

Fronthaul Compression and Transmit Beamforming Optimization for Multi-Antenna Uplink C-RAN

Yuhan Zhou, *Member, IEEE* and Wei Yu, *Fellow, IEEE*

Abstract—This paper considers the joint fronthaul compression and transmit beamforming design for the uplink cloud radio access network (C-RAN), in which multi-antenna user terminals communicate with a cloud-computing based centralized processor (CP) through multi-antenna base-stations (BSs) serving as relay nodes. A compress-and-forward relaying strategy, named the virtual multiple-access channel (VMAC) scheme, is employed, in which the BSs can either perform single-user compression or Wyner-Ziv coding to quantize the received signals and send the quantization bits to the CP via capacity-limited fronthaul links; the CP performs successive decoding with either successive interference cancellation (SIC) receiver or linear minimum-mean-square-error (MMSE) receiver. Under this setup, this paper investigates the joint optimization of the transmit beamformers at the users and the quantization noise covariance matrices at the BSs for maximizing the network utility. A novel weighted minimum-mean-square-error successive convex approximation (WMMSE-SCA) algorithm is first proposed for maximizing the weighted sum rate under the user transmit power and fronthaul capacity constraints with single-user compression. Assuming a heuristic decompression order, the proposed algorithm is then adapted for optimizing the transmit beamforming and fronthaul compression under Wyner-Ziv coding. This paper also proposes a low-complexity separate design consisting of optimizing transmit beamformers for the Gaussian vector multiple-access channel along with per-antenna quantizers with uniform quantization noise levels across the antennas at each BS. Numerical results show that with optimized beamforming and fronthaul compression, C-RAN can significantly improve the overall performance of conventional cellular networks. Majority of the performance gain comes from the implementation of SIC at the central receiver. Furthermore, the low complexity separate design already performs very close to the optimized joint design in regime of practical interest.

Index Terms—Cloud radio access network, fronthaul compression, transmit beamforming, compress-and-forward, linear MMSE receiver, SIC receiver, network MIMO.

I. INTRODUCTION

To meet the exponentially increasing data demand in wireless communication driven by smartphones, tablets, and video

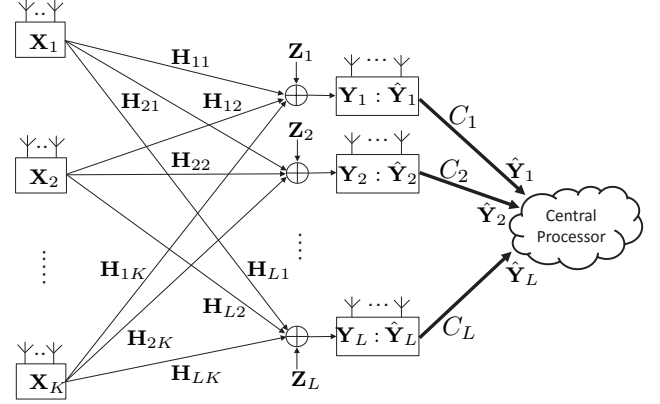


Fig. 1. An uplink C-RAN system with capacity-limited fronthaul.

streaming, modern cellular communication systems are moving towards densely deployed heterogeneous networks consisting of base-stations (BSs) covering progressively smaller areas. As a consequence, inter-cell interference becomes the dominant performance limiting factor. Cloud radio access network (C-RAN) is a novel mobile network architecture that offers an efficient way for managing inter-cell interference [1]. In a C-RAN architecture, the baseband and higher-layers operations of the BSs are migrated to a cloud-computing based centralized processor (CP). By allowing coordination and joint signal processing across multiple cells, C-RAN provides a platform for implementing network multiple-input multiple-output (network MIMO), also known as coordinated multi-point (CoMP), which can achieve significantly higher data rates than conventional cellular networks [2].

This paper focuses on the uplink C-RAN architecture as shown in Fig. 1, where multi-antenna mobile users communicate with the CP with multi-antenna BSs serving as relay nodes. The BSs are connected with the CP via digital fronthaul links with finite capacities. We consider a two-stage compress-and-forward relaying strategy, referred to as the virtual multiple access channel (VMAC) scheme, in which the BSs quantize the received signals using either single-user compression or Wyner-Ziv coding and send the compressed bits to the CP. The CP performs successive decoding to decode the quantization codewords first, then the user messages sequentially. Under the VMAC scheme, this paper studies the optimization of the transmit beamforming vectors and quantization noise covariance matrices for maximizing the weighted sum rate of the C-RAN system. Being different from the conventional multicell cellular systems, in which the optimal transmit beamforming only depends on the interfering signal strength and the channel gain matrices, in C-RAN,

Copyright (c) 2015 IEEE. Personal use of this material is permitted. However, permission to use this material for any other purposes must be obtained from the IEEE by sending a request to pubs-permissions@ieee.org.

Manuscript received October 6, 2015; revised March 6, 2016; accepted April 15, 2016. This work was supported in part by Huawei Technologies Canada Co., Ltd., and in part by the Natural Sciences and Engineering Research Council (NSERC) of Canada. This paper was presented in part at IEEE Globecom Workshop, Austin, Texas, USA, December 2014. The associate editor coordinating the review of this manuscript and approving it for publication was Dr. Paolo Banelli.

Y. Zhou was with the The Edward S. Rogers Sr. Department of Electrical and Computer Engineering, University of Toronto, Toronto, ON M5S 3G4 Canada. He is now with Qualcomm Technologies Inc., San Diego, CA 92121 USA (email: yzhou@ece.utoronto.ca).

W. Yu is with The Edward S. Rogers Sr. Department of Electrical and Computer Engineering, University of Toronto, Toronto, ON M5S 3G4 Canada (e-mail: weiyu@ece.utoronto.ca).

the finite fronthaul capacity also needs to be taken into account in the beamforming design. This paper proposes a novel weighted minimum-mean-square-error successive convex approximation (WMMSE-SCA) algorithm to find the optimal transmit beamformers and quantization noise covariance matrices for maximizing the weighted sum rate of C-RAN. Moreover, a simple separate design consisting of optimizing transmit beamformers for the Gaussian vector multiple-access channel and per-antenna quantizers with uniform quantization noise levels across the antennas at each BS is also developed, under the assumption that the signal-to-quantization-noise ratio (SQNR) is high and successive interference cancellation (SIC) is applied at the receiver. Numerical simulations show that the proposed separate design already performs very close to the optimized joint design in the SQNR regime of practical interest.

This paper considers two different fronthaul compression strategies for C-RAN, namely *single-user compression* and *Wyner-Ziv coding*. In single-user compression, which is also referred to as point-to-point compression in the literature [3], each BS uses vector quantization to compress the received signals but ignores the correlation between the received signals across different BSs. In contrast, Wyner-Ziv coding fully utilizes the correlation of the received signals for higher compression efficiency, thereby achieving better overall performance. The optimization strategy proposed in this paper is first developed for single-user compression, then for the more complex Wyner-Ziv coding, assuming a heuristic ordering for decompression of the quantized signals at the BSs. The performance of the VMAC schemes with single-user compression and Wyner-Ziv coding are evaluated for practical multicell networks under linear *minimum-mean-square-error (MMSE) receiver* and *SIC receiver* respectively. It is shown that the implementation of SIC receiver significantly improves the performance achieved by linear MMSE receiver under both single-user compression and Wyner-Ziv coding. Furthermore, although single-user compression with SIC receiver can already realize majority of the benefit brought by the C-RAN architecture, Wyner-Ziv coding can further improve upon single-user compression when the fronthaul capacity is limited.

To precisely quantify the advantage of the C-RAN architecture, this paper further evaluates the performance of optimized beamforming and fronthaul compression under two different types of BS clustering strategies: *disjoint clustering* and *user-centric clustering*. In disjoint clustering scheme, the entire network is divided into non-overlapping clusters and the BSs in each cluster jointly serve all the users within the coverage area [4], [5]. In user-centric clustering, each user is served by an individually selected subset of neighboring BSs; different clusters for different users may overlap. The performance of user-centric clustering has been evaluated for the downlink of cooperative cellular networks [6] and C-RAN systems [7]. This paper further shows numerically that in uplink C-RAN, with optimized beamforming and fronthaul compression, the user-centric clustering strategy significantly outperforms the disjoint clustering strategy, because the cell edges are effectively eliminated.

A. Related Work

One of the main issues in the implementation of C-RAN is how to optimally utilize the capacity-limited fronthaul links to efficiently reap the benefit of multicell processing. Substantial research works have made progress towards this direction [8], [9], [10]. Under the assumption of compress-and-forward relaying strategy at the BSs, the largest achievable rate for uplink C-RAN is given by the joint decompression and decoding strategy, in which quantization codewords and user messages are decoded simultaneously [11], [12], [13]. However, the complexity of joint decompression and decoding is very high, which prevents it from practical implementation. In [14], a virtual multiple access channel (VMAC) scheme, which is a compress-and-forward strategy based on the successive decoding of quantization codewords followed by user messages, is proposed for the single-input single-output (SISO) C-RAN architecture. As compared to the joint decompression and decoding scheme, the VMAC scheme has lower decoding complexity and shorter decoding delay, which makes it more desirable for practical implementation. Furthermore, it is shown in [15] that with Wyner-Ziv coding the successive decoding based VMAC scheme actually achieves the same maximum sum rate as the joint decompression and decoding strategy for the uplink C-RAN model under a sum fronthaul constraint.

This paper studies the linear transceiver and fronthaul compression design in the VMAC scheme for the uplink multiple-input-multiple-output (MIMO) C-RAN model. As a generalization of [14] which considers the SISO case only, this paper considers the MIMO case where both the users and the BSs are equipped with multiple antennas. The main difference between the SISO case and the MIMO case is the impact of transmitter optimization at the user terminals. For example, in the SISO case, to maximize the sum rate of the VMAC scheme for a single-cluster C-RAN system, all the users within the cluster need to transmit at their maximum powers (assuming that SIC receiver is implemented at the CP). However, in the MIMO case, the users are capable of doing transmit beamforming, so the optimal transmit beamforming design is more involved.

The fronthaul compression problem for the uplink C-RAN model has been considered extensively in the literature. Various algorithms such as alternating convex optimization [14], gradient projection [16], and the robust fronthaul compression approach [17] have been developed for maximizing the (weighted) sum rate under the fronthaul constraints. All of these algorithms focus only on the optimization of quantization noise covariance matrices across the BSs, with fixed transmit beamformers. This paper goes one step further by considering the joint transmit beamformer and quantization noise covariance matrix optimization problem. Accounting for both the transmit beamforming and the quantization design problem together in the optimization framework is nontrivial, because the two are coupled through the fronthaul constraints. To tackle this problem, this paper proposes a novel WMMSE-SCA algorithm for efficiently finding a local optimum solution to the weighted sum rate maximization problem. The proposed

algorithm integrates the well-known WMMSE beamforming design strategy [18], [19], with the successive convex approximation technique [20], [21], to arrive at a stationary point of the maximization problem. The performance of optimized beamforming vectors and quantization noise covariance matrices for both Wyner-Ziv coding and single-user compression are evaluated under practical multicell networks with different receive beamforming schemes, i.e., the linear MMSE receiver and the SIC receiver. Simulation results show that the performance improvement of the SIC receiver as compared to the linear MMSE receiver is much larger than that of Wyner-Ziv coding as compared to single-user compression. Most of the performance gain brought by C-RAN can thus be obtained by single-user compression together with SIC receiver.

From a broader perspective, this paper studies the radio resource allocation optimization for uplink fronthaul-constrained C-RAN [5]. As related work, we mention [22] which proposes to utilize the signal sparsity in C-RAN to improve the performance of the fronthaul compression and user detection. The fronthaul compression can also be designed for enhancing synchronization in C-RAN [23]. Additionally, we mention briefly that the latency issue in the C-RAN design has been studied in [24], which introduces a delay-optimal fronthaul allocation strategy for the latency control. Finally, we point out that with proper modification the design idea of fronthaul-aware beamforming proposed in this paper can be extended to the more recently proposed heterogeneous C-RAN [25] and Fog RAN [26] architectures.

B. Paper Organization and Notation

The rest of the paper is organized as follows. Section II introduces the system model and the VMAC scheme. Section III considers the joint design of beamforming and fronthaul compression under single-user compression, where a novel WMMSE based successive convex optimization algorithm is proposed. The proposed joint design scheme is developed further in Section IV for maximizing weighted sum rate under Wyner-Ziv coding. Section V is devoted to a low-complexity separate beamforming and fronthaul compression design, which is shown to be near-optimal at high SQNR regime. The proposed algorithms are evaluated numerically for practical multicell and multicluster networks in Section VI. Conclusions are drawn in Section VII.

The notations used in this paper are as follows. Boldface lower-case letters denote column vectors. Boldface upper-case letters denote vector random variables or matrices, where context should make the distinction clear. Superscripts $(\cdot)^T$, $(\cdot)^\dagger$, and $(\cdot)^{-1}$ denote transpose, Hermitian transpose, and matrix inverse operators; $\mathbb{E}(\cdot)$, $\|\cdot\|_F$, $\text{Tr}(\cdot)$, and $\text{rank}(\cdot)$ denote expectation, Frobenius norm, matrix trace, and matrix rank operators. For a vector/matrix \mathbf{X} , \mathbf{X}_S denotes a vector/matrix with elements whose indices are elements of S . Given matrices $\{\mathbf{X}_1, \dots, \mathbf{X}_L\}$, $\text{diag}(\{\mathbf{X}_\ell\}_{\ell=1}^L)$ denotes the block diagonal matrix formed with \mathbf{X}_ℓ on the diagonal. We let $\mathcal{K} = \{1, \dots, K\}$ and $\mathcal{L} = \{1, \dots, L\}$.

II. PRELIMINARIES

A. System Model

This paper considers the uplink C-RAN, where K multi-antenna mobile users communicate with a CP through L multi-antenna BSs serving as relay nodes, as shown in Fig. 1. The noiseless fronthaul links connecting the BSs with the CP have per-link capacity of C_ℓ . Each user terminal is equipped with M antennas; each BS is equipped with N antennas.

We consider the VMAC scheme [14] applied to such a C-RAN system, in which the BSs quantize the received signals using either Wyner-Ziv coding or single-user compression, then forward the compressed bits to the CP for decoding. In single-user compression, the compression process only involves the conventional vector quantizers, one for each BS, while in Wyner-Ziv coding, the correlation between the received signals across the BSs are fully utilized for higher compression efficiency. At the CP side, a two-stage successive decoding strategy is employed, where the quantization codewords are first decoded, and then the user messages are decoded sequentially.

We assume that the wireless channel between the users and the BSs are quasi-static, with large coherence time/bandwidth. Furthermore, perfect channel state information (CSI) is assumed to be available at the BSs. The BSs forward the CSI to the CP through fronthaul links. The CP determines the optimal quantization noise covariance matrices, then feeds them back to the BSs. The CP also determines the user transmit beamformers and feeds them back to the users. Under the quasi-static channel, the overheads due to the CSI transmission and feedback between the BSs and the CP can be amortized within the channel coherence time and is ignored in this paper for simplicity.

Define $\mathbf{H}_{\ell,k}$ as the $N \times M$ complex channel matrix between the k th user and the ℓ th BS. It is assumed that each user intends to transmit d parallel data streams to the CP. (Throughout this paper we assume that d is fixed.) Let $\mathbf{V}_k \in \mathbb{C}^{M \times d}$ denote the linear transmit beamformer that user k utilizes to transmit message signal $\mathbf{s}_k \in \mathbb{C}^{d \times 1}$ to the central receiver. We assume that each message signal \mathbf{s}_k intended for user k is taken from a Gaussian codebook so that we have $\mathbf{s}_k \sim \mathcal{CN}(\mathbf{0}, \mathbf{I})$. Then the transmit signal at user k is given by $\mathbf{x}_k = \mathbf{V}_k \mathbf{s}_k$ with $\mathbb{E}[\mathbf{x}_k \mathbf{x}_k^\dagger] = \mathbf{V}_k \mathbf{V}_k^\dagger$. The transmit beamformers are subjected to per-user power constraints, i.e., $\text{Tr}(\mathbf{V}_k \mathbf{V}_k^\dagger) \leq P_k$ for $k \in \mathcal{K}$. The received signal at BS ℓ , \mathbf{y}_ℓ , can be expressed as

$$\mathbf{y}_\ell = \sum_{k=1}^K \mathbf{H}_{\ell,k} \mathbf{V}_k \mathbf{s}_k + \mathbf{z}_\ell, \quad \forall \ell \in \mathcal{L}, \quad (1)$$

where $\mathbf{z}_\ell \sim \mathcal{CN}(\mathbf{0}, \mathbf{\Lambda}_\ell)$ represents the additive Gaussian noise for BS ℓ . Assuming Gaussian quantization test channel, the quantized received signal $\hat{\mathbf{y}}_\ell$ for the ℓ th BS is given by

$$\hat{\mathbf{y}}_\ell = \mathbf{y}_\ell + \mathbf{q}_\ell \quad (2)$$

where $\mathbf{q}_\ell \sim \mathcal{CN}(\mathbf{0}, \mathbf{Q}_\ell)$ represents the Gaussian quantization noise for the ℓ th BS.

The above Gaussian quantization test channel model (2) is sufficiently general to encompass the possibility of per-

forming receive beamforming at the BSs prior to quantization. For a given quantization noise covariance matrix \mathbf{Q}_ℓ , let $\mathbf{Q}_\ell = \mathbf{B}_\ell \Phi_\ell \mathbf{B}_\ell^\dagger$ be its eigenvalue composition, where $\Phi_\ell = \text{diag}(\{\varphi_\ell^i\}_{i=1}^N)$ with φ_ℓ^i as the i th eigenvalue of \mathbf{Q}_ℓ , and \mathbf{B}_ℓ is the $N \times N$ unitary matrix whose i th column is the eigenvector \mathbf{b}_ℓ^i corresponding to eigenvalue φ_ℓ^i . Then, setting quantization noise covariance to be \mathbf{Q}_ℓ is equivalent to beamforming with \mathbf{B}_ℓ^\dagger followed by independent quantization, i.e.,

$$\hat{\mathbf{y}}'_\ell = \mathbf{B}_\ell^\dagger \mathbf{y}_\ell + \tilde{\mathbf{q}}_\ell \quad (3)$$

where $\mathbf{B}_\ell = [\mathbf{b}_\ell^1, \dots, \mathbf{b}_\ell^N]$ is the receive beamforming at the ℓ th BS, and $\tilde{\mathbf{q}}_\ell \sim \mathcal{CN}(\mathbf{0}, \Phi_\ell)$ represents the quantization noise after the beamforming.

B. Achievable Rate of the VMAC scheme

The rate region of the VMAC scheme is characterized by that of a multiple-access channel, in which multiple users send information to a common CP. Following the results in [14], assuming that the linear MMSE receiver is applied at the CP, the transmission rate R_k for user k for the VMAC scheme is given by

$$\begin{aligned} R_k &\leq I(\mathbf{X}_k; \mathbf{Y}_1, \dots, \mathbf{Y}_L) \\ &= \log \left| \mathbf{I} + \mathbf{V}_k^\dagger \mathbf{H}_{\mathcal{L},k}^\dagger \mathbf{J}_k^{-1} \mathbf{H}_{\mathcal{L},k} \mathbf{V}_k \right| \end{aligned} \quad (4)$$

where

$$\mathbf{J}_k = \mathbf{J}_k^{\text{LE}} = \sum_{j \neq k} \mathbf{H}_{\mathcal{L},j} \mathbf{V}_j \mathbf{V}_j^\dagger \mathbf{H}_{\mathcal{L},j}^\dagger + \mathbf{\Lambda} + \mathbf{Q}, \quad (5)$$

with $\mathbf{\Lambda} = \text{diag}(\{\mathbf{\Lambda}_\ell\}_{\ell=1}^L)$ and $\mathbf{Q} = \text{diag}(\{\mathbf{Q}_\ell\}_{\ell=1}^L)$. To achieve higher throughput, the SIC scheme can also be combined with the linear MMSE receiver. In this case, assuming without loss of generality a decoding order of user messages $1, 2, \dots, K$, the matrix $\mathbf{J}_k = \mathbf{J}_k^{\text{LE}}$ in (4) is replaced by $\mathbf{J}_k^{\text{SIC}}$ expressed as

$$\mathbf{J}_k = \mathbf{J}_k^{\text{SIC}} = \sum_{j > k} \mathbf{H}_{\mathcal{L},j} \mathbf{V}_j \mathbf{V}_j^\dagger \mathbf{H}_{\mathcal{L},j}^\dagger + \mathbf{\Lambda} + \mathbf{Q}. \quad (6)$$

The compression rates at the BSs should also satisfy the fronthaul link capacity constraints. Using information-theoretic formulation, the fronthaul constraints under single-user compression can be written as

$$I(\mathbf{Y}_\ell; \hat{\mathbf{Y}}_\ell) \leq C_\ell, \quad \forall \ell \in \mathcal{L}. \quad (7)$$

Evaluating the above mutual information expression with Gaussian input and Gaussian quantization noise, the fronthaul constraint (7) becomes [27]

$$\log \frac{\left| \sum_{k=1}^K \mathbf{H}_{\ell,k} \mathbf{V}_k \mathbf{V}_k^\dagger \mathbf{H}_{\ell,k}^\dagger + \mathbf{\Lambda}_\ell + \mathbf{Q}_\ell \right|}{|\mathbf{Q}_\ell|} \leq C_\ell, \quad (8)$$

for all $\ell = 1, 2, \dots, L$. When Wyner-Ziv coding is implemented at BSs, the fronthaul constraints are given by the following mutual information expressions [15], [28]

$$I(\mathbf{Y}_S; \hat{\mathbf{Y}}_S | \hat{\mathbf{Y}}_{S^c}) \leq \sum_{\ell \in S} C_\ell, \quad \forall S \subseteq \mathcal{L}. \quad (9)$$

Utilizing the chain rule on mutual information and the Gaussian assumption, one can express the fronthaul constraint (9) for Wyner-Ziv coding as follows,

$$\begin{aligned} &\log \frac{\left| \sum_{k=1}^K \mathbf{H}_{\mathcal{L},k} \mathbf{V}_k \mathbf{V}_k^\dagger \mathbf{H}_{\mathcal{L},k}^\dagger + \text{diag}(\{\mathbf{\Lambda}_\ell + \mathbf{Q}_\ell\}_{\ell \in \mathcal{L}}) \right|}{\left| \sum_{k=1}^K \mathbf{H}_{S^c,k} \mathbf{V}_k \mathbf{V}_k^\dagger \mathbf{H}_{S^c,k}^\dagger + \text{diag}(\{\mathbf{\Lambda}_\ell + \mathbf{Q}_\ell\}_{\ell \in S^c}) \right|} \\ &\quad - \sum_{\ell \in S} \log |\mathbf{Q}_\ell| \leq \sum_{\ell \in S} C_\ell, \quad \forall S \subseteq \mathcal{L}. \end{aligned} \quad (10)$$

III. JOINT BEAMFORMING AND COMPRESSION DESIGN UNDER SINGLE-USER COMPRESSION

A. Weighted Sum Rate Maximization

This section investigates the joint beamforming and fronthaul compression design for the VMAC scheme with single-user compression. As shown in the achievable rate expression (4) and the fronthaul constraint expression (8), the beamforming vectors and quantization noise covariance matrices are coupled, and the two together determine the overall performance of a C-RAN system. To characterize the tradeoff between the achievable rates for the users and the system resources, we formulate the following weighted sum rate maximization problem:

$$\max_{\mathbf{V}_k, \mathbf{Q}_\ell} \sum_{k=1}^K \alpha_k \log \left| \mathbf{I} + \mathbf{V}_k^\dagger \mathbf{H}_{\mathcal{L},k}^\dagger \mathbf{J}_k^{-1} \mathbf{H}_{\mathcal{L},k} \mathbf{V}_k \right| \quad (11a)$$

$$\text{s.t.} \quad \mathbf{J}_k = \mathbf{J}_k^{\text{LE}} \quad \text{or} \quad \mathbf{J}_k = \mathbf{J}_k^{\text{SIC}},$$

$$\log \frac{\left| \sum_{k=1}^K \mathbf{H}_{\ell,k} \mathbf{V}_k \mathbf{V}_k^\dagger \mathbf{H}_{\ell,k}^\dagger + \mathbf{\Lambda}_\ell + \mathbf{Q}_\ell \right|}{|\mathbf{Q}_\ell|} \leq C_\ell, \quad (11b)$$

$$\mathbf{Q}_\ell \succeq \mathbf{0}, \quad \forall \ell \in \mathcal{L},$$

$$\text{Tr}(\mathbf{V}_k \mathbf{V}_k^\dagger) \leq P_k, \quad \forall k \in \mathcal{K}$$

where α_k 's are the weights representing the priorities associated with the mobile users typically determined from upper layer protocols. When SIC receiver is implemented, to maximize the weighed sum rate, the user with larger weight should be decoded last. Without loss of generality, we assume $0 \leq \alpha_1 \leq \alpha_2 \leq \dots \leq \alpha_K$, which results in the decoding order of user messages $1, 2, \dots, K$.

Due to the non-convexity of both the objective function and the fronthaul capacity constraints in problem (11), finding the global optimum solution of (11) is challenging. We point out here that the present formulation (11) actually implicitly includes the user scheduling strategy. More specifically, one can consider a weighted sum rate maximization problem over all the users in the network, where the beamformers for the users are set to be the zero vector if they are not scheduled. For simplicity in the following development, we focus on the active users only and assume that user scheduling is done prior to solving problem (11). Implicit scheduling is discussed later in the simulation part of the paper.

B. The WMMSE-SCA Algorithm

In this section, we propose a novel algorithm to find a stationary point of the problem (11). The main difficulty in

solving (11) comes from the fact that the objective function and fronthaul capacity constraints are both nonconvex functions with respect to the optimization variables. Inspired by the recent work of using the WMMSE approach for beamforming design [18], [19], we first reformulate the objective function in problem (11) as a convex function with respect to the MMSE matrix given by the user's target signal s_k and decoded signal \hat{s}_k . We then linearize the convex objective function and the compression rate expressions in the fronthaul constraints of (11) to obtain a convex approximation of the original problem. Finally we successively approximate the optimal solution by optimizing this convex approximation. The idea of convex approximation is rooted from modern optimization techniques including block successive minimization method and minorize-maximization algorithm, which have been previously applied for solving related problems in wireless communications [6], [29].

Before presenting the proposed algorithm, we first state the following lemma, which is a direct consequence of concavity of the $\log|\cdot|$ function.

Lemma 1: For positive definite Hermitian matrices $\Omega, \Sigma \in \mathbb{C}^{N \times N}$,

$$\log|\Omega| \leq \log|\Sigma| + \text{Tr}(\Sigma^{-1}\Omega) - N \quad (12)$$

with equality if and only if $\Omega = \Sigma$.

By applying Lemma 1 to the first log-determinant term in the fronthaul constraint expression (8) or (11b) and by setting

$$\Omega_\ell = \sum_{k=1}^K \mathbf{H}_{\ell,k} \mathbf{V}_k \mathbf{V}_k^\dagger \mathbf{H}_{\ell,k}^\dagger + \mathbf{A}_\ell + \mathbf{Q}_\ell, \quad (13)$$

we can approximate the fronthaul constraint (8) or (11b) with the following convex constraint:

$$\log|\Sigma_\ell| + \text{Tr}\left(\Sigma_\ell^{-1} \left(\sum_{k=1}^K \mathbf{H}_{\ell,k} \mathbf{V}_k \mathbf{V}_k^\dagger \mathbf{H}_{\ell,k}^\dagger + \mathbf{A}_\ell + \mathbf{Q}_\ell\right)\right) - \log|\mathbf{Q}_\ell| \leq C_\ell + N \quad (14)$$

for $\ell = 1, 2, \dots, L$. It is not hard to see that the fronthaul constraint (8) or (11b) is always feasible when the convex constraint (14) is feasible. The two constraints are equivalent when

$$\Sigma_\ell^* = \sum_{k=1}^K \mathbf{H}_{\ell,k} \mathbf{V}_k \mathbf{V}_k^\dagger \mathbf{H}_{\ell,k}^\dagger + \mathbf{A}_\ell + \mathbf{Q}_\ell. \quad (15)$$

Now we approximate the objective function (11a) using the WMMSE approximation. Let $\mathbf{U}_k \in \mathbb{C}^{N \times d}$ be the linear receiver applied at the CP for recovering s_k . The transmission rate R_k in (4) can be expressed as the following [18] [19],

$$R_k = \max_{\mathbf{U}_k} \log|\mathbf{E}_k^{-1}| \quad (16)$$

where

$$\mathbf{E}_k = (\mathbf{I} - \mathbf{U}_k^\dagger \mathbf{H}_{\mathcal{L},k} \mathbf{V}_k)(\mathbf{I} - \mathbf{U}_k^\dagger \mathbf{H}_{\mathcal{L},k} \mathbf{V}_k)^\dagger + \mathbf{U}_k^\dagger \left(\sum_{j \neq k} \mathbf{H}_{\mathcal{L},j} \mathbf{V}_j \mathbf{V}_j^\dagger \mathbf{H}_{\mathcal{L},j}^\dagger + \mathbf{A} + \mathbf{Q} \right) \mathbf{U}_k. \quad (17)$$

By applying Lemma 1 again, we rewrite rate expression (16) as

$$R_k = \max_{\mathbf{W}_k, \mathbf{U}_k} (\log|\mathbf{W}_k| - \text{Tr}(\mathbf{W}_k \mathbf{E}_k) + d) \quad (18)$$

where \mathbf{W}_k is the weight matrix introduced by the WMMSE method. The optimal \mathbf{W}_k is given by

$$\mathbf{W}_k^* = \mathbf{E}_k^{-1} = \mathbf{I} + \mathbf{V}_k^\dagger \mathbf{H}_{\mathcal{L},k}^\dagger \mathbf{U}_k^*, \quad (19)$$

where \mathbf{U}_k^* is the MMSE receive beamformer given by

$$\mathbf{U}_k^* = \left(\sum_{j \neq k} \mathbf{H}_{\mathcal{L},j} \mathbf{V}_j \mathbf{V}_j^\dagger \mathbf{H}_{\mathcal{L},j}^\dagger + \mathbf{A} + \mathbf{Q} \right)^{-1} \mathbf{H}_{\mathcal{L},k} \mathbf{V}_k. \quad (20)$$

Using (18) and (14) to replace the objective function and the fronthaul constraints in problem (11), we reformulate the weighted sum-rate maximization problem as follows

$$\begin{aligned} \max_{\mathbf{V}_k, \mathbf{Q}_\ell, \mathbf{U}_k, \mathbf{W}_k, \Sigma_\ell, \Theta_\ell} \quad & \sum_{k=1}^K \alpha_k (\log|\mathbf{W}_k| - \text{Tr}(\mathbf{W}_k \mathbf{E}_k)) \\ & - \rho \sum_{\ell=1}^L \|\mathbf{Q}_\ell - \Theta_\ell\|_F^2 \\ \text{s.t.} \quad & \log|\Sigma_\ell| + \text{Tr}(\Sigma_\ell^{-1} \Omega_\ell) - \log|\mathbf{Q}_\ell| \leq C'_\ell, \quad \forall \ell \in \mathcal{L}, \\ & \mathbf{Q}_\ell \succeq \mathbf{0}, \quad \Theta_\ell \succeq \mathbf{0}, \quad \forall \ell \in \mathcal{L}, \\ & \text{Tr}(\mathbf{V}_k \mathbf{V}_k^\dagger) \leq P_k, \quad \forall k \in \mathcal{K}, \end{aligned} \quad (21)$$

where $\Omega_\ell = \sum_{k=1}^K \mathbf{H}_{\ell,k} \mathbf{V}_k \mathbf{V}_k^\dagger \mathbf{H}_{\ell,k}^\dagger + \mathbf{A}_\ell + \mathbf{Q}_\ell$, ρ is some positive constant, and $C'_\ell = C_\ell + N$. Note that the last term in the objective function which involves a summation of Frobenius norms is a quadratic regularization term. It makes the optimization problem (21) strictly convex with respect to each optimization variable.

It is easy to verify that problem (21) is convex with respect to any one of the optimization variables when the other optimization variables are fixed. Specifically, when the other variables are fixed, the optimal values of Σ_ℓ , \mathbf{W}_k , and \mathbf{U}_k are given by equations (15), (19), and (20) respectively. The optimal values of Φ_ℓ are given by $\Theta_\ell = \mathbf{Q}_\ell$. When Σ_ℓ , \mathbf{U}_k , and \mathbf{W}_k are fixed, the optimal values of \mathbf{V}_k and \mathbf{Q}_ℓ are solutions to the following optimization problem:

$$\begin{aligned} \min_{\mathbf{V}_k, \mathbf{Q}_\ell} \quad & \sum_{k=1}^K \alpha_k \text{Tr}(\mathbf{W}_k \mathbf{E}_k) + \rho \sum_{\ell=1}^L \|\mathbf{Q}_\ell - \Theta_\ell\|_F^2 \\ \text{s.t.} \quad & \text{Tr}(\Sigma_\ell^{-1} \Omega_\ell) - \log|\mathbf{Q}_\ell| \leq C'_\ell - \log|\Sigma_\ell|, \quad \forall \ell \in \mathcal{L}, \\ & \mathbf{Q}_\ell \succeq \mathbf{0}, \quad \forall \ell \in \mathcal{L}, \\ & \text{Tr}(\mathbf{V}_k \mathbf{V}_k^\dagger) \leq P_k, \quad \forall k \in \mathcal{K}, \end{aligned} \quad (22)$$

where $\Omega_\ell = \sum_{k=1}^K \mathbf{H}_{\ell,k} \mathbf{V}_k \mathbf{V}_k^\dagger \mathbf{H}_{\ell,k}^\dagger + \mathbf{A}_\ell + \mathbf{Q}_\ell$. The above problem is convex with respect to \mathbf{V}_k and \mathbf{Q}_ℓ , and can be solved efficiently with polynomial complexity. Standard convex optimization solver such as CVX [30] can be used for solving problem (22) numerically. We summarize the proposed WMMSE-SCA algorithm for single-user compression as Algorithm 1.

Algorithm 1 WMMSE-SCA Algorithm

-
- 1: Initialize \mathbf{Q}_ℓ and \mathbf{V}_k such that $\text{Tr}(\mathbf{V}_k \mathbf{V}_k^\dagger) = P_k$.
 - 2: $\Sigma_\ell \leftarrow \sum_{k=1}^K \mathbf{H}_{\ell,k} \mathbf{V}_k \mathbf{V}_k^\dagger \mathbf{H}_{\ell,k}^\dagger + \mathbf{\Lambda}_\ell + \mathbf{Q}_\ell$.
 - 3: $\mathbf{U}_k \leftarrow \left(\sum_{j \neq k} \mathbf{H}_{\mathcal{L},j} \mathbf{V}_j \mathbf{V}_j^\dagger \mathbf{H}_{\mathcal{L},j}^\dagger + \mathbf{\Lambda} + \mathbf{Q} \right)^{-1} \mathbf{H}_{\mathcal{L},k} \mathbf{V}_k$.
 - 4: $\mathbf{W}_k \leftarrow \mathbf{I} + \mathbf{V}_k^\dagger \mathbf{H}_{\mathcal{L},k}^\dagger \mathbf{U}_k$ and $\Theta_\ell \leftarrow \mathbf{Q}_\ell$.
 - 5: Fix Σ_ℓ , \mathbf{U}_k , \mathbf{W}_k , and Θ_ℓ solve the convex optimization problem (22). Set $(\mathbf{V}_k, \mathbf{Q}_\ell)$ to be its optimal solution.
 - 6: Repeat Steps 2–5, until convergence.
-

C. Convergence and Complexity Analysis

The WMMSE-SCA algorithm yields a nondecreasing sequence of objective values for problem (11). So the algorithm is guaranteed to converge. Moreover, it converges to a stationary point of the optimization problem. The convergence result is stated in Theorem 1.

Theorem 1: From any initial point $(\mathbf{V}_k^{(0)}, \mathbf{Q}_\ell^{(0)})$, the proposed WMMSE-SCA algorithm is guaranteed to converge. The limit point $(\mathbf{V}_k^*, \mathbf{Q}_\ell^*)$ generated by the WMMSE-SCA algorithm is a stationary point of the weighted sum-rate maximization problem (11).

Proof: See Appendix A. ■

We point out here that Theorem 1 can also be proved following a similar procedure as that for demonstrating the convergence of WMMSE algorithm [19]. Specifically, it follows from the general optimization theory [31, Theorem 2.7.1] that the WMMSE-SCA algorithm, which does block coordinate descent on the reformulated problem (21), converges to a stationary point of (21). Then one can show every stationary point of (21) is also a stationary point of the original maximization problem (11), thereby establishing the claim in Theorem 1. However such a proof is not as simple as the proof presented in this paper which utilizes the convergence result of the successive convex approximation algorithm [32]. We also emphasize the importance of the regularization term involving sum of Frobenius norms in the objective function of (21). The regularization term makes the objective function in (21) a strongly convex function with respect to $(\mathbf{V}_k, \mathbf{Q}_\ell)$, therefore guaranteeing the convergence of Algorithm 1.

Assuming a typical network with $K > L > N > M$, the computational complexity of the proposed WMMSE-SCA algorithm is dominated by the joint optimization of $(\mathbf{V}_k, \mathbf{Q}_\ell)$, i.e. Step 5 of Algorithm 1. Step 5 solves a convex optimization problem, which can be efficiently implemented by primal-dual interior point method with approximate complexity of $\mathcal{O}((KMd + LN)^{3.5})$ [33]. Suppose that Algorithm 1 takes T total number of iterations to converge, the overall computational complexity of Algorithm 1 is therefore $\mathcal{O}((KMd + LN)^{3.5}T)$.

IV. JOINT BEAMFORMING AND COMPRESSION OPTIMIZATION UNDER WYNER-ZIV CODING

In single-user compression, the compression and decompression across different BSs take place independently. This separate processing neglects the key fact that the received

signals \mathbf{y}_ℓ in (1) are statistically correlated across the BS index ℓ , since they are noisy observations of the same transmitted signals \mathbf{x}_k . Based on this fact, Wyner-Ziv coding, which jointly decompresses the signals at the CP, is expected to be superior to the pre-link single-user compression in utilizing the limited fronthaul capacities. With fixed transmitters, the advantages of Wyner-Ziv coding have been demonstrated in [3], [14]. We take one step further in this section to study the problem of jointly optimizing transmit beamforming vectors and Wyner-Ziv quantization noise covariance matrices for the VMAC scheme in uplink C-RAN.

In the implementation of Wyner-Ziv coding, we decompress the quantization codeword $\hat{\mathbf{y}}_\ell$ sequentially from one BS to the other. To this end, we need to determine a decompression order on the BS indices $\{1, 2, \dots, L\}$. The decompression order generally affects the achievable performance of the VMAC scheme and should be optimized. However, in order to determine the optimal order that results in the largest weighted sum rate (or the maximum network utility) for the uplink C-RAN model shown in Fig. 1, we need to exhaustively search over $L!$ different decompression orders, which is impractical for large L . To tackle this problem, we propose a heuristic order of decompressing first the signals from the BS with larger value of

$$C_\ell - \log \left| \mathbf{H}_{\ell,\mathcal{K}} \tilde{\mathbf{K}} \mathbf{H}_{\ell,\mathcal{K}}^\dagger + \mathbf{\Lambda}_\ell \right|, \quad \forall \ell \in \mathcal{L}, \quad (23)$$

where $\tilde{\mathbf{K}} = \text{diag}(\{P_k \mathbf{I}\}_{k=1}^K)$ represents the transmit signal covariance matrix with all the users emitting independent signals across the antennas at their maximum powers. The rationale of this approach is to let signals from the BSs with either larger fronthaul capacity or lower received signal power be recovered first, then the recovered signals can serve as side information in helping the decompression of signals from other BSs. This decompression order attempts to make the quantization noise levels across the BSs small. It is shown by simulation in the later section that the proposed heuristic approach works rather well for implementing Wyner-Ziv coding in practical uplink C-RAN when the fronthaul capacities or the received signal powers at the BSs are different.

Assume that π is the decompression order of $\hat{\mathbf{y}}_\ell$ given by the heuristic approach. Denote the index set by $\mathcal{T}_\ell = \{\pi(1), \dots, \pi(\ell)\}$, where $\pi(\ell)$ represents the ℓ th component in π . Let $\mathbf{Q}_{\mathcal{T}_\ell} = \text{diag}(\{\mathbf{Q}_\ell\}_{\ell \in \mathcal{T}_\ell})$. The weighted sum rate maximization problem under Wyner-Ziv coding can be formulated as follows:

$$\begin{aligned} \max_{\mathbf{V}_k, \mathbf{Q}_\ell} \quad & \sum_{k=1}^K \alpha_k \log \left| \mathbf{I} + \mathbf{V}_k^\dagger \mathbf{H}_{\mathcal{L},k}^\dagger \mathbf{J}_k^{-1} \mathbf{H}_{\mathcal{L},k} \mathbf{V}_k \right| \\ \text{s.t.} \quad & \log \frac{|\mathbf{\Upsilon}_{\mathcal{T}_\ell} + \mathbf{Q}_{\mathcal{T}_\ell}|}{|\mathbf{\Upsilon}_{\mathcal{T}_{\ell-1}} + \mathbf{Q}_{\mathcal{T}_{\ell-1}}|} - \log |\mathbf{Q}_{\pi(\ell)}| \leq C_{\pi(\ell)}, \quad \forall \ell \in \mathcal{L}, \\ & \mathbf{Q}_\ell \succeq \mathbf{0}, \quad \forall \ell \in \mathcal{L}, \\ & \text{Tr}(\mathbf{V}_k \mathbf{V}_k^\dagger) \leq P_k, \quad \forall k \in \mathcal{K}, \end{aligned} \quad (24)$$

where $\mathbf{\Upsilon}_{\mathcal{T}_\ell} = \sum_{k=1}^K \mathbf{H}_{\mathcal{T}_\ell,k} \mathbf{V}_k \mathbf{V}_k^\dagger \mathbf{H}_{\mathcal{T}_\ell,k}^\dagger + \text{diag}(\{\mathbf{\Lambda}_\ell\}_{\ell \in \mathcal{T}_\ell})$, α_k 's are the weights associated with the users, and \mathbf{J}_k is given by either equation (5) for the linear MMSE receiver or equation (6) for the SIC receiver.

Algorithm 2 Beamforming and Fronthaul Compression Optimization under Wyner-Ziv coding

- 1: Determine a decomposition order π of $\hat{\mathbf{y}}_\ell$'s according to $C_\ell - \log |\mathbf{H}_{\ell,\mathcal{K}} \tilde{\mathbf{K}} \mathbf{H}_{\ell,\mathcal{K}}^\dagger + \mathbf{\Lambda}_\ell|$.
 - 2: Solve the optimization problem (26) using Algorithm 1. Set $(\mathbf{V}_k, \mathbf{Q}_\ell)$ to be its optimal solution.
-

The above problem is again non-convex, which makes finding its global optimum challenging. To efficiently solve problem (24), we again utilize the successive convex approximation approach proposed in the WMMSE-SCA algorithm. An obstacle to applying the convex approximation procedure directly to problem (24) lies in the Wyner-Ziv fronthaul constraint, which contains three log-determinant functions. To facilitate the utilization of the WMMSE-SCA algorithm, we reformulate problem (24) as an equivalent problem as follows,

$$\begin{aligned} \max_{\mathbf{V}_k, \mathbf{Q}_\ell} \quad & \sum_{k=1}^K \alpha_k \log |\mathbf{I} + \mathbf{V}_k^\dagger \mathbf{H}_{\mathcal{L},k}^\dagger \mathbf{J}_k^{-1} \mathbf{H}_{\mathcal{L},k} \mathbf{V}_k| \\ \text{s.t.} \quad & \log |\mathbf{Y}_{\mathcal{T}_\ell} + \mathbf{Q}_{\mathcal{T}_\ell}| - \sum_{\ell \in \mathcal{T}_\ell} \log |\mathbf{Q}_\ell| \leq \sum_{\ell \in \mathcal{T}_\ell} C_\ell, \quad \forall \ell \in \mathcal{L}, \\ & \mathbf{Q}_\ell \succeq \mathbf{0}, \quad \forall \ell \in \mathcal{L}, \\ & \text{Tr}(\mathbf{V}_k \mathbf{V}_k^\dagger) \leq P_k, \quad \forall k \in \mathcal{K}, \end{aligned} \quad (25)$$

The advantage of reformulation (25) is that it has similar format as (11), so the successive convex approximation procedure can again be used directly. Similar to the single-user case, by approximating the objective function and the fronthaul constraints in (25) with (12) and (18) respectively, problem (25) can be rewritten as

$$\begin{aligned} \max_{\substack{\mathbf{V}_k, \mathbf{Q}_\ell, \mathbf{U}_k, \\ \mathbf{W}_k, \mathbf{\Theta}_\ell, \mathbf{\Sigma}_{\mathcal{T}_\ell}}} \quad & \sum_{k=1}^K \alpha_k (\log |\mathbf{W}_k| - \text{Tr}(\mathbf{W}_k \mathbf{E}_k)) \\ & - \rho \sum_{\ell=1}^L \|\mathbf{Q}_\ell - \mathbf{\Theta}_\ell\|_F^2 \\ \text{s.t.} \quad & \log |\mathbf{\Sigma}_{\mathcal{T}_\ell}| + \text{Tr}(\mathbf{\Sigma}_{\mathcal{T}_\ell}^{-1} \mathbf{\Omega}_{\mathcal{T}_\ell}) - \log |\mathbf{Q}_{\mathcal{T}_\ell}| \leq C'_{\mathcal{T}_\ell}, \\ & \mathbf{Q}_\ell \succeq \mathbf{0}, \quad \mathbf{\Theta}_\ell \succeq \mathbf{0}, \quad \forall \ell \in \mathcal{L}, \\ & \text{Tr}(\mathbf{V}_k \mathbf{V}_k^\dagger) \leq P_k, \quad \forall k \in \mathcal{K}, \end{aligned} \quad (26)$$

where $\rho > 0$ is a constant, $C'_{\mathcal{T}_\ell} = \sum_{\ell \in \mathcal{T}_\ell} (C_\ell + N)$, and $\mathbf{\Omega}_{\mathcal{T}_\ell} = \sum_{k=1}^K \mathbf{H}_{\mathcal{T}_\ell,k} \mathbf{V}_k \mathbf{V}_k^\dagger \mathbf{H}_{\mathcal{T}_\ell,k}^\dagger + \text{diag}(\{\mathbf{\Lambda}_\ell + \mathbf{Q}_\ell\}_{\ell \in \mathcal{T}_\ell})$. Clearly, the proposed WMMSE-SCA algorithm can be applied for solving the above optimization problem. We summarize the beamforming and fronthaul compression scheme for Wyner-Ziv coding as Algorithm 2.

V. SEPARATE DESIGN OF BEAMFORMING AND COMPRESSION

Although locally optimal transmit beamformers and quantization noise covariance matrices can be found using the WMMSE-SCA algorithm for any fixed user schedule, user priority, and channel condition, the implementation of WMMSE-SCA in practice can be computationally intensive, especially when the channels are under fast fading or when the scheduled

users in the time-frequency slots change frequently. In this section, we aim at deriving near optimal transmit beamformers and quantization noise covariance matrices in the high SQNR regime. The main result of this section is that a simple separate design which involves optimizing transmit beamformers for the Gaussian vector multiple-access channel at the user side and using quantizers with uniform quantization noise levels across the antennas at each BS is approximately optimal under appropriate conditions. This leads to an efficient transmit beamforming and fronthaul compression design for practical uplink C-RAN systems.

A. Quantization Noise Design Under High SQNR

The proposed approximation scheme is derived by considering the sum rate maximization problem assuming single-user compression and assuming that SIC is implemented at the central receiver. Denote the transmit signal covariance matrix for the j th user as $\mathbf{K}_j = \mathbf{V}_j \mathbf{V}_j^\dagger$, and let $\mathbf{K}_\mathcal{K} = \text{diag}(\{\mathbf{K}_j\}_{j=1}^K)$. The sum rate maximization problem can be formulated as follows,

$$\begin{aligned} \max_{\mathbf{K}_j, \mathbf{Q}_\ell} \quad & \log \frac{|\mathbf{H}_{\mathcal{L},\mathcal{K}} \mathbf{K}_\mathcal{K} \mathbf{H}_{\mathcal{L},\mathcal{K}}^\dagger + \mathbf{\Lambda} + \mathbf{Q}|}{|\mathbf{\Lambda} + \mathbf{Q}|} \\ \text{s.t.} \quad & \log \frac{|\sum_{j=1}^K \mathbf{H}_{\ell,j} \mathbf{K}_j \mathbf{H}_{\ell,j}^\dagger + \mathbf{\Lambda}_\ell + \mathbf{Q}_\ell|}{|\mathbf{Q}_\ell|} \leq C_\ell, \quad \forall \ell \in \mathcal{L}, \\ & \mathbf{Q}_\ell \succeq \mathbf{0}, \quad \forall \ell \in \mathcal{L}, \\ & \text{Tr}(\mathbf{K}_j) \leq P_j, \quad \forall j \in \mathcal{K}, \end{aligned} \quad (27)$$

where $\mathbf{\Lambda} = \text{diag}(\{\mathbf{\Lambda}_\ell\}_{\ell=1}^L)$ and $\mathbf{Q} = \text{diag}(\{\mathbf{Q}_\ell\}_{\ell=1}^L)$.

In the following, we provide a justification that the optimal quantization noise levels should be set to be uniform across the antennas at each BS for maximizing the sum rate under high SQNR. By high SQNR, we require at least that the received signals across all the BS antennas occupy the entire space of receive dimensions, so implicitly enough number of users need to be scheduled, e.g., when $Kd = LN$. Further, the received signal strength needs to be much larger than the combined quantization and background noise level.

Mathematically, the required condition can be obtained by examining the Karush-Kuhn-Tucker (KKT) condition for the optimization problem (27). Form the Lagrangian

$$\begin{aligned} L(\mathbf{K}_j, \mathbf{Q}_\ell, \lambda_\ell, \mu_j) = & \log |\mathbf{H}_{\mathcal{L},\mathcal{K}} \mathbf{K}_\mathcal{K} \mathbf{H}_{\mathcal{L},\mathcal{K}}^\dagger + \mathbf{\Lambda} + \mathbf{Q}| \\ & - \log |\mathbf{\Lambda} + \mathbf{Q}| - \sum_{\ell=1}^L \lambda_\ell \log \left| \sum_{j=1}^K \mathbf{H}_{\ell,j} \mathbf{K}_j \mathbf{H}_{\ell,j}^\dagger + \mathbf{\Lambda}_\ell + \mathbf{Q}_\ell \right| \\ & + \sum_{\ell=1}^L \lambda_\ell \log |\mathbf{Q}_\ell| - \sum_{j=1}^K \mu_j \text{Tr}(\mathbf{K}_j), \end{aligned} \quad (28)$$

where λ_ℓ is the Lagrangian dual variable associated with the ℓ th fronthaul constraint, and μ_j is Lagrangian multiplier for the j th transmit power constraint.

Setting $\partial L / \partial \mathbf{Q}_\ell$ to zero, we obtain the optimality condition as follows,

$$\mathbf{F}_\ell \left(\mathbf{H}_{\mathcal{L},\mathcal{K}} \mathbf{K}_\mathcal{K} \mathbf{H}_{\mathcal{L},\mathcal{K}}^\dagger + \mathbf{\Lambda} + \mathbf{Q} \right)^{-1} \mathbf{F}_\ell^T - (\mathbf{\Lambda}_\ell + \mathbf{Q}_\ell)^{-1} - \lambda_\ell \left(\sum_{j=1}^K \mathbf{H}_{\ell,j} \mathbf{K}_j \mathbf{H}_{\ell,j}^\dagger + \mathbf{\Lambda}_\ell + \mathbf{Q}_\ell \right)^{-1} + \lambda_\ell \mathbf{Q}_\ell^{-1} = \mathbf{0}, \quad (29)$$

where $\mathbf{F}_\ell = [\mathbf{0}, \dots, \mathbf{0}, \mathbf{I}_N, \mathbf{0}, \dots, \mathbf{0}]$ with only the ℓ th $N \times N$ block being nonzero. Assuming that $Kd = LN$, the received signal covariance matrix $\mathbf{H}_{\mathcal{L},\mathcal{K}} \mathbf{K}_\mathcal{K} \mathbf{H}_{\mathcal{L},\mathcal{K}}^\dagger$ is full rank. Furthermore, if the overall system is to operate at reasonably high spectral efficiency, the received signal-to-noise ratios (SNRs) are likely to be high and the fronthaul capacities are likely to be large. In this case, we must have $\mathbf{H}_{\mathcal{L},\mathcal{K}} \mathbf{K}_\mathcal{K} \mathbf{H}_{\mathcal{L},\mathcal{K}}^\dagger + \mathbf{\Lambda} + \mathbf{Q} \gg \mathbf{\Lambda} + \mathbf{Q}$ and $\sum_{j=1}^K \mathbf{H}_{\ell,j} \mathbf{K}_j \mathbf{H}_{\ell,j}^\dagger + \mathbf{\Lambda}_\ell + \mathbf{Q}_\ell \gg \mathbf{\Lambda}_\ell + \mathbf{Q}_\ell$. Under this high SQNR condition, we argue that $\mathbf{F}_\ell \left(\mathbf{H}_{\mathcal{L},\mathcal{K}} \mathbf{K}_\mathcal{K} \mathbf{H}_{\mathcal{L},\mathcal{K}}^\dagger + \mathbf{\Lambda} + \mathbf{Q} \right)^{-1} \mathbf{F}_\ell^T \ll (\mathbf{\Lambda}_\ell + \mathbf{Q}_\ell)^{-1}$ and $\left(\sum_{j=1}^K \mathbf{H}_{\ell,j} \mathbf{K}_j \mathbf{H}_{\ell,j}^\dagger + \mathbf{\Lambda}_\ell + \mathbf{Q}_\ell \right)^{-1} \ll \mathbf{Q}_\ell^{-1}$, so that the optimality condition becomes

$$\mathbf{Q}_\ell \approx \frac{\lambda_\ell}{1 - \lambda_\ell} \mathbf{\Lambda}_\ell \quad (30)$$

where $\lambda_\ell \in [0, 1)$ is chosen to satisfy the fronthaul capacity constraints for single-user compression. Following the same analysis, similar conclusion can also be obtained for the sum rate maximization under Wyner-Ziv coding.

The above result implies that per-antenna quantizers with uniform quantization noise levels across the antennas at each BS are nearly optimal at high SQNR, although the quantization noise level may differ from BS to BS depending on the background noise levels and the fronthaul constraints. Note that this line of reasoning is very similar to the corresponding condition for the SISO case derived in [14].

It is worth emphasizing that in order to satisfy the high SQNR condition, the number of user data streams scheduled in the system should be at least as large as the number of receive spatial dimensions, and all these data streams must transmit at high rate. In scenarios where the number of data streams is less (i.e., some spatial dimensions are used for diversity instead of multiplexing), receive beamforming at the BSs prior to quantization may be beneficial. For example, MMSE beamforming or maximum ratio-combining may be applied at each BS in order to reduce the number of received dimensions before quantization.

B. Beamforming Design Under High SQNR

We next consider the optimal transmit beamforming and power allocation under high SQNR. Intuitively speaking, for maximizing the sum rate, each user should align its signaling direction with the strongest eigenmode of the effective channel and allocate power along this direction in a “water-filling” fashion. For this, we need to whiten the combined quantization and background noise and interference, then diagonalize the resulting channel to find its eigenmodes, and iteratively

perform the water-filling process among the users [34]. As seen from (30), at high SQNR, the optimal quantization noise covariance matrices are proportional to the background noise covariance matrices. Further, if $d = NL/K$, i.e., if the total number of user data streams is equal to the number of degrees of freedom in the system, then multiuser interference would be reasonably contained.

Based on the above intuition, we propose a simple beamforming design in which each user selects its transmit beamformers by ignoring the affect of fronthaul capacity limitation. Specifically, we consider the following weighted sum rate maximization problem for a Gaussian vector multiple-access channel:

$$\begin{aligned} \max_{\mathbf{K}_j} \quad & \sum_{k=1}^K \alpha_k \log \frac{\left| \sum_{j=k}^K \mathbf{H}_{\mathcal{L},j} \mathbf{K}_j \mathbf{H}_{\mathcal{L},j}^\dagger + \mathbf{\Lambda} \right|}{\left| \sum_{j>k}^K \mathbf{H}_{\mathcal{L},j} \mathbf{K}_j \mathbf{H}_{\mathcal{L},j}^\dagger + \mathbf{\Lambda} \right|} \\ \text{s.t.} \quad & \text{Tr}(\mathbf{K}_j) \leq P_j, \quad \forall j \in \mathcal{K}, \\ & \mathbf{K}_j \succeq \mathbf{0}, \quad \forall j \in \mathcal{K}, \\ & \text{rank}(\mathbf{K}_j) = d, \quad \forall j \in \mathcal{K}. \end{aligned} \quad (31)$$

If the rank constraints on transmit signal covariance matrices \mathbf{K}_j are ignored, the above problem becomes a convex optimization, and its optimum solution \mathbf{K}_j^* can be efficiently found through the interior-point method [35]. With the rank constraints applied, the problem is no longer convex, and we propose to find a set of approximately optimal transmit beamformers for user j as follows. First, solve (31) with the rank constraints removed; let the optimal solution be \mathbf{K}_j^* . Let γ_j^i represent the i th largest eigenvalue of \mathbf{K}_j^* and Ψ_j^i represent its normalized eigenvector. Then an approximately optimal transmit beamforming matrix \mathbf{V}_j^* for user j is given by

$$\mathbf{V}_j^* = \left[\sqrt{\frac{P_j \gamma_j^1}{\Gamma_j}} \Psi_j^1, \dots, \sqrt{\frac{P_j \gamma_j^d}{\Gamma_j}} \Psi_j^d \right] \quad (32)$$

where $\Gamma_j = \sum_{i=1}^d \gamma_j^i$ represents the sum of d largest of eigenvalues \mathbf{K}_j^* .

When linear MMSE receiver is employed, simply ignoring the rank constraints in the weighted sum rate maximization problem does not make it a convex optimization. In this case, one can rewrite $\mathbf{K}_j = \mathbf{V}_j \mathbf{V}_j^\dagger$ and use the WMMSE method [18], [19] to find the optimal beamforming vector \mathbf{V}_j^* .

C. Separate Beamforming and Compression Design

The above beamforming strategy together with per-antenna scalar quantizer provide us a low-complexity separate design for transmit beamforming and fronthaul compression. With single-user compression, define

$$C^{\text{SU}}(\beta_\ell) = \log \frac{\left| \sum_{k=1}^K \mathbf{H}_{\ell,k} \mathbf{V}_k \mathbf{V}_k^\dagger \mathbf{H}_{\ell,k}^\dagger + (1 + \beta_\ell) \mathbf{\Lambda}_\ell \right|}{|\beta_\ell \mathbf{\Lambda}_\ell|}. \quad (33)$$

To fully utilize the fronthaul capacities, the bisection search is employed to find the optimal β_ℓ such that $C^{\text{SU}}(\beta_\ell) = C_\ell$ for $\ell = 1, \dots, L$.

Algorithm 3 Separate Beamforming and Fronthaul Compression Design

- 1: Solve problem (31) with the rank constraints removed and set \mathbf{K}_j^* to be its optimal solution.
- 2: Perform eigenvalue decomposition on \mathbf{K}_j^* to obtain its normalized eigenvalues γ_j^i and eigenvectors Ψ_j^i . Set $\mathbf{V}_j^* = \left[\sqrt{\frac{P_j \gamma_j^1}{\Gamma_j}} \Psi_j^1, \dots, \sqrt{\frac{P_j \gamma_j^d}{\Gamma_j}} \Psi_j^d \right]$.
- 3: Under single-user compression, use bisection method in $[\beta_\ell^{\min}, \beta_\ell^{\max}]$ to solve for β_ℓ in $C_\ell^{\text{SU}}(\beta_\ell) = C_\ell$ independently for $\ell = 1, \dots, L$; Under Wyner-Ziv coding, use bisection in $[\beta_j^{\min}, \beta_j^{\max}]$ to solve for β_j in $C_j^{\text{WZ}}(\beta_1, \dots, \beta_{j-1}, \beta_j) = \sum_{\ell=1}^j C_\ell$ for $j = 1, \dots, L$ successively with the values of $\beta_1, \dots, \beta_{j-1}$ fixed and as determined by the previous $j-1$ bisection searches.
- 4: Set $\mathbf{Q}_\ell = \beta_\ell \mathbf{A}_\ell$ for $\ell = 1, \dots, L$.

With Wyner-Ziv coding, assuming without loss of generality a decoding order of \hat{y}_ℓ from 1 to L , define

$$C_j^{\text{WZ}}(\beta_1, \dots, \beta_j) = \frac{\log \left| \sum_{k=1}^K \mathbf{H}_{\mathcal{T}_j, k} \mathbf{V}_k \mathbf{V}_k^\dagger \mathbf{H}_{\mathcal{T}_j, k}^\dagger + \text{diag}(\{(1 + \beta_\ell) \mathbf{A}_\ell\}_{\ell \in \mathcal{T}_j}) \right|}{\log |\text{diag}(\{\beta_\ell \mathbf{A}_\ell\}_{\ell \in \mathcal{T}_j})|} \quad (34)$$

where $\mathcal{T}_j = \{1, \dots, j\}$. Different from the single-user case, the optimal β_ℓ in Wyner-Ziv coding is determined one after another in a successive fashion. Specifically, for $j = 1, \dots, L$, we use bisection search to find the optimal β_j such that $C_j^{\text{WZ}}(\beta_1, \dots, \beta_{j-1}, \beta_j) = \sum_{\ell=1}^j C_\ell$ assuming that the values of $\beta_1, \dots, \beta_{j-1}$ are fixed and as determined by the previous $j-1$ bisection searches.

The separate transmit beamforming and fronthaul compression design scheme is summarized as Algorithm 3. In the Step 3 of Algorithm 3, the values of β_ℓ^{\min} and β_ℓ^{\max} can be found using the same way in [14]. Specifically, under single-user compression, initialize $\beta_\ell = 1$ for $\ell = 1, \dots, L$, and keep updating $\beta_\ell = 2\beta_\ell$ until $C_\ell^{\text{SU}}(\beta_\ell) \leq C_\ell$ is satisfied. Then, we set $\beta_\ell^{\min} = 0$ and $\beta_\ell^{\max} = \beta_\ell$ for $\ell = 1, \dots, L$. Similar procedure can also be used in the case of Wyner-Ziv coding.

There are two differences between the joint design scheme and the separate design scheme. First, in the joint design, transmit beamforming are chosen to be fronthaul-aware, while the impact of limit fronthaul is ignored in the separate design. Second, in the joint design, the quantization is performed on the received signal vector across all the receive antennas at each BS while separate design adopts per-antenna quantization on each receive antenna of the BSs. It is shown by simulation in the next section that the separate design performs very well in the high SQNR regime. In other regimes, the difference between the joint design and separate design represents a tradeoff between complexity and performance in implementing uplink C-RAN.

We remark that when the high SQNR condition is not satisfied, the optimal beamforming in uplink C-RAN should

TABLE I
MULTICELL NETWORK SYSTEM PARAMETERS

Cellular Layout	Hexagonal, 19-cell, 3 sectors/cell
BS-to-BS Distance	500 m
Frequency Reuse	1
Channel Bandwidth	10 MHz
Number of Users per Sector	20
Total Number of Users	420
Max Transmit Power	23 dBm
Antenna Gain	14 dBi
Background Noise	-169 dBm/Hz
Noise Figure	7 dB
Tx/Rx Antenna No.	2×2
Distance-dependent Path Loss	$128.1 + 37.6 \log_{10}(d)$
Log-normal Shadowing	8 dB standard deviation
Shadow Fading Correlation	0.5
Cluster Size	7 cells (21 sectors)
Scheduling Strategy	WMMSE based scheduling

be fronthaul-aware. For example, consider a two-layer heterogeneous C-RAN system with both pico BSs and macro BSs serving as relay nodes. The fronthaul capacity of the macro BS is typically much larger than that of the pico BS. Therefore, users are more likely to form their transmit beamformer pointing toward the receive antennas at the macro BSs rather than the pico BSs. Under this scenario, both of the channel strength between the users and the BSs and the fronthaul capacities between the BSs and the CP should be taken into account in the beamforming design in order to maximize the network throughput.

From the computational complexity point of view, the separate design is significantly superior to the joint design. Algorithm 3 involves solving a single convex optimization problem (31) plus a bisection search, as compared to iteratively solving a series of convex optimization problems (22) or (26) as in the WMMSE-SCA algorithm.

VI. SIMULATION RESULTS

A. Single-Cluster Network

In this section, the performances of the proposed WMMSE-SCA schemes with different compression strategies (i.e., Wyner-Ziv coding and single-user compression) and different receiving schemes (i.e., linear MMSE receiver and SIC receiver) are evaluated on a 19-cell 3-sector/cell wireless network setup with center 7 cells (i.e., 21 sectors) forming a cooperating cluster. The users are randomly located and associated with the strongest BS. The proposed WMMSE-SCA algorithm is applied to all the users within the cluster, which automatically schedules the users with non-zero beamforming vectors. Each BS is equipped with $N = 2$ antennas, each user is equipped with $M = 2$ antennas, and each user sends one data stream (i.e., $d = 1$) to the CP. Perfect channel estimation is assumed, and the CSI is made available to all BSs and to the CP. Various algorithms are run on fixed set of channels. Detailed system parameters are outlined in Table I.

Under single-user compression, Fig. 2 and Fig. 3 compare the performance of the WMMSE-SCA and separate design schemes implemented either with SIC (labeled as “SIC receiver” in the figures) or without SIC (labeled as “linear

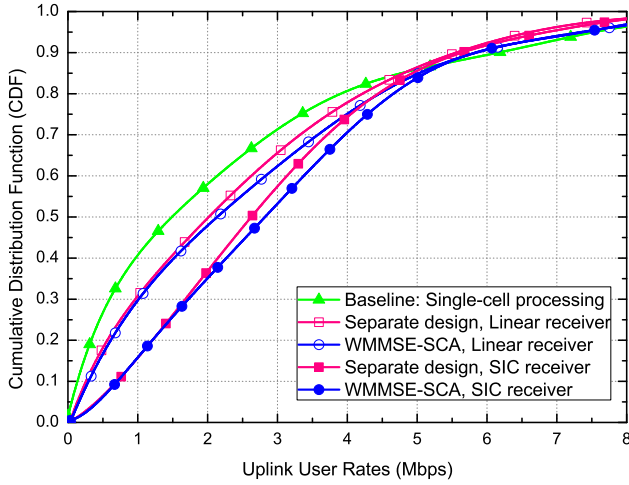


Fig. 2. Cumulative distribution of user rates with single-user compression for a 19-cell network with center 7 cells forming a single cluster under the fronthaul capacity of 120Mbps per sector.

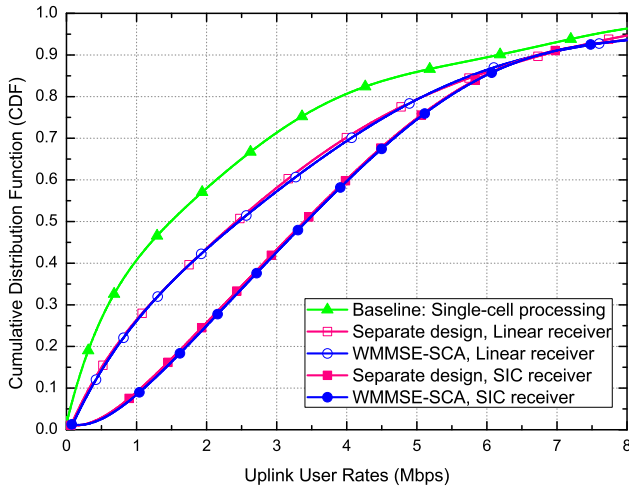


Fig. 3. Cumulative distribution of user rates with single-user compression for a 19-cell network with center 7 cells forming a single cluster under the fronthaul capacity of 320Mbps per sector.

receiver” in the figures) at the receiver under two different fronthaul constraints. It is shown that both the WMMSE-SCA scheme and the separate design scheme significantly outperform the baseline scheme without multicell processing. Fig. 2 and Fig. 3 show that the SIC receiver achieves significant gain as compared to the linear receiver. The performance improvement is more significant for the users with low rate.

To further compare the performance of the proposed two schemes, Fig. 4 plots the average per-cell sum rate of the WMMSE-SCA scheme and the low-complexity separate design as a function of the fronthaul capacity. As the fronthaul capacity increases, the performance gap between these two schemes becomes smaller. This demonstrates the approximate optimality for separate design of transmit beamforming and fronthaul compression in the high SQNR regime.

Fig. 5 and Fig. 6 show the CDF curves of user rates for the WMMSE-SCA scheme implemented with four different choices of coding schemes: with either single-user or Wyner-

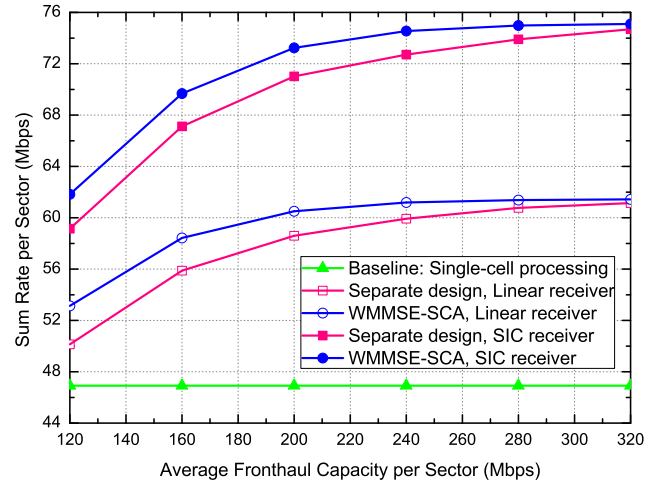


Fig. 4. Per-cell sum rate vs. average per-sector fronthaul capacity for single-user compression with linear receiver and with SIC receiver for a 19-cell network with center 7 cells forming a single cluster.

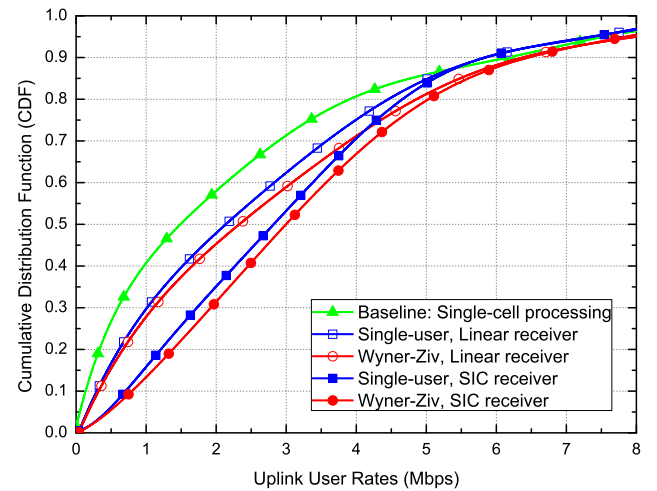


Fig. 5. Cumulative distribution of user rates with either single-user compression or Wyner-Ziv coding for a 19-cell network with center 7 cells forming a single cluster under the fronthaul capacity of 120Mbps per sector.

Ziv compression at the BSs and with either linear MMSE or SIC receiver at the CP. It can be seen from Fig. 5 that under the fronthaul capacity of 120Mbps, single-user compression with SIC receiver significantly improves the performance of linear MMSE receiver. Further gain on performance can be obtained if one replaces single-user compression by Wyner-Ziv coding. As the capacity of fronthaul increases to 320Mbps, as shown in Fig. 6, the gain due to Wyner-Ziv coding becomes negligible. In this high fronthaul scenario, SIC receiver still achieves a very large gain.

In order to quantify the performance gain brought by Wyner-Ziv coding and SIC receiver, Fig. 7 shows the average per-cell sum rate obtained by different schemes as the average capacity of fronthaul increases. It is observed that, under fronthaul capacity of 320Mbps, SIC receiver and Wyner-Ziv coding outperform the linear receiver and single-user compression respectively. But the performance improvement of SIC receiver upon linear receiver is much larger than the gain of Wyner-Ziv

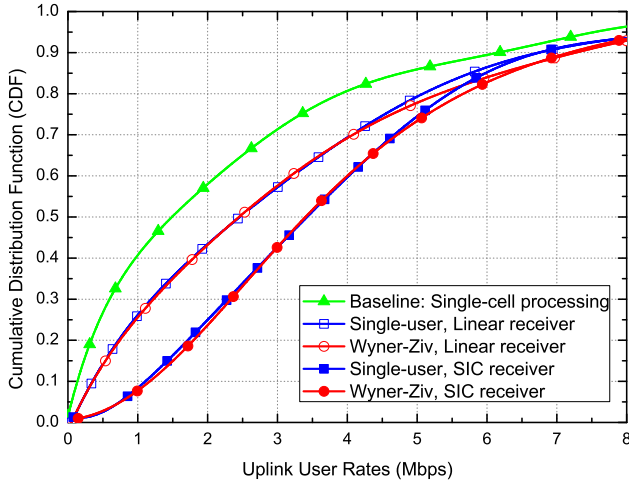


Fig. 6. Cumulative distribution of user rates with either single-user compression or Wyner-Ziv coding using WMMSE-SCA algorithm for a 19-cell network with center 7 cells forming a single cluster under the fronthaul capacity of 320Mbps per sector.

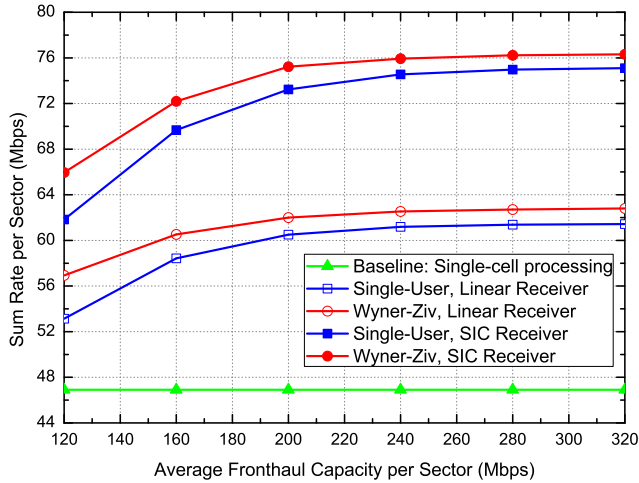


Fig. 7. Per-cell sum rate vs. average per-cell fronthaul capacity with either single-user compression or Wyner-Ziv coding using WMMSE-SCA algorithm for a 19-cell network with center 7 cells forming a single cluster.

coding over single-user compression.

B. Multi-Cluster Network

The performance of the proposed WMMSE-SCA scheme is further evaluated for a large-scale multicell network with 65 cells and 10 mobile users randomly located within each cell. The BS-to-BS distance is set to be 200m, each user is equipped with 2 transmit antennas, and each BS is equipped with 4 receive antennas. The channel is assumed to be flat-fading. Round-robin user scheduling is used on a per-cell basis and system is operated with loading factor 0.5, i.e., in each time slot, BS schedules two users. Detailed system parameters are outlined in Table II. Two different clustering strategies, i.e., disjoint clustering [4], [5] and user-centric clustering [6], [7], are applied to form clusters within the network. Disjoint clustering partitions the BSs in the network into nonoverlapping sets of cooperating clusters. In user-

TABLE II
MULTI-CLUSTER NETWORK PARAMETERS

Cellular Layout	Hexagonal
BS-to-BS Distance	200 m
Frequency Reuse	1
Channel Bandwidth	10 MHz
Number of Users per Cell	10
Number of Cells	65
Total Number of Users	650
Max Transmit Power	23 dBm
Antenna Gain	14 dBi
Background Noise	-169 dBm/Hz
Noise Figure	7 dB
Tx Antenna No.	2
Rx Antenna No.	4
Distance-dependent Path Loss	$128.1 + 37.6 \log_{10}(d)$
Log-normal Shadowing	8 dB standard deviation
Shadow Fading Correlation	0.5
Scheduling Strategy	Round-robin

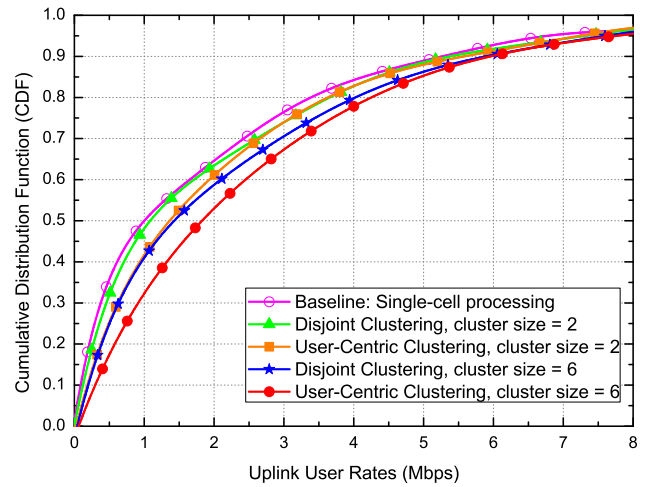


Fig. 8. Cumulative distribution of user rates for the WMMSE-SCA algorithm with single-user compression under the average fronthaul capacity of 120Mbps with either disjoint or user-centric clustering for a multi-cluster network.

centric clustering, each user chooses a set of nearest BSs to form a cooperation cluster, and cooperating clusters overlap, which makes the implementation of Wyner-Ziv coding and SIC receiver under fronthaul capacity constraints (10) more difficult. Therefore, for fair comparison, we only consider here the case where single-user compression and linear MMSE receiver are employed.

Fig. 8 and Fig. 9 show the CDF plots of user rates achieved with both disjoint clustering and user-centric clustering with WMMSE-SCA. It is clear that with optimized beamforming and fronthaul compression, the user-centric clustering significantly improves over disjoint clustering, and both of these two schemes improve as the cluster size increases. As the capacity of fronthaul links increases from 120Mbps to 360Mbps, the performance gap between the two clustering schemes becomes larger. Further, for disjoint clustering, increasing the cluster size from 2 to 6 achieves 60% performance improvement for the 50-percentile rate. This gain doubles when we further replace disjoint clustering with user-centric clustering.

Fig. 10 plots the average per-cell sum rate as the fronthaul

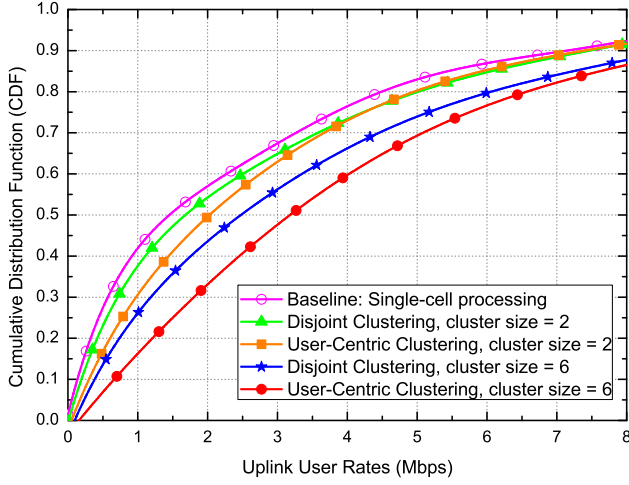


Fig. 9. Cumulative distribution of user rates for the WMMSE-SCA algorithm with single-user compression under the average fronthaul capacity of 360Mbps with either disjoint or user-centric clustering for a multi-cluster network.

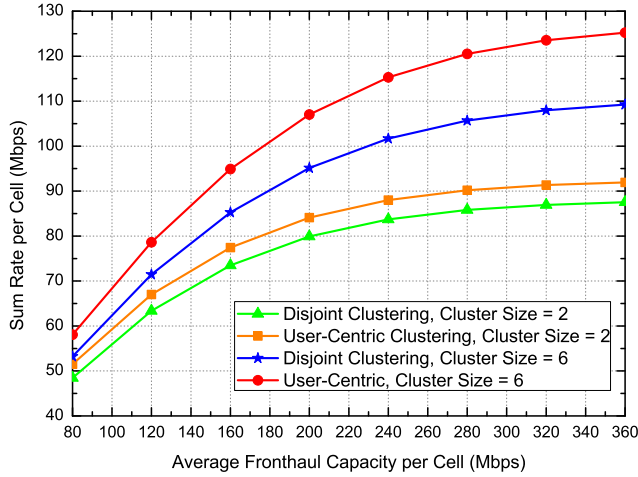


Fig. 10. Per-cell sum rate vs. average per-cell fronthaul capacity of the WMMSE-SCA algorithm with single-user compression for a multi-cluster network under different clustering strategies and different cluster size.

capacity increases. The result again shows that user-centric clustering achieves significant performance gain over disjoint clustering. When cluster size increases to 6, to achieve per-cell sum rate of 110Mbps, disjoint clustering needs fronthaul capacity of 360Mbps, while user-centric needs 220Mbps, which is more than 60% improvement on the fronthaul requirement.

Finally, the performance of the two different clustering strategies are compared as a function of cluster size in Fig. 11. It is shown that for both disjoint clustering and user-centric clustering, the average per-cell sum rate increases as either the cluster size or fronthaul capacity increases. As expected, user-centric clustering always outperforms disjoint clustering. If we compare the performance of disjoint clustering with fronthaul capacity of 360Mbps with user-centric clustering with fronthaul capacity of 240Mbps, we see that even with 120Mbps lower fronthaul capacity, user-centric clustering already achieves higher per-cell sum rate. This improvement on per-cell sum rate becomes larger as the cluster size increases.

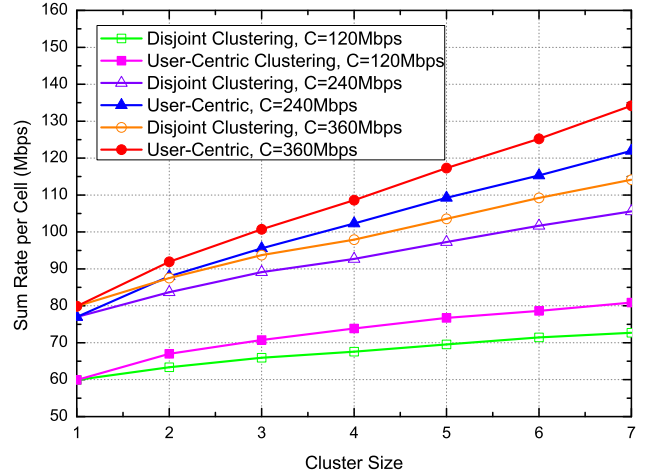


Fig. 11. Per-cell sum rate vs. cluster size for the WMMSE-SCA algorithm with single-user compression for a multi-cluster network under different clustering strategies and different fronthaul capacity constraints.

VII. CONCLUSION

This paper studies the fronthaul compression and transmit beamforming design for an uplink MIMO C-RAN system. From algorithm design perspective, we propose a novel WMMSE-SCA algorithm to efficiently optimize the transmit beamformer and quantization noise covariance matrix jointly for maximizing the weighted sum rate with either Wyner-Ziv coding or single-user compression. Further, we propose a separate design consisting of transmit beamforming optimized for the Gaussian vector multiple-access channel without accounting for compression together with scalar quantization with uniform quantization noise levels across the antennas at each BS. This low-complexity separate design is shown to be near optimal for maximizing the weighted sum rate when the SQNR is high. The performance of optimized beamforming and fronthaul compression is evaluated for practical multicell networks with different compression strategies, different receiving schemes, and different clustering methods. Numerical results show that, with optimized beamforming and fronthaul compression, C-RAN can significantly improve the overall performance of MIMO cellular networks. Most of the performance gain are due to the implementation of SIC at the central receiver. Finally, user-centric clustering significantly outperforms disjoint clustering in terms of fronthaul capacity saving.

APPENDIX A PROOF OF THEOREM 1

The proof of Theorem 1 is a direct application of the convergence result of the successive convex approximation algorithm [32]. Let $\mathbf{V} = \text{diag}(\{\mathbf{V}_k\}_{k=1}^K)$. Define the objective

function and fronthaul constraints in problem (11) to be

$$f(\mathbf{V}, \mathbf{Q}) = \sum_{k=1}^K \alpha_k \log \left| \mathbf{I} + \mathbf{V}_k^\dagger \mathbf{H}_{\mathcal{L},k}^\dagger \mathbf{J}_k^{-1} \mathbf{H}_{\mathcal{L},k} \mathbf{V}_k \right|,$$

$$g_\ell(\mathbf{V}, \mathbf{Q}) = \log \frac{\left| \sum_{k=1}^K \mathbf{H}_{\mathcal{L},k} \mathbf{V}_k \mathbf{V}_k^\dagger \mathbf{H}_{\mathcal{L},k}^\dagger + \mathbf{\Lambda}_\ell + \mathbf{Q}_\ell \right|}{|\mathbf{Q}_\ell|} - C_\ell,$$

where $\mathbf{J}_k = \sum_{j \neq k}^K \mathbf{H}_{\mathcal{L},j} \mathbf{V}_j \mathbf{V}_j^\dagger \mathbf{H}_{\mathcal{L},j}^\dagger + \mathbf{\Lambda} + \mathbf{Q}$ for the linear receiver or $\mathbf{J}_k = \sum_{j > k}^K \mathbf{H}_{\mathcal{L},j} \mathbf{V}_j \mathbf{V}_j^\dagger \mathbf{H}_{\mathcal{L},j}^\dagger + \mathbf{\Lambda} + \mathbf{Q}$ for the SIC receiver.

At the t th iteration, assume that the output of WMMSE-SCA algorithm is $(\mathbf{V}^t, \mathbf{Q}^t)$. Putting $(\mathbf{V}^t, \mathbf{Q}^t)$ into equations (15) and (19) gives

$$\Sigma_\ell^t = \sum_{k=1}^K \mathbf{H}_{\mathcal{L},k} \mathbf{V}_k^t (\mathbf{V}_k^t)^\dagger \mathbf{H}_{\mathcal{L},k}^\dagger + \mathbf{\Lambda}_\ell + \mathbf{Q}_\ell^t,$$

$$\mathbf{W}_k^t = \mathbf{I} + (\mathbf{V}_k^t)^\dagger \mathbf{H}_{\mathcal{L},k}^\dagger \mathbf{U}_k^t,$$

where

$$\mathbf{U}_k^t = \left(\sum_{j \neq k} \mathbf{H}_{\mathcal{L},j} \mathbf{V}_j^t (\mathbf{V}_j^t)^\dagger \mathbf{H}_{\mathcal{L},j}^\dagger + \mathbf{\Lambda} + \mathbf{Q}^t \right)^{-1} \mathbf{H}_{\mathcal{L},k} \mathbf{V}_k^t.$$

Then the objective function and fronthaul constraints in problem (21) can be written as

$$\tilde{f}(\{\mathbf{V}, \mathbf{Q}\}, \{\mathbf{V}^t, \mathbf{Q}^t\}) = \sum_{k=1}^K \alpha_k (\log |\mathbf{W}_k^t| - \text{Tr}(\mathbf{W}_k^t \mathbf{E}_k))$$

$$- \rho \sum_{\ell=1}^L \|\mathbf{Q}_\ell - \mathbf{Q}_\ell^t\|_F^2,$$

$$\tilde{g}_\ell(\{\mathbf{V}, \mathbf{Q}\}, \{\mathbf{V}^t, \mathbf{Q}^t\}) = \log |\Sigma_\ell^t| + \text{Tr}((\Sigma_\ell^t)^{-1} \mathbf{\Omega}_\ell)$$

$$- \log |\mathbf{Q}_\ell| - C_\ell - N,$$

where

$$\mathbf{E}_k = (\mathbf{I} - (\mathbf{U}_k^t)^\dagger \mathbf{H}_{\mathcal{L},k} \mathbf{V}_k) (\mathbf{I} - (\mathbf{U}_k^t)^\dagger \mathbf{H}_{\mathcal{L},k} \mathbf{V}_k)^\dagger$$

$$+ (\mathbf{U}_k^t)^\dagger \left(\sum_{j \neq k} \mathbf{H}_{\mathcal{L},j} \mathbf{V}_j \mathbf{V}_j^\dagger \mathbf{H}_{\mathcal{L},j}^\dagger + \mathbf{\Lambda} + \mathbf{Q} \right) \mathbf{U}_k^t,$$

and $\mathbf{\Omega}_\ell = \sum_{k=1}^K \mathbf{H}_{\mathcal{L},k} \mathbf{V}_k \mathbf{V}_k^\dagger \mathbf{H}_{\mathcal{L},k}^\dagger + \mathbf{\Lambda}_\ell + \mathbf{Q}_\ell$.

We now observe that the WMMSE-SCA algorithm is actually a special case of the general successive convex approximation (SCA) method, with \tilde{f} and \tilde{g}_ℓ being the convex approximation functions of f and g_ℓ respectively. Furthermore, based on the fact that \tilde{f} is strictly convex over (\mathbf{V}, \mathbf{Q}) and the result of [36, Lemma 3.1], it can be shown that \tilde{f} is uniformly strongly convex over (\mathbf{V}, \mathbf{Q}) . We point out here that the regularization term, $\rho \sum_{\ell=1}^L \|\mathbf{Q}_\ell - \mathbf{Q}_\ell^t\|_F^2$, plays a key role in making \tilde{f} strongly convex.

Define

$$\mathcal{X} \triangleq \left\{ (\mathbf{V}, \mathbf{Q}) \left| \begin{array}{l} \mathbf{Q}_\ell \succeq \mathbf{0}, \quad \forall \ell \in \mathcal{L} \\ \text{Tr}(\mathbf{V}_k \mathbf{V}_k^\dagger) \leq P_k, \quad \forall k \in \mathcal{K} \end{array} \right. \right\} \quad (35)$$

and

$$\mathcal{Y} \triangleq \left\{ (\mathbf{V}, \mathbf{Q}) \left| \begin{array}{l} g_\ell(\mathbf{V}, \mathbf{Q}) \leq 0, \quad \forall \ell \in \mathcal{L} \\ (\mathbf{V}, \mathbf{Q}) \in \mathcal{X} \end{array} \right. \right\} \quad (36)$$

We summarize the conditions that are satisfied for the functions f , g_ℓ , \tilde{f} and \tilde{g}_ℓ as follows:

- 1) \mathcal{X} is closed and convex (and nonempty);
- 2) f and g_ℓ are continuous and differentiable on \mathcal{X} , and ∇f is Lipschitz continuous on \mathcal{X} ;
- 3) $\tilde{f}(\cdot, \mathbf{y})$ is uniformly strongly convex on \mathcal{X} for all $\mathbf{y} \in \mathcal{Y}$ with some positive constant;
- 4) $\tilde{f}(\cdot, \cdot)$ is continuous on $\mathcal{X} \times \mathcal{Y}$ and $\nabla_{\mathbf{y}} \tilde{f}(\mathbf{y}, \mathbf{y}) = \nabla_{\mathbf{y}} f(\mathbf{y})$, for all $\mathbf{y} \in \mathcal{Y}$;
- 5) $\tilde{g}_\ell(\cdot, \mathbf{y})$ is convex on \mathcal{X} for all $\mathbf{y} \in \mathcal{Y}$, and $\tilde{g}_\ell(\mathbf{y}, \mathbf{y}) = g_\ell(\mathbf{y})$, for all $\mathbf{y} \in \mathcal{Y}$;
- 6) $g_\ell(\mathbf{x}) \leq \tilde{g}_\ell(\mathbf{x}, \mathbf{y})$ for all $\mathbf{x} \in \mathcal{X}$ and $\mathbf{y} \in \mathcal{Y}$;
- 7) $\tilde{g}_\ell(\cdot, \cdot)$ is continuous on $\mathcal{X} \times \mathcal{Y}$ and $\nabla_{\mathbf{y}} \tilde{g}_\ell(\mathbf{y}, \mathbf{y}) = \nabla_{\mathbf{y}} g_\ell(\mathbf{y})$, for all $\mathbf{y} \in \mathcal{Y}$;
- 8) All feasible points of problem (11) are regular (see, e.g. [32]).

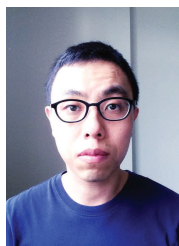
where $\nabla_{\mathbf{y}} \tilde{f}(\mathbf{y}, \mathbf{y})$ and $\nabla_{\mathbf{y}} \tilde{g}_\ell(\mathbf{y}, \mathbf{y})$ denote the (partial) gradients of \tilde{f} and \tilde{g}_ℓ respectively, which are with respect to the first argument evaluated at \mathbf{y} (the second argument is kept fixed at \mathbf{y}).

Based on the above conditions, it is shown in [32, Theorem 2] that the SCA algorithm converges to a stationary point of the nonconvex problem (11). Therefore, we conclude that each of the limit points generated by the proposed WMMSE-SCA algorithm is also a stationary point of problem (11), which completes the proof of Theorem 1.

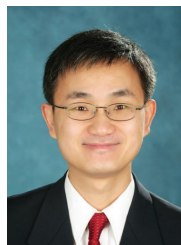
REFERENCES

- [1] O. Simeone, A. Maeder, M. Peng, O. Sahin, and W. Yu, "Cloud radio access network: Virtualizing wireless access for dense heterogeneous systems," *J. Commun. Netw.*, to appear 2016. [Online]. Available: <http://arxiv.org/pdf/1512.07743v1.pdf>
- [2] D. Gesbert, S. Hanly, H. Huang, S. Shamai, O. Simeone, and W. Yu, "Multi-cell MIMO cooperative networks: A new look at interference," *IEEE J. Sel. Areas Commun.*, vol. 28, no. 9, pp. 1380–1408, Dec. 2010.
- [3] S.-H. Park, O. Simeone, O. Sahin, and S. Shamai, "Fronthaul compression for cloud radio access networks: Signal processing advances inspired by network information theory," *IEEE Signal Process. Mag.*, vol. 31, no. 6, pp. 69–79, 2014.
- [4] 3GPP, "Coordinated multi-point operation for LTE physical layer aspects," 3rd Generation Partnership Project (3GPP), TR 36.819, Sep. 2011. [Online]. Available: <http://www.3gpp.org/Specs/36819-b10.pdf>
- [5] M. Peng, C. Wang, V. Lau, and H. V. Poor, "Fronthaul-constrained cloud radio access networks: Insights and challenges," *IEEE Wireless Commun. Mag.*, vol. 22, no. 2, pp. 152–160, Apr. 2015.
- [6] C. T. K. Ng and H. Huang, "Linear precoding in cooperative MIMO cellular networks with limited coordination clusters," *IEEE J. Sel. Areas Commun.*, vol. 28, no. 9, pp. 1446–1454, Sep. 2010.
- [7] B. Dai and W. Yu, "Sparse beamforming and user-centric clustering for downlink cloud radio access network," *IEEE Access*, vol. 2, pp. 1326–1339, Oct. 2014.
- [8] A. Sanderovich, S. Shamai, Y. Steinberg, and G. Kramer, "Communication via decentralized processing," *IEEE Trans. Inf. Theory*, vol. 54, no. 7, pp. 3008–3023, Jul. 2008.
- [9] L. Zhou and W. Yu, "Uplink multicell processing with limited backhaul via per-base-station successive interference cancellation," *IEEE J. Sel. Areas Commun.*, vol. 31, no. 10, pp. 1981–1993, Oct. 2013.
- [10] Y. Zhou and W. Yu, "Optimized beamforming and backhaul compression for uplink MIMO cloud radio access networks," in *Proc. IEEE GLOBECOM Workshops*, Dec. 2014, pp. 1487–1492.
- [11] A. Sanderovich, O. Somekh, H. V. Poor, and S. Shamai, "Uplink macro diversity of limited backhaul cellular network," *IEEE Trans. Inf. Theory*, vol. 55, no. 8, pp. 3457–3478, Aug. 2009.
- [12] S. H. Lim, Y.-H. Kim, A. El Gamal, and S.-Y. Chung, "Noisy network coding," *IEEE Trans. Inf. Theory*, vol. 57, no. 5, pp. 3132–3152, May 2011.

- [13] S.-H. Park, O. Simeone, O. Sahin, and S. Shamai, "Joint decompression and decoding for cloud radio access networks," *IEEE Signal Process. Lett.*, vol. 20, no. 5, pp. 503–506, May 2013.
- [14] Y. Zhou and W. Yu, "Optimized backhaul compression for uplink cloud radio access network," *IEEE J. Sel. Areas Commun.*, vol. 32, no. 6, pp. 1295–1307, Jun. 2014.
- [15] Y. Zhou, Y. Xu, J. Chen, and W. Yu, "Optimality of Gaussian fronthaul compression for uplink MIMO cloud radio access networks," in *Proc. IEEE Int. Symp. Inf. Theory (ISIT)*, Jun. 2015, pp. 2241–2245.
- [16] A. del Coso and S. Simoes, "Distributed compression for MIMO coordinated networks with a backhaul constraint," *IEEE Trans. Wireless Commun.*, vol. 8, no. 9, pp. 4698–4709, Sep. 2009.
- [17] S.-H. Park, O. Simeone, O. Sahin, and S. Shamai, "Robust and efficient distributed compression for cloud radio access networks," *IEEE Trans. Veh. Technol.*, vol. 62, no. 2, pp. 692–703, Feb. 2013.
- [18] S. S. Christensen, R. Agarwal, E. Carvalho, and J. M. Cioffi, "Weighted sum-rate maximization using weighted MMSE for MIMO-BC beamforming design," *IEEE Trans. Wireless Commun.*, vol. 7, no. 12, pp. 4792–4799, Jul. 2008.
- [19] 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 Process.*, vol. 59, no. 9, pp. 4331–4340, Sep. 2011.
- [20] D. R. Hunter and K. Lange, "A tutorial on MM algorithms," *Am. Stat.*, vol. 58, no. 1, pp. 30–37, Jan. 2004.
- [21] M. Razaviyayn, M. Hong, and Z.-Q. Luo, "A unified convergence analysis of block successive minimization methods for nonsmooth optimization," *SIAM J. Optim.*, vol. 23, no. 2, pp. 1126–1153, Feb. 2013.
- [22] X. Rao and V. K. Lau, "Distributed fronthaul compression and joint signal recovery in Cloud-RAN," *IEEE Trans. Signal Process.*, vol. 63, no. 4, pp. 1056–1065, Feb. 2015.
- [23] E. Heo, O. Simeone, and H. Park, "Optimal fronthaul compression for synchronization in the uplink of cloud radio access networks," Oct. 2015. [Online]. Available: <http://arxiv.org/abs/1510.01545>
- [24] W. Wang, V. K. Lau, and M. Peng, "Delay-optimal fronthaul allocation via perturbation analysis in cloud radio access networks," in *Proc. IEEE International Conf. Comm. (ICC)*, Jun. 2015, pp. 3999–4004.
- [25] M. Peng, Y. Li, J. Jiang, J. Li, and C. Wang, "Heterogeneous cloud radio access networks: a new perspective for enhancing spectral and energy efficiencies," *IEEE Trans. Wireless Commun.*, vol. 21, no. 6, pp. 126–135, 2014.
- [26] M. Peng, S. Yan, K. Zhang, and C. Wang, "Fog computing based radio access networks: Issues and challenges," *IEEE Netw.*, to appear 2016. [Online]. Available: <http://arxiv.org/abs/1506.04233>
- [27] Y. Zhou and W. Yu, "Approximate bounds for limited backhaul uplink multicell processing with single-user compression," in *Proc. IEEE Canadian Workshop Inf. Theory (CWIT)*, Jun. 2013, pp. 113–116.
- [28] Y. Zhou, W. Yu, and D. Toupakaris, "Uplink multi-cell processing: Approximate sum capacity under a sum backhaul constraint," in *Proc. IEEE Inf. Theory Workshop (ITW)*, Sep. 2013, pp. 569–573.
- [29] S.-H. Park, O. Simeone, O. Sahin, and S. Shamai, "Joint precoding and multivariate backhaul compression for the downlink of cloud radio access networks," *IEEE Trans. Signal Process.*, vol. 61, no. 22, pp. 5646–5658, Nov. 2013.
- [30] M. Grant and S. Boyd, "CVX: Matlab software for disciplined convex programming, version 2.1," Jun. 2015. [Online]. Available: <http://cvxr.com/cvx/doc/CVX.pdf>
- [31] D. P. Bertsekas, *Nonlinear programming*, 2nd ed. Athena scientific, 1999.
- [32] G. Scutari, F. Facchinei, L. Lampariello, and P. Song, "Distributed methods for constrained nonconvex multi-agent optimization-part I: Theory," *submitted to IEEE Trans. Signal Process.*, 2014. [Online]. Available: <http://arxiv.org/abs/1410.4754>
- [33] F. A. Potra and S. J. Wright, "Interior-point methods," *J. Comput. Appl. Math.*, vol. 124, no. 1–2, pp. 281–302, Dec. 2000.
- [34] W. Yu, W. Rhee, S. Boyd, and J. M. Cioffi, "Iterative water-filling for Gaussian vector multiple-access channels," *IEEE Trans. Inf. Theory*, vol. 50, no. 1, pp. 145–152, Jan. 2004.
- [35] L. Vandenberghe, S. Boyd, and S.-P. Wu, "Determinant maximization with linear matrix inequality constraints," *SIAM J. Matrix Anal. Appl.*, vol. 19, no. 2, pp. 499–533, 1998.
- [36] A. Beck, A. Ben-Tal, and L. Tetruashvili, "A sequential parametric convex approximation method with applications to nonconvex truss topology design problems," *J. Global Optim.*, vol. 47, no. 1, pp. 29–51, May 2010.



Yuhuan Zhou (S'08) received the B.E. degree in Electronic and Information Engineering from Jilin University, Jilin, China, in 2005, the M.A.Sc. degree from the University of Waterloo, ON, Canada, in 2009, and the Ph.D. degree from the University of Toronto, ON, Canada, in 2016, both in Electrical and Computer Engineering. Since 2016, he has been with Qualcomm Technologies Inc., San Diego, CA, USA. His research interests include wireless communications, network information theory, and convex optimization.



Wei Yu (S'97-M'02-SM'08-F'14) received the B.A.Sc. degree in Computer Engineering and Mathematics from the University of Waterloo, Waterloo, Ontario, Canada in 1997 and M.S. and Ph.D. degrees in Electrical Engineering from Stanford University, Stanford, CA, in 1998 and 2002, respectively. Since 2002, he has been with the Electrical and Computer Engineering Department at the University of Toronto, Toronto, Ontario, Canada, where he is now Professor and holds a Canada Research Chair (Tier 1) in Information Theory and Wireless Communica-

tions. His main research interests include information theory, optimization, wireless communications and broadband access networks.

Prof. Wei Yu currently serves on the IEEE Information Theory Society Board of Governors (2015-17). He is an IEEE Communications Society Distinguished Lecturer (2015-16). He served as an Associate Editor for IEEE TRANSACTIONS ON INFORMATION THEORY (2010-2013), as an Editor for IEEE TRANSACTIONS ON COMMUNICATIONS (2009-2011), as an Editor for IEEE TRANSACTIONS ON WIRELESS COMMUNICATIONS (2004-2007), and as a Guest Editor for a number of special issues for the IEEE JOURNAL ON SELECTED AREAS IN COMMUNICATIONS and the EURASIP JOURNAL ON APPLIED SIGNAL PROCESSING. He was a Technical Program co-chair of the IEEE Communication Theory Workshop in 2014, and a Technical Program Committee co-chair of the Communication Theory Symposium at the IEEE International Conference on Communications (ICC) in 2012. He was a member of the Signal Processing for Communications and Networking Technical Committee of the IEEE Signal Processing Society (2008-2013). Prof. Wei Yu received a Steacie Memorial Fellowship in 2015, an IEEE Communications Society Best Tutorial Paper Award in 2015, an IEEE ICC Best Paper Award in 2013, an IEEE Signal Processing Society Best Paper Award in 2008, the McCharles Prize for Early Career Research Distinction in 2008, the Early Career Teaching Award from the Faculty of Applied Science and Engineering, University of Toronto in 2007, and an Early Researcher Award from Ontario in 2006. He was named a Highly Cited Researcher by Thomson Reuters in 2014.