# Centralized and Distributed Sparsification for Low-Complexity Message Passing Algorithm in C-RAN Architectures

Alessandro Brighente and Stefano Tomasin

Department of Information Engineering, University of Padova

via Gradenigo 6/B, 35131 Padova, Italy.

Email: alessandro.brighente@studenti.unipd.it, tomasin@dei.unipd.it

*Abstract*—**Cloud radio access network (C-RAN) is a promising technology for fifth-generation (5G) cellular systems. However the burden imposed by the huge amount of data to be collected (in the uplink) from the radio remote heads (RRHs) and processed at the base band unit (BBU) poses serious challenges. In order to reduce the computation effort of minimum mean square error (MMSE) receiver at the BBU the Gaussian message passing (MP) together with a suitable sparsification of the channel matrix can be used. In this paper we propose two sets of solutions, either centralized or distributed ones. In the centralized solutions, we propose different approaches to sparsify the channel matrix, in order to reduce the complexity of MP. However these approaches still require that all signals reaching the RRH are conveyed to the BBU, therefore the communication requirements among the backbone network devices are unaltered. In the decentralized solutions instead we aim at reducing both the complexity of MP at the BBU and the requirements on the RRHs-BBU communication links by pre-processing the signals at the RRH and convey a reduced set of signals to the BBU.**

*Index Terms*—**Cellular Systems; cloud radio access network (C-RAN); message passing (MP); Uplink.**

## I. Introduction

The fifth-generation (5G) of mobile communication systems has ambitious targets in terms (among others) of data rate, latency, number of supported users. Among the technologies envisioned to this end, cloud radio access network (C-RAN) may provide the flexibility in the deployment and planning of the network, combined with powerful energy-efficient computational resources [1].

Indeed, since the signal processing of multiple cells is implemented in the centralized facility of the base band unit (BBU), the computational resources are allocated on demand to the areas that have instantaneously more users, also with a better handling of inference and hand-off capabilities. On the other hand the need to process signals of many radio remote heads (RRHs) poses significant challenges to the BBU. Various approaches have been proposed to reduce the huge amount of data that is exchanged in this centralized approach, including suitable quantization of either the received signal [2] or the log-likelihood ratios (LLRs) [3]. On the other hand, also the signal processing itself at the BBU is very challenging, since even a minimum mean square error (MMSE) receiver requires the inversion of very large matrices. Similar problems are encountered in massive-multiple input multiple output (MIMO) systems with a huge number of users. About the reduction of signal processing burden in up-link detection, it has been proposed in [4] to cluster both users and RRHs based on the distance of terminals from RRH thus parallelizing MMSE operations into small size matrix operations. A further step forward has been done in [5] where it is proposed to implement the MMSE receiver by the message passing (MP). By exploiting the Gaussian distribution of the noise, a simple solution is obtained where the complexity per unit network area remains constant with growing network sizes. In particular [5] combines MP with the sparsification approach of [4], i.e., a first selection of users based on their distance from RRH reduces the size of the equivalent channel matrix before MP is applied.

In this paper we leverage on the results of [5] to propose two sets of solutions, either centralized or distributed ones. In the centralized solutions, we propose different approaches to sparsify the channel matrix, in order to reduce the complexity of MP. However these approaches still require that all signals reaching the RRH are conveyed to the BBU, therefore the communication requirements among the backbone network devices are unaltered. In the decentralized solutions instead we aim at reducing both the complexity of MP at the BBU and the requirements on the RRHs-BBU communication links by pre-processing the signals at the RRH and conveying a reduced set of signals to the BBU.

The rest of the paper is organized as follows. We first introduce the system model in Section II. Then we propose the centralized sparsification techniques in Section III. The decentralized sparsification methods are discussed in Section IV. Numerical results are presented in Section V, before conclusions are obtained in Section VI.

Notation: matrices and vectors are denoted in boldface. $x^T$ and $x^H$ denote the transpose and Hermitian of vector $x$, respectively.

## II. System Model

We consider the up-link of a cellular network with $N_c$ cells, each one containing a base station (BS) equipped with $N_a$ omnidirectional receive antennas (RRHs). Each cell is populated by $N_u$ mobile terminals (MTs) uniformly distributed

over the entire cell area, each one equipped with a single antenna and transmitting with power $P$.

The overall network can be seen as a MIMO system, where the unit-power column vector $\boldsymbol{x}$ of size $K = N_c N_u$ comprises the data signals of MTs scaled by $\sqrt{P}$ before transmission, whereas column vector $\boldsymbol{y}$ of size $N = N_c N_a$ comprises all signals received by RRHs. The MIMO channel model of the up-link from MTs to the RRHs can be written as

$$\boldsymbol{y} = \sqrt{P}\boldsymbol{H}\boldsymbol{x} + \boldsymbol{w}\,, \tag{1}$$

where $\boldsymbol{H}$ is the $N \times K$ channel matrix with entries $[\boldsymbol{H}]_{i,j}$ and $\boldsymbol{w}$ is the additive white Gaussian noise (AWGN) vector with independent and identically distributed (i.i.d.) complex Gaussian entries with zero-mean and variance $N_0$.

The signals received by the RRHs are forwarded to the BBU that aims at performing the MMSE receiver, i.e., computing

$$\hat{\boldsymbol{x}} = P^{\frac{1}{2}}\boldsymbol{H}^H(P\boldsymbol{H}\boldsymbol{H}^H + N_0\boldsymbol{I})^{-1}\boldsymbol{y}. \tag{2}$$

### A. Randomized Gaussian MP decoder

The MP algorithm can be used to solve the interference problem over sparse factor graphs [6], therefore providing the solution of the MMSE receiver (2). Since the received signal is affected by Gaussian noise we can use the Gaussian message-passing (GMP) solution, and in particular we focus on the randomized randomized GMP (RGMP) of [5] which has been shown to have better convergence properties. In order to obtain the MMSE estimate of the transmitted signal $\boldsymbol{x}$ the proposed RGMP Algorithm exploits the knowledge of the statistical description of all the elements in (1) and iteratively updates the values of mean and variance of all components of both $\boldsymbol{x}$ and $\boldsymbol{y}$ vectors. The Algorithm stops updating these values when a stopping criterion is satisfied and the MMSE estimate of $\boldsymbol{x}$ is returned.

The computational complexity of the RGMP Algorithm is $\mathcal{O}(NK^2)$, hence it depends on the number of users (growing quadratically with it) and receiving antennas of the system. In large systems, with many MTs and RRHs, the decoding process is therefore prohibitively complex. An approach to reduce the complexity is to reduce the number of non-zero entries in $\boldsymbol{H}$ over which the MP is run, i.e. applying MP on a sparsified version of $\boldsymbol{H}$. Note that the sparsification on the one side will reduce the complexity, while on the other side provides an approximation of $\hat{\boldsymbol{x}}$, thus reducing the ASR (ASR) of the system.

Different approaches will be analysed in the following sections: a centralised approach, where sparsification is performed at the BBU pool before RGMP decoding, and a distributed approach, where sparsification is applied as pre-coding operations at each BS.

### III. CENTRALIZED SPARSIFICATION METHODS

With centralized sparsification methods the decoding process is entirely demanded to the central BBU pool. Then the signal received at the RRH, down-converted to base-band and converted to the digital domain, is entirely forwarded

to the BBU. Hence no local processing is performed at the BS. Since no pre-processing operation is done at the BS in order to reduce the computational complexity of the decoding process, this latter task is demanded to the central BBU. We here introduce and discuss different approaches to sparsify the channel matrix by performing operations on its entries at the BBU.

### A. Sparsification based on the received power (CRPS)

The first approach is based on the received power. In particular, we set to zero the channel matrix coefficients having power below a threshold value $P_{\min}$.

We thus obtain matrix $\hat{\boldsymbol{H}}$ with entries

$$[\hat{\boldsymbol{H}}]_{i,j} = \begin{cases} [\boldsymbol{H}]_{i,j} & \text{if } |[\boldsymbol{H}]_{i,j}|^2 \geq P_{\min} \\ 0 & \text{otherwise.} \end{cases} \tag{3}$$

The neglected coefficients can be accounted for as additional noise into the system. In particular, defining the error matrix $\tilde{\boldsymbol{H}} = \boldsymbol{H} - \hat{\boldsymbol{H}}$ the statistical power of noise and error $N_0$ becomes

$$\hat{N}_0 = N_0 + \frac{1}{N}\sum_{n=1}^{N}\sum_{k=1}^{K}|[\tilde{\boldsymbol{H}}]_{n,k}|^2 \tag{4}$$

RGMP is then run over channel $\hat{\boldsymbol{H}}$ and considers as noise power $\hat{N}_0$.

### B. Sparsification based on semi-orthogonality (MCOS)

The second proposed approach is based on MT channels semi-orthogonality. Let us consider singularly each BS: we notice that MTs having orthogonal channels do not interfere. Now, assuming that each MT signal is mainly detected by the antennas of its cell, we can ignore the contribution of the external MTs since they will not significantly contribute to the computation of the MMSE.

In formulas, let us consider the channel row vector $\boldsymbol{h}_{k_1} = [\boldsymbol{H}]_{n_1,k_1}, [\boldsymbol{H}]_{n_2,k_1}, ..., [\boldsymbol{H}]_{n_{N_a},k_1}$ from MT $k_1$ to all RRHs belonging to a certain BS with indexes in the set $\mathcal{A} = \{n_1, n_2, ..., n_{N_a}\}$. The orthogonality among channels toward the same BS is established by the internal product of the channels and we consider that two channels are semi-orthogonal if the product is below a threshold $T_{\text{prod}}$, i.e.,

$$|\boldsymbol{h}_{k_1}\boldsymbol{h}_{k_2}^H|^2 < T_{\text{prod}}. \tag{5}$$

If MTs $k_1$ outside the cell $i$ is semi-orthogonal to all MTs inside the cell, then entries of channel matrix $\boldsymbol{H}$ corresponding to the link between MT $k_1$ and all RRHs of BS $i$ are set to zero.
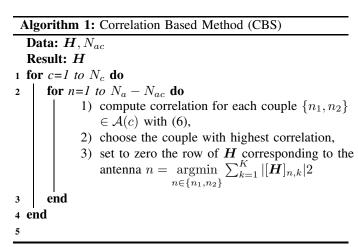
### C. Sparsification based on the correlation

The idea is to reduce the number of rows of the channel matrix by selecting the subset $\mathcal{S}$ of the antennas $\mathcal{A}(c)$ located in cell $c$. In order to chose a suitable subset and, hence, which rows to delete, we exploit the algorithms presented in [7], i.e. correlation based sparsification (CBS) and mutual information based sparsification (MIBS). We denote by $N_{ac}$ the number of

antennas, and hence the number of rows of the channel matrix relative to $c$ used for decoding.

In formulas, we consider couples $\{n_1, n_2\}$ of antennas and channel matrix rows $\boldsymbol{g}_{n \in \mathcal{A}(c)} = [[\boldsymbol{H}]_{n,1}, [\boldsymbol{H}]_{n,2}, ..., [\boldsymbol{H}]_{n,K}]$, belonging to set $\mathcal{A}$ of cell $c$ and measure their correlation as

$$c_{n_1, n_2} = |\boldsymbol{g}_{n_1} \boldsymbol{g}_{n_2}^H|^2. \tag{6}$$

For each cell the correlation between couples of antenna channels belonging to the considered cell is computed. Then the couple with highest correlation is selected and the antenna channel with lowest power is discarded. Its corresponding row in the channel matrix is hence set to zero. This procedure is repeated until we set to zero a number of rows equal to $N_a - N_{ac}$. A description of this method is provided in Algorithm 1.

---

**Algorithm 1:** Correlation Based Method (CBS)

**Data:** $\boldsymbol{H}, N_{ac}$
**Result:** $\boldsymbol{H}$
1 **for** *c=1 to $N_c$* **do**
2     **for** *n=1 to $N_a - N_{ac}$* **do**
        1) compute correlation for each couple $\{n_1, n_2\}$ $\in \mathcal{A}(c)$ with (6),
        2) choose the couple with highest correlation,
        3) set to zero the row of $\boldsymbol{H}$ corresponding to the antenna $n = \underset{n \in \{n_1, n_2\}}{\operatorname{argmin}} \sum_{k=1}^{K} |[\boldsymbol{H}]_{n,k}|2$
3     **end**
4 **end**
5

---

### D. Sparsification based on the mutual information

This antenna selection approach, MIBS, behaves similarly to Algorithm 1, except that correlation in step 1 is substituted by the normalized mutual information. The mutual information for a couple $\{n_1, n_2\} \in \mathcal{A}(c)$ is computed as

$$I(n_1, n_2) = \log_2 \left( \frac{\|\boldsymbol{h}_{n_1}\|^2 \|\boldsymbol{h}_{n_2}\|^2}{\|\boldsymbol{h}_{n_1}\|^2 \|\boldsymbol{h}_{n_2}\|^2 + |\boldsymbol{h}_{n_1} \boldsymbol{h}_{n_2}^H|^2} \right) \tag{7}$$

whereas its normalized version is

$$I_0(n_1, n_2) = \frac{I(n_1, n_2)}{\min\{|\log_2 \|\boldsymbol{h}_{n_1}\|^2|, |\log_2 \|\boldsymbol{h}_{n_2}\|^2|\}}. \tag{8}$$

In Algorithm 1 we replace (6) with (8). In both CBS and MIBS the noise power $N_0$ is not modified as in (4), because, when deleting an antenna channel (and hence a channel matrix row), we assume that its information is contained in the other rows of the considered couple.

## IV. DISTRIBUTED SPARSIFICATION METHODS

The centralized sparsification approach has the drawback that the entire received signal is forwarded from RRHs to the central BBU. Since the requirements for a front-haul link are very stringent (multi-gigabit-per-second-capacity and few-milliseconds latency [8]) and this amount of data turns out to be prohibitively high for satisfying this requirements, we consider distributed sparsification solutions, which aim together at reducing both the decoding computational complexity and the amount of data flowing through the front-haul.

In this section we will discuss sparsification applied as pre-coding at the BS of each cell before forwarding the received signals to the BBU. Let $\boldsymbol{y}_c$ be the received $N_a$-size column vector signal at the BS of cell $c$. If we consider a pre-coding $N_r \times N_a$ matrix $\boldsymbol{B}$ for cell $c$ and we multiply it by the received signal we obtain

$$\boldsymbol{B}\boldsymbol{y}_c = \boldsymbol{B}\boldsymbol{H}_{N_r}\sqrt{P}\tilde{\boldsymbol{x}} + \boldsymbol{B}\boldsymbol{H}_{\bar{N}_r}\sqrt{P}\boldsymbol{i} + \boldsymbol{B}\boldsymbol{w}, \tag{9}$$

where $\tilde{\boldsymbol{x}}$ is the vector containing signals coming from MTs in set $\mathcal{M}$ (as later discussed), $\boldsymbol{H}_{N_r}$ is the sub-channel matrix composed by the columns of $\boldsymbol{H}$ for users considered in $\mathcal{M}$, $\boldsymbol{H}_{\bar{N}_r}$ is the sub-channel matrix composed by the column of $\boldsymbol{H}$ for users $\notin \mathcal{M}$ and is the vector containing signals coming from users $\notin \mathcal{M}$.

Pre-coding matrix $\boldsymbol{B}$ can assume different forms and consider different number and types of users. In particular, we let $\boldsymbol{G}$ be the sub-channel matrix of users in $\mathcal{M}$. Then we set $\boldsymbol{B} = \boldsymbol{G}^H$, i.e. $\boldsymbol{B}$ assumes to form of the matched matrix to the considered channel. A second option provides that $\boldsymbol{B}$ is the zero-forcing matrix, i.e. $\boldsymbol{B} = \boldsymbol{G}^H(\boldsymbol{G}\boldsymbol{G}^H)^{-1}$. In the following we define different strategies to select $\mathcal{M}$.

### A. Selection based on the position (PSS)

We first assume the knowledge of users location and, in particular, we know the cell each user belongs to. Then $\mathcal{M}$ is the set of users located in cell $c$, with $|\mathcal{M}| = N_u$. Matrix $\boldsymbol{B}$ will hence be a $N_u \times N_a$ dimesnion matrix. Such a pre-coding operation hence reduces the number of rows of the sub-channel matrix of each cell from $N_a$ (the number of antennas of the considered BS) to $N_u$. We notice that, with the pre-coding operation, noise vector entries are correlated and that the MP algorithm must be modified. Since noise power remains the same in all branches the noise level depends on $n$ and becomes

$$N_0(n) = N_0 \sum_{k=1}^{N_u} |[\boldsymbol{B}]_{n,k}|^2, \tag{10}$$

with $n \in \{1, ..., N_u\}$, which takes into account correlation introduced by matrix $\boldsymbol{B}$ in each receiver branch. This new version of RGMP will be considered as default for henceforth presented methods. Note that this approach is sub-optimal respect to MMSE as the MP solution in this case neglects the correlation among the noise components.

### B. Selection based on received power (DRPS)

In this approach MTs are selected according to the received power. We select the $N_p$ users with highest power reaching the BS of cell $c$, i.e. given the channel from user $k$ to the BS in $c$, we compute the received power (11) for each user in the cellular network,

$$p(k) = \sum_{n \in \mathcal{A}(c)} |[\boldsymbol{H}]_{n,k}|^2 \tag{11}$$

and consider the $N_p$ users with highest $p(k)$ toward the BS of cell $c$. The channel matrix columns of this set of users will then compose the columns of matrix $G$ for cell $c$.

### C. Selection based on mixed criterion (MSS)

The third approach is a mix of the first two. In fact matrix $G$ collects columns of both users located in cell $c$ and the $N_p$ most powerful users, i.e. with highest $p(k)$, located outside cell $c$.

## V. NUMERICAL RESULTS

We here first present the ASR results obtained for all the sparsification methods introduced in previous sections and then discuss their computational complexity. Mostly the trade-off between ASR and computational complexity is analyzed. We consider a scenario with $N_c = 16$ cells, each one equipped with a BS with $N_a = 8$ RRHs. Each cell contains $N_u = 4$ users and each user is allocated the same transmitting power $P = 1$. Noise power is chosen to have a border cell signal to noise ratio (SNR) of 0 dB. In the following we assume that $H$ is affected by both path loss (with coefficient $\alpha = 2$) and Rayleigh fading, so that each entry is a zero-mean complex Gaussian random variable with variance equal to the inverse of the distance from the considered MT and the considered antenna of the BS. Channel matrix entries are i.i.d.

The RGMP Algorithm is stopped when the mean of the transmitted signal does not change more than 1% in one iteration. Each method has been compared both in terms of sparsification level, i.e. the number of entries of the channel matrix $H$ different from zero after sparsification, and channel ASR. All results have been compared with those of pure RGMP, i.e. without channel sparsification.

### A. Centralized sparsification

We consider first the centralized sparsification. Fig. 1 reports the mean ASR values vs. SNR for two parameter values of each centralized sparsification method and for RGMP without channel sparsification. ASR results for MIBS are analogous to the ones obtained with CBS, and are not reported here for brevity. With all the presented methods we can obtain good results in terms of ASR values, comparable or equal to that obtained with RGMP without channel sparsification.

### B. Distributed sparsification

Distributed sparsification has been implemented for both matched and zero forcing matrix $B$. Fig. 2 reports mean ASR values vs. SNR obtained for the maximum and minimum considered users by distributed sparsification methods and for RGMP without channel sparsification. We denoted the different methods with their acronym followed by the number of considered users. We can see that the matched implementation of $B$ outperforms the zero-forcing implementation is terms of mean ASR. Furthermore the matched implementation of all methods considering the maximum number of users, allows a better exploitation of the channel for low SNR values obtaining mean ASR values equal to the ones obtained with RGMP without channel sparsification.
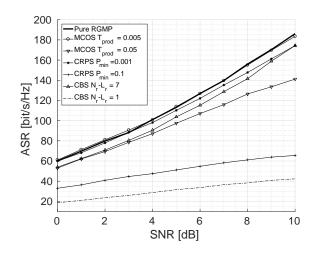


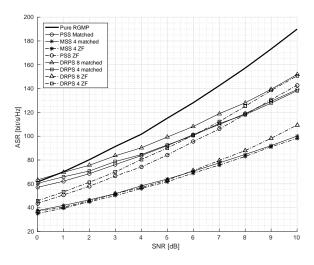Fig. 1. Mean ASR vs. SNR for centralized sparsification.



Fig. 2. Mean ASR vs $SNR$ for distributed sparsification.

### C. Computational complexity analysis

We now analyse the computational complexity of the different approaches in terms of number of decoding operations after sparsification. This depends on the number of entries of $\hat{H} \neq 0$ as each requires two sums over the total number of users $K$, operations that are repeated until the stopping criterion is satisfied. Hence the total number of operations is

$$N_{op} = 2\,K\,s\,I, \tag{12}$$

where $s$ denotes the number of channel matrix entries different from 0, and $I$ the number of message passing iterations needed to satisfy the stopping criterion. Fig. 3 shows the ASR vs. the number of operations needed for the decoding process for the centralized sparsification methods with an SNR level of 0 dB. We notice that with semi-orthogonal-based sparsification we obtain the best performing system, with an achievable sum rate of 58 bit/s/Hz and a computational complexity of $9.2 \cdot 10^5$ operations. However notice that this implementation is not the

| Sparsification method | Sparsification level | # of operations | ASR [bit/s/Hz] |
|---|---|---|---|
| Pure RGMP | 8192 | 2097152 | 60 |
| CRPS, $P_{\min} = 0.001$ | 4121 | 1582464 | 60.19 |
| MCOS, $T_{\mathrm{prod}} = 0.001$ | 4288 | 1097728 | 58.83 |
| CBS, $L_r = 1$ | 7168 | 1835008 | 55.72 |
| MIBS, $L_r = 1$ | 7168 | 917504 | 58.73 |
| PSS, $\boldsymbol{B} = \boldsymbol{G}^H$ | 4096 | 1835008 | 57 |
| MSS, 6 usr. $\boldsymbol{B} = \boldsymbol{G}^H$ | 6144 | 1835008 | 63.5 |
| DRPS , 4 usr. $\boldsymbol{B} = \boldsymbol{G}^H$ | 4096 | 1966080 | 61.28 |
| DRPS, 4 usr. $\boldsymbol{B} = \boldsymbol{G}^H(\boldsymbol{G}\boldsymbol{G}^H)^{-1}$ | 4096 | 393216 | 46.2 |
| DRPS , 8 usr. $\boldsymbol{B} = \boldsymbol{G}^H$ | 8192 | 4082131 | 63.25 |

best performing in terms of achievable sum rate, instead it is the best compromise between computational complexity and ASR. Notice that RGMP without channel sparsification obtains an ASR of 60 bit/s/Hz with a computational complexity of $2.1 \cdot 10^6$ operations. Hence the reduction of $1 \cdot 10^6$ operations comes with an ASR loss of 2 bit/s/Hz.
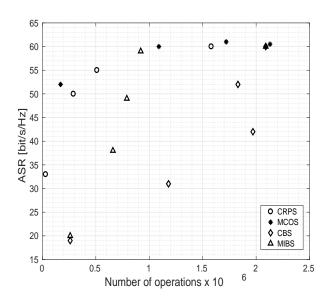


Fig. 3. Trade-off between computational complexity and achievable sum-rate for centralized sparsification methods: 0 dB SNR.

Fig. 4 reports the ASR vs. the number of operations needed for the decoding process for the centralized sparsification methods with an SNR level of 0 dB. We notice that the best compromise between ASR and computational complexity is obtained for MSS with matched matrix, which presents an ASR of approximately 63 bit/s/Hz with a computational complexity of $2 \cdot 10^5$ operations.

Table I reports the obtained computational complexity and ASR values for the best performing parameter of each method when SNR value is 0 dB. A trade-off can be obtained, since we
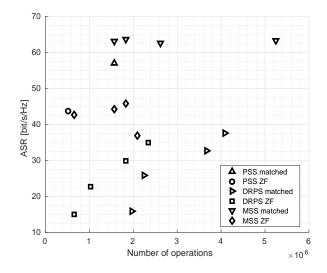


Fig. 4. Trade-off between computational complexity and achievable sum-rate for distributed sparsification methods: 0dB SNR

want to maximize the ASR while maintaining a low computational complexity. We can hence state that all methods present a channel ASR comparable to the one obtained with pure RGMP, but generally need a significantly lower number of decoding operations. The best performing among all presented methods in terms of both computational complexity and ASR is MIBS sparsification when SNR value is 0 dB. This method needs less than half of the number of operations required by pure RGMP with an ASR loss of approximately 2 bit/s/Hz.

## VI. CONCLUSIONS

For a C-RAN system where signals coming from many RRHs we have considered the problem of implementing a MMSE receiver at the BBU. In order to decrease the computational complexity a RGMP algorithm has been considered, and suitable sparsifications of the channel matrix have been introduced. We considered both centralized approaches, performed at the BBU and requiring a complete transfer of received signals from the RRHs and decentralized solutions

where a pre-processing is performed at the BS. This latter solution not only has been shown to be effective in terms of reduction of the computational complexity of the decoding process, but also of the amount of data flowing from the BSs to the BBU, and hence of the front-haul network capacity as well as the centralization overhead. Numerical results have shown a variety of trade-off between complexity and performance (in terms of ASR) confirming that the proposed solutions are promising for an implementation of these approaches in 5G C-RAN systems.

## REFERENCES

[1] C. Liu, K. Sundaresan, M. Jiang, S. Rangarajan and G. K. Chang, "The case for re-configurable backhaul in cloud-RAN based small cell networks," in *Proc. IEEE INFOCOM 2013*, pp. 1124-1132.

[2] P. Baracca, S. Tomasin and N. Benvenuto, "Constellation Quantization in Constrained Backhaul Downlink Network MIMO," *IEEE Trans. on Commun.*, vol. 60, no. 3, pp. 830-839, March 2012.

[3] K. Miyamoto, S. Kuwano, J. Terada and A. Otaka, "Uplink Joint Reception with LLR Forwarding for Optical Transmission Bandwidth Reduction in Mobile Fronthaul," in *Proc.2015 IEEE 81st Vehicular Technology Conference (VTC Spring)*, Glasgow, 2015.

[4] C. Fan, Y. J. Zhang, and X. Yuan, "Dynamic nested clustering for parallel PHY-layer processing in cloud-RANs," *IEEE Trans. Wireless Commun.*, vol. 15, no. 3, pp. 1881-1894, Mar. 2016.

[5] C.Fan,X.Yuan,Y.J.A.Zhang *Scalable Uplink Signal Detection in C-RANs via Randomized Gaussian Message Passing*, arXiv:1511.09024, May 2016.

[6] F. R. Kschischang, B. J. Frey, and H.-A. Loeliger, "Factor graphs and the sum-product algorithm," *IEEE Trans. Inf. Theory*, vol. 47, no. 2, pp. 498519, Feb. 2001.

[7] A. F. Molish, M. Z. Win, Y. S. Choi, J. H. Winters, "Capacity of MIMO systems with antenna selection," *IEEE Trans. Wireless Commun.*, vol. 4, no. 4, July 2005.

[8] J. Bartlet, P. Rost, D. Wubben, J. Lessmann, B. Melis, G. Fettweis, "Fronthaul and Backhaul Requirements of Flexibly Centralized Radio Access", *IEEE Wireless Communications*, October 2015.