \quad (10)$$

with  $\hat{r}_k^{-1} = 0, \forall k$ ;  $\mathbf{f}^t$  is an approximation of the partial derivative  $\nabla_{\theta} g(\bar{\mathbf{r}}(\theta, \Omega))$ , which is updated recursively as

$$\mathbf{f}^t = (1 - \rho_t) \mathbf{f}^{t-1} + \rho_t \mathbf{J}_{\theta}(\theta^t, \mathbf{x}(tT_s + 1); \mathbf{H}^t) \nabla_{\bar{\mathbf{r}}} g(\hat{\mathbf{r}}^t),$$

with  $\mathbf{f}^{-1} = \mathbf{0}$ , where  $\rho_t \in (0, 1]$  is a sequence satisfying  $\frac{1}{\rho_t} \leq O(t^\kappa)$  for some  $\kappa \in (0.5, 1)$ ,  $\mathbf{J}_{\theta}(\theta, \mathbf{x}; \mathbf{H})$  is the Jacobian matrix of the rate vector  $\mathbf{r}(\theta, \mathbf{x}; \mathbf{H})$  w.r.t.  $\theta$  and its expression is derived in [6]. The weight vector  $\mu$  is updated as

$$\mu^{t+1} = (1 - \gamma_t) \mu^t + \gamma_t \bar{\mu}^t. \quad (11)$$

with  $\bar{\mu}^t \triangleq \nabla_{\bar{\mathbf{r}}} g(\hat{\mathbf{r}}^t)$ , where  $\gamma_t \in (0, 1]$  is a sequence satisfying  $\sum_t \gamma_t = \infty$ ,  $\sum_t (\gamma_t)^2 < \infty$  and  $\lim_{t \rightarrow \infty} \gamma_t / \rho_t = 0$ .

In Step 2c, the optimal solution  $\bar{\theta}^t$  of the following quadratic optimization problem is solved:

$$\bar{\theta}^t = \operatorname{argmax}_{\theta \in \Theta} \bar{f}^t(\theta), \quad (12)$$

which has closed-form solution  $\bar{\theta}^t = \mathbb{P}_{\Theta} \left[ \theta^t + \frac{\mathbf{f}^t}{2\tau} \right]$ , where  $\mathbb{P}_{\Theta}[\cdot]$  denotes the projection on to the box feasible region  $\Theta$ . Finally,  $\theta$  is updated according to

$$\theta^{t+1} = (1 - \gamma_t) \theta^t + \gamma_t \bar{\theta}^t. \quad (13)$$

Then the above iteration is carried out until convergence. In the full version in [6], we established the convergence of the BC-SSCA algorithm to stationary solutions.

### B. Short-term Block-Coordinate Algorithm

We first transform the WSRMP  $\mathcal{P}_S(\mu, \theta, \mathbf{H})$  to the following weighted minimum mean square error (WMMSE) problem

$$\begin{aligned} \min_{\beta, \mathbf{v}, \mathbf{d}, \mathbf{u}, \mathbf{w}} \quad & \sum_{k=1}^K \mu_k (w_k \eta_k - \log w_k) \\ \text{s.t.} \quad & \mathbf{d} \geq \mathbf{0}, \quad (2) \text{ and } (3), \end{aligned} \quad (14)$$

where  $\mathbf{w} = [w_1, \dots, w_K]$  with  $w_k > 0 : \forall k$  is a weight vector for MSE,  $\beta = [\beta_1, \dots, \beta_K]^T$  with  $|\beta_k|^2 = p_k$  and

$$\begin{aligned} \eta_k & \triangleq \mathbb{E} \left[ |s_k - \hat{s}_k|^2 | \mathbf{H} \right] \\ & = \left| 1 - \mathbf{u}_k^H \tilde{\mathbf{V}}^H \mathbf{h}_k \beta_k \right|^2 + \sum_{l \neq k} \left| \mathbf{u}_k^H \tilde{\mathbf{V}}^H \mathbf{h}_l \beta_l \right|^2 \\ & \quad + \mathbf{u}_k^H \tilde{\mathbf{V}}^H \tilde{\mathbf{V}} \mathbf{u}_k + \mathbf{u}_k^H \mathbf{Q}(\theta, \mathbf{p}, \mathbf{v}, \mathbf{d}) \mathbf{u}_k, \end{aligned}$$

is the MSE of user  $k$ . Following similar proof to that of Theorem 1 in [7], it can be shown that Problem  $\mathcal{P}_S(\mu, \theta, \mathbf{H})$  is equivalent to (14). Therefore, we shall focus on designing a BC algorithm to find a stationary point of (14). In the proposed BC algorithm, starting from an initial point, the short-term control variables  $\beta, \mathbf{v}, \mathbf{d}, \mathbf{u}, \mathbf{w}$  are optimized in an alternating way by solving a convex subproblem w.r.t. each variable. The update equation for each variable is elaborated below.

When fixing the other short-term variables, the optimal  $\mathbf{u}$  is given by the MMSE receiver

$$\mathbf{u}_k = \left( \sum_{l=1}^K \tilde{\mathbf{V}}^H \mathbf{h}_l |\beta_l|^2 \mathbf{h}_l^H \tilde{\mathbf{V}} + \tilde{\mathbf{V}}^H \tilde{\mathbf{V}} + \mathbf{Q} \right)^{-1} \tilde{\mathbf{V}}^H \mathbf{h}_k \beta_k,$$

$\forall k$ , where  $\mathbf{Q}$  is an abbreviation for  $\mathbf{Q}(\theta, \mathbf{p}, \mathbf{v}, \mathbf{d})$ ; the optimal  $w_k, \forall k$  is given by  $w_k = \left( 1 - \mathbf{u}_k^H \tilde{\mathbf{V}}^H \mathbf{h}_k \beta_k \right)^{-1}$ ; and the optimal  $\beta$  is given by  $\beta_k = \beta_k^*(\lambda_k), \forall k$  with

$$\begin{aligned} \beta_k^*(\lambda_k) & = \mu_k w_k \operatorname{Re} \left[ \mathbf{u}_k^H \tilde{\mathbf{V}}^H \mathbf{h}_k \right] \\ & \quad \times \left( \sum_{l=1}^K 2\mu_l w_l \mathbf{h}_k^H \tilde{\mathbf{V}} \mathbf{u}_l \mathbf{u}_l^H \tilde{\mathbf{V}}^H \mathbf{h}_k + \nu_k + 2\lambda_k \right)^{-1}, \end{aligned}$$

where  $\nu_k = \sum_{n,l} \frac{6}{4^{a_{n,l}}} |u_{k,n,l}|^2 |\mathbf{h}_{n,k}^H \tilde{\mathbf{v}}_{n,l}|^2$ ,  $u_{k,n,l}$  is the  $((n-1)N + s)$ -th element of  $\mathbf{u}_k$ ,  $\lambda_k$  is zero if  $|\beta_k^*(0)|^2 \leq P_k$  and chosen to satisfy  $|\beta_k^*(\lambda_k)|^2 = P_k$  otherwise.

When fixing the other short-term variables, we solve the following modified subproblem w.r.t.  $\mathbf{v}$  by adding a proximal regularization term  $\epsilon \left\| \mathbf{v} - \mathbf{v}' \right\|^2$  with  $\epsilon > 0$  a small number:

$$\min_{\mathbf{v}} \sum_{k=1}^K \mu_k (w_k \eta_k - \log w_k) + \epsilon \left\| \mathbf{v} - \mathbf{v}' \right\|^2, \quad (15)$$

where  $\mathbf{v}'$  is the current digital filter. By solving (15), we obtain the updated digital filter as follows

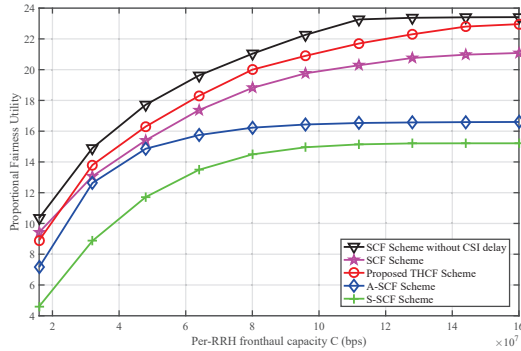
$$\mathbf{v} = (\mathbf{B} + \epsilon \mathbf{I})^{-1} (\mathbf{J} + \epsilon \mathbf{v}'), \quad (16)$$

where  $\mathbf{B} = [\mathbf{B}_{1,1}, \dots, \mathbf{B}_{N,L}]^T$  and  $\mathbf{B}_{n,l} = [\mathbf{B}_{1,1,n,l}, \dots, \mathbf{B}_{N,L,n,l}]^T$  with

$$\mathbf{B}_{n',l',n,l} = \begin{cases} \sum_{k=1}^K \mu_k w_k |u_{k,n,l}|^2 \left( \frac{3}{4^{a_{n,l}}} + 1 \right) \mathbf{D}_n, & n' = n, l' = l, \\ \sum_{k=1}^K \mu_k w_k u_{k,n,l}^* u_{k,n',l'} \mathbf{D}_n, & n' = n, l' \neq l, \\ \sum_{k=1}^K \mu_k w_k u_{k,n,l}^* u_{k,n',l'} \mathbf{D}_{n,n'}, & n' \neq n, \end{cases}$$

$$\mathbf{D}_n = \mathbf{F}_n^H \mathbf{F}_n + \sum_{k=1}^K \beta_k^2 \mathbf{F}_n^H \mathbf{h}_{n,k} \mathbf{h}_{n,k}^H \mathbf{F}_n,$$

$$\mathbf{D}_{n,n'} = \sum_{k=1}^K \beta_k^2 \mathbf{F}_n^H \mathbf{h}_{n,k} \mathbf{h}_{n',k}^H \mathbf{F}_{n'},$$

Figure 3: PFS utility versus per-RRH fronthaul capacity  $C$ 

and  $\mathbf{J} = [\mathbf{J}_{1,1}, \dots, \mathbf{J}_{N,L}]^T$  with  $\mathbf{J}_{n,l} = \sum_{k=1}^K \mu_k w_k \beta_k u_{k,n,l}^* \mathbf{F}_n^H \mathbf{h}_{n,k}, \forall n, l$ .

Finally, the optimal quantization bits allocation is obtained by solving the KKT conditions as

$$d_{n,l}^*(\lambda_n) = \left[ \frac{\log 2B\lambda_n - \log(\log 4 \sum_{k=1}^K \mu_k w_k s_{k,n,l})}{\log 4} \right]^+,$$

$\forall n, l$ , where  $s_{k,n,l} = 3|u_{k,n,l}|^2 (\sum_{k=1}^K p_k |\mathbf{h}_{n,k}^H \tilde{\mathbf{v}}_{n,l}|^2 + \|\tilde{\mathbf{v}}_{n,l}\|^2)$  and the optimal Lagrange multiplier  $\lambda_n \geq 0$  is chosen such that  $2B \sum_{l=1}^L d_{n,l}^*(\lambda_n) = C_n$ .

## V. SIMULATION RESULTS AND DISCUSSIONS

Consider a C-RAN with 4 RRHs placed in a circle cell of radius 500 m. There are 8 users randomly distributed in the cell. The channel bandwidth is 1 MHz. We adopt the geometry-based channel model in [8] for simulations. Unless otherwise specified, we consider  $M = 64$  antennas,  $S = 16$  RF chains for each RRH. There are  $T_s = 10$  time slots in each frame and the slot size is 1 ms. The coherence time for the channel statistics is 10 s. As in [9], we assume that the CSI delay is proportional to the dimension of the channel vector that is required at the BS. The carrier frequency is 2.14 GHz and the velocity of users is 3 Km/h. The CSI delay for the full channel matrix is set to be  $\tau = 4$  ms except for Fig. 4. We consider PFS utility. Three baseline schemes are considered for comparison: the spatial-compression-and-forward (SCF) scheme in [3], the analog SCF (A-SCF) scheme in [4] and the slow-timescale SCF (S-SCF) obtained by removing the short-term optimization in the proposed scheme.

Fig. 3 shows the performance comparison of different schemes versus per-RRH fronthaul capacity  $C$  varies from  $C = 16$  Mbps to  $C = 160$  Mbps. It can be observed that the best PFS utility is achieved by the SCF scheme without CSI delay, followed by the proposed THCF scheme. Furthermore, the proposed THCF scheme achieves significant gain over A-SCF and S-SCF, which demonstrates the importance of hybrid analog-and-digital processing and two-timescale joint optimization. Finally, it is observed that the performance of SCF is inferior to the proposed THCF since the full-CSI delay is larger than the effective-CSI delay. In Fig. 4, we plot the PFS utility versus the CSI delay, where the per-RRH fronthaul capacity is fixed as  $C = 64$  Mbps. We can

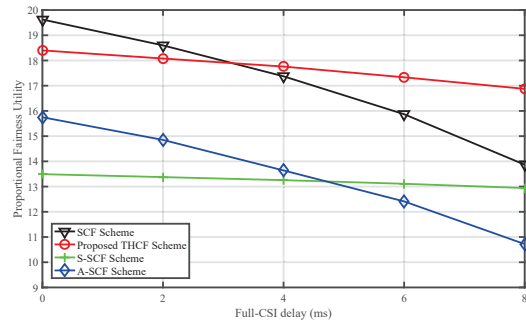


Figure 4: PFS utility versus the CSI delay.

see that as the CSI delay increases, the PFS of all schemes decreases gradually. It is observed that the PFS achieved with the proposed THCF scheme is higher than that achieved by the other schemes for moderate and large full-CSI delay. This is because the performance of the proposed THCF scheme is insensitive to the full-CSI delay. Although the performance of the S-SCF scheme is also insensitive to the full-CSI delay, its performance is still much worse than the proposed THCF scheme due to the lack of optimal power control and quantization bits allocation.

## VI. CONCLUSION

We propose a two-timescale hybrid compression and forward (THCF) scheme to reduce the fronthaul consumption in Massive MIMO aided C-RAN. We formulate the optimization of THCF as a general utility maximization problem, and propose a BC-SSCA algorithm to find stationary solutions of this two-stage non-convex stochastic optimization problem. Simulations verify that the proposed BC-SSCA algorithm achieves significant gain over existing solutions.

## REFERENCES

- [1] N. Chen, B. Rong, X. Zhang, and M. Kadoch, "Scalable and flexible massive MIMO precoding for 5G H-CRAN," *IEEE Wireless Communications*, vol. 24, no. 1, pp. 46–52, February 2017.
- [2] Y. Zhou and W. Yu, "Optimized backhaul compression for uplink cloud radio access network," *IEEE Journal on Selected Areas in Communications*, vol. 32, no. 6, pp. 1295–1307, June 2014.
- [3] L. Liu and R. Zhang, "Optimized uplink transmission in multi-antenna C-RAN with spatial compression and forward," *IEEE Transactions on Signal Processing*, vol. 63, no. 19, pp. 5083–5095, Oct 2015.
- [4] L. Combi and U. Spagnolini, "Hybrid beamforming in RoF fronthauling for millimeter-wave radio," in *2017 European Conference on Networks and Communications (EuCNC)*, June 2017, pp. 1–5.
- [5] X. Zhang, A. Molisch, and S.-Y. Kung, "Variable-phase-shift-based RF-baseband codesign for MIMO antenna selection," *IEEE Trans. Signal Processing*, vol. 53, no. 11, pp. 4091–4103, Nov. 2005.
- [6] A. Liu, X. Chen, W. Yu, V. Lau, and M.-J. Zhao, "Two-timescale hybrid compression and forward for massive MIMO aided C-RAN," *in preparation*, 2018.
- [7] Q. Shi, M. Razaviyayn, Z.-Q. Luo, and C. He, "An iteratively weighted MMSE approach to distributed sum-utility maximization for a MIMO interfering broadcast channel," *IEEE Trans. Signal Processing*, vol. 59, no. 9, pp. 4331–4340, Sept. 2011.
- [8] S. Park, J. Park, A. Yazdan, and R. W. Heath, "Exploiting spatial channel covariance for hybrid precoding in massive MIMO systems," *IEEE Trans. Signal Processing*, vol. 65, no. 14, pp. 3818–3832, July 2017.
- [9] A. Liu and V. K. N. Lau, "Phase only RF precoding for massive MIMO systems with limited RF chains," *IEEE Trans. Signal Processing*, vol. 62, no. 17, pp. 4505–4515, Sept. 2014.