

Capacity Bounds for One-Bit MIMO Gaussian Channels With Analog Combining

Neil Irwin Bernardo¹, Graduate Student Member, IEEE, Jingge Zhu², Member, IEEE,
Yonina C. Eldar³, Fellow, IEEE, and Jamie Evans⁴, Senior Member, IEEE

Abstract—The use of 1-bit analog-to-digital converters (ADCs) is seen as a promising approach to significantly reduce the power consumption and hardware cost of multiple-input multiple-output (MIMO) receivers. However, the nonlinear distortion due to 1-bit quantization fundamentally changes the optimal communication strategy and also imposes a capacity penalty to the system. In this paper, the capacity of a Gaussian MIMO channel in which the antenna outputs are processed by an analog linear combiner and then quantized by a set of zero threshold ADCs is studied. A new capacity upper bound for the zero threshold case is established that is tighter than the bounds available in the literature. In addition, we propose an achievability scheme which configures the analog combiner to create parallel Gaussian channels with phase quantization at the output. Under this class of analog combiners, an algorithm is presented that identifies the analog combiner and input distribution that maximize the achievable rate. Numerical results are provided showing that the rate of the achievability scheme is tight in the low signal-to-noise ratio (SNR) regime. Finally, a new 1-bit MIMO receiver architecture which employs analog temporal and spatial processing is proposed. The proposed receiver attains the capacity in the high SNR regime.

Index Terms—MIMO communications, quantization, analog combining, capacity.

I. INTRODUCTION

THE use of multiple-input multiple-output (MIMO) technology has attracted considerable attention during the

Manuscript received 8 April 2022; revised 28 July 2022; accepted 12 September 2022. Date of publication 22 September 2022; date of current version 18 November 2022. This work was supported in part by the Australian Research Council under Project DE210101497, in part by the European Research Council (ERC) under the European Union's Horizon 2020 Research and Innovation Program under Grant 101000967, and in part by the Israel Science Foundation under Grant 536/22. The work of Neil Irwin Bernardo was supported in part by the Melbourne Research Scholarship of The University of Melbourne and in part by the Department of Science and Technology-Engineering Research and Development for Technology (DOST-ERDT) Faculty Development Fund of the Republic of the Philippines. The associate editor coordinating the review of this article and approving it for publication was N. Liu. (*Corresponding author: Neil Irwin Bernardo.*)

Neil Irwin Bernardo is with the Department of Electrical and Electronic Engineering, The University of Melbourne, Parkville, VIC 3010, Australia, and also with the Electrical and Electronics Engineering Institute, University of the Philippines Diliman, Quezon City 1101, Philippines (e-mail: bernardon@student.unimelb.edu.au).

Jingge Zhu and Jamie Evans are with the Department of Electrical and Electronic Engineering, The University of Melbourne, Parkville, VIC 3010, Australia (e-mail: jingge.zhu@unimelb.edu.au; jse@unimelb.edu.au).

Yonina C. Eldar is with the Faculty of Mathematics and Computer Science, Weizmann Institute of Science, Rehovot 7610001, Israel (e-mail: yonina.eldar@weizmann.ac.il).

Color versions of one or more figures in this article are available at <https://doi.org/10.1109/TCOMM.2022.3208632>.

Digital Object Identifier 10.1109/TCOMM.2022.3208632

last two decades due to the significant capacity enhancement it offers. This is evident in the number of wireless broadband standards that have incorporated MIMO technology into their specifications. However, even though its theoretical gains are well-established in the literature, some practical concerns (e.g. hardware cost, power consumption) arise as more antenna elements are connected to the communication device [1]. Furthermore, recent advancements in next generation mobile technology have pushed forward the use of much wider transmission bandwidth and much larger antenna arrays to achieve its performance targets [2], [3], [4]. Consequently, there is an increasing demand for new MIMO receiver designs that are energy-efficient and are able to *reliably* support high data-rate applications.

High-speed and high-resolution analog-to-digital converters (ADCs) are one of the primary contributors to the power consumption of wireless receivers. The ADC power consumption scales linearly with the sampling rate and exponentially with the number of quantization bits per sample, regardless of the topology [5]. Thus, one straightforward design methodology in implementing low-power MIMO receivers is to simply replace the high-resolution ADCs connected to each radio frequency (RF) chain with low-resolution counterparts (usually 1-bit). We shall refer to this as the *conventional* low-resolution ADC design. This receiver design is particularly attractive for massive MIMO systems due to the hardware scaling law – that is, the impact of hardware imperfections on the overall spectral efficiency diminishes as the number of antenna elements grows [6]. In fact, some early results suggest that the power savings and cost reduction obtained from using conventional 1-bit ADC MIMO systems may outweigh the rate loss caused by severe quantization; even with simple linear detection schemes [1], [7], [8]. Error rate analyses of single-input multiple-output (SIMO) fading channels with low-resolution output quantization have established that, under certain conditions, the diversity order of optimal detectors is nonzero and improves linearly with the number of receive antennas [9], [10]. Furthermore, other receiver functionalities (e.g. timing recovery, channel estimation) have been shown to work with acceptable performance even if the receiver is equipped with low-resolution ADCs [11], [12], [13], [14].

The notable benefits of the conventional low-resolution receiver design have sparked interest in investigating the information-theoretic limits of communication channels with output quantization. The work of Singh *et al.* [15] appears to be the first to present exact results on the capacity of real additive white Gaussian noise (AWGN) channel with

low-resolution ADCs. Several works that followed characterized the capacity-achieving input of various single-input single-output (SISO) channels with output quantization [16], [17], [18], [19], [20], [21], [22]. Yet, while the main motivation of using low-resolution ADCs is for scalable implementation of massive MIMO systems, the aforementioned information-theoretic results have not been extended to the MIMO setting. To date, the capacity of quantized MIMO remains elusive and is analyzed through capacity bounds [23], [24], [25], [26]. These bounds, however, can be loose in certain signal-to-noise ratio (SNR) regimes and problem settings. For instance, the tightness of the finite SNR upper bound and channel inversion lower bound established in [24] depends on the row rank and condition number of the channel. The additive quantization noise model (AQNM) capacity lower bound in [23] and [25] is also shown to be loose in the high SNR regime; thus hinting at the suboptimality of Gaussian signaling for quantized MIMO channels.

Recent literature surveys [27], [28], [29] have examined new receiver architectures that are energy-efficient and cost-effective but incur less performance degradation than conventional low-resolution ADC systems. These architectures include the mixed-ADC receivers [30], hybrid analog/digital receivers [31], [32], [33], [34], and machine learning (ML)-based receivers [35], [36]. In [37], the authors suggest two new receiver designs, namely the *hybrid blockwise receiver* and *adaptive threshold receiver*, that perform analog spatial and temporal processing prior to 1-bit quantization. The idea of using analog linear combiners prior to 1-bit quantization, called the *hybrid one-shot receiver*, was initially proposed in [38] as a means to compare performance of different MIMO systems under various output quantization constraints. Analog temporal processing was then incorporated in this quantized MIMO framework to demonstrate a fundamental tradeoff between latency and maximum achievable rate of quantized MIMO channels in the high SNR regime [39].

In this paper, we first look at the performance of the hybrid one-shot receiver with zero threshold ADCs. In other words, we only observe the signs of the analog linear combiner outputs. The zero threshold ADC setup is expected to yield a lower capacity than its non-zero threshold counterpart. Yet, the zero threshold case is still interesting since it eliminates the need for an automatic gain control (AGC), which is otherwise required in multi-level quantization to match the dynamic range of the received signal [40]. In addition, the hybrid one-shot receiver with zero threshold ADCs is essentially the 1-bit ADC version of the hybrid analog-digital acquisition system used for task-based quantization [33], [34]. We then introduce analog temporal processing to the receiver. A new MIMO receiver architecture is proposed which uses analog domain pipelining and adaptive phase shifters to attain higher achievable rate than the hybrid one-shot receiver with zero threshold ADCs. Our main contributions are summarized as follows:

- We provide an achievability scheme for the hybrid one-shot receiver with zero threshold ADCs. The achievability scheme formulates the capacity problem as a nonconvex resource allocation problem. To this

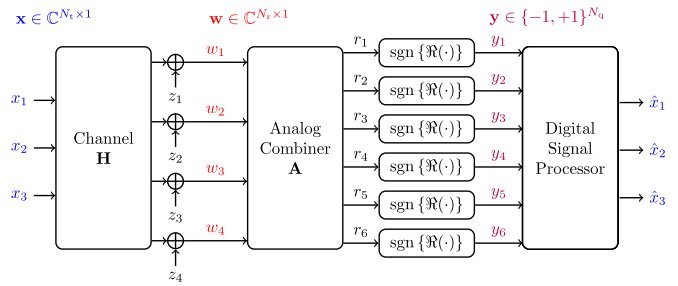


Fig. 1. System Model of MIMO with Analog Combining prior to Output Quantization. Here, $N_t = 3$, $N_r = 4$, and $N_q = 6$.

end, an alternating optimization approach is presented to obtain a local optimal solution to this nonconvex problem.

- Using the data processing inequality (DPI), we establish a new capacity upper bound in the finite SNR regime for the hybrid one-shot receiver with zero threshold ADCs. More precisely, by showing that the amplitude information is discarded in the zero threshold ADC case, we obtain an upper bound that is tighter than the truncated Shannon capacity (TSC) used in [37].
- Through numerical evaluation, we show that the output produced by our alternating optimization approach is tight in the low SNR regime. However, a gap between the capacity upper bound and achievability scheme exists in the high SNR regime. We characterize this gap as a function of the number of eigenchannels and number of sign quantizers.
- To close the gap mentioned above, we introduce a new ADC mechanism, called *pipelined phase ADC*, which incorporates analog temporal processing in the quantization process. By incorporating this ADC mechanism to the receiver, we show that the high SNR capacity can be attained. We compare the achievable rate of our proposed receiver to that of the adaptive threshold receiver in [37]. Numerical results are presented showing that the proposed receiver outperforms the adaptive threshold receiver when the eigenvalues of the channel are equal.

The rest of the paper is organized as follows: Section II formulates the system model and the problem we aim to address. Sections III and IV present the details of the achievability scheme and the derivation of the capacity upper bound, respectively, for the hybrid one-shot receiver with zero threshold ADCs. Section V provides numerical results and analysis for the achievability scheme and the capacity upper bound established in Sections III and IV. Sections VI and VII discuss the proposed receiver along with analysis of its achievable rate. We also compare its performance with existing work. Finally, Section VIII concludes the paper.

II. PROBLEM FORMULATION

We consider a discrete-time memoryless MIMO system shown in Figure 1 with N_t transmit antennas and N_r receive antennas. The input vector \mathbf{x} satisfies the power constraint $\mathbb{E}[\|\mathbf{x}\|^2] \leq P$. The input-output relationship between transmitted symbol $\mathbf{x} \in \mathbb{C}^{N_t \times 1}$ and received

symbol $\mathbf{w} \in \mathbb{C}^{N_r \times 1}$ is

$$\mathbf{w} = \mathbf{H}\mathbf{x} + \mathbf{z}, \quad (1)$$

where $\mathbf{H} \in \mathbb{C}^{N_r \times N_t}$ is the fixed MIMO channel gain known at the transmitter and receiver, and $\mathbf{z} \sim \mathcal{CN}(\mathbf{0}, \sigma^2 \mathbf{I}_{N_r})$ is a circular-symmetric zero-mean complex Gaussian noise vector with noise variance σ^2 . We further assume that \mathbf{H} is a full rank matrix. The receiver is equipped with N_q sign quantizers and an analog linear combiner $\mathbf{A} \in \mathbb{C}^{N_q \times N_r}$ preprocesses the signal before quantization. This receiver structure has been widely adopted in various applications, such as in task-based quantization [33], [34]. The output vector $\mathbf{y} \in \mathbb{C}^{N_q \times 1}$ can be written as

$$\mathbf{y} = \text{sign}(\Re\{\mathbf{A}\mathbf{w}\}) = \text{sign}(\Re\{\mathbf{A}\mathbf{H}\mathbf{x} + \mathbf{z}'\}), \quad (2)$$

where $\mathbf{z}' = \mathbf{A}\mathbf{z}$ and the $\text{sign}\{\cdot\}$ function is applied to every element of the real vector.

We note the differences between our problem setup and the system model of the hybrid one-shot receiver formulated in [37] and [41]. In their work, the quantized output vector \mathbf{y} is expressed as $\mathbf{y} = \text{sign}(\mathbf{A}\mathbf{w} + \mathbf{t})$, where $\mathbf{A} \in \mathbb{R}^{N_q \times N_r}$, $\mathbf{w} \in \mathbb{R}^{N_r \times 1}$, and $\mathbf{t} \in \mathbb{R}^{N_q \times 1}$. In contrast, the entries of \mathbf{x} , \mathbf{H} , and \mathbf{A} in our problem setup are complex-valued quantities. The threshold vector \mathbf{t} is also set to be the all-zero vector. The $\mathbf{t} = \mathbf{0}$ case is particularly appealing from a practical viewpoint since receivers employing zero-threshold ADCs do not require AGC during the data detection phase [1]; thus further reducing the hardware complexity and cost. Moreover, by restricting the setup to $\mathbf{t} = \mathbf{0}$, we are able to derive tighter capacity bounds for this problem setup than simply plugging in $\mathbf{t} = \mathbf{0}$ to the capacity bounds established in [37] and [41] for general \mathbf{t} .

For a fixed configuration of the analog combiner \mathbf{A} , the capacity expression of model (2) can be written as

$$C(\mathbf{A}) = \max_{F_{\mathbf{x}}} I_{\mathbf{A}}(\mathbf{x}; \mathbf{y}), \quad (3)$$

where we used the subscript \mathbf{A} in (3) to indicate that the mutual information between the transmitted signals and the sign quantizer outputs are induced by a choice of \mathbf{A} . With slight abuse of notation, we use $F_{\mathbf{x}}$ to refer to $F_{\mathbf{x}}(\mathbf{x})$. We are interested in the largest input-output mutual information over all configurations of the analog linear combiner. Mathematically, we seek for the capacity C , which is defined as

$$C = \max_{\mathbf{A}} C(\mathbf{A}), \quad (4)$$

as well as the optimal \mathbf{A} and $F_{\mathbf{x}}$ that achieve this capacity.

It is known that, under general assumptions, the maximizing input distribution of a channel with finite output cardinality is discrete with finite number of mass points [42]. By combining the concavity of the mutual information over the input distribution and the discreteness of the optimal input distribution, numerical algorithms [43], [44] can be used to obtain the input distribution that solves problem (3) in the general case. The computation, however, may involve multi-dimensional integration, and the complexity would grow exponentially with the number of sign quantizers. For the special case of $N_r = N_q$ and $\mathbf{A} = \mathbf{I}_{N_q \times N_q}$ (i.e. no analog preprocessing), [24] established capacity bounds which are tight in some SNR regimes and under certain channel conditions.

While $C(\mathbf{A})$ in problem (3) is known to be a concave maximization problem, it is not clear whether C in problem (4) can be solved by some provably optimal method in the quantized setting. This problem formulation was considered in [41], which showed that C can be attained in the infinite SNR regime by choosing an \mathbf{A} that maximizes the partitions of the transmit signal space. Moreover, C in problem (4) can be upper bounded by the TSC = $\min\{C_{\text{MIMO-AWGN}}, N_q\}$ [37], [45]. Here, $C_{\text{MIMO-AWGN}}$ is the MIMO AWGN channel capacity without quantization and N_q is the number of 1-bit quantizers. The TSC bound serves as a *universal* upper bound for the capacity of *any* discrete-time memoryless MIMO channel with *any* N_q -bit quantization at the output. This bound, however, does not utilize the information about \mathbf{t} .

In the following section, we present an achievability scheme that frames (4) as a resource allocation problem of the transmit power and sign quantizers over all eigenchannels of \mathbf{H} . The resulting transmit power and sign quantizer allocation correspond to specific choices of $F_{\mathbf{x}}$ and \mathbf{A} , respectively, which are not necessarily optimal. Effectively, the rate of this scheme gives a lower bound on the solution of (4). We also show that the rate of this achievability scheme is tight in the low SNR regime. However, it does not attain the high SNR capacity when $\min\{N_t, N_r\} < 2N_q$.

III. ACHIEVABILITY SCHEME FOR THE HYBRID ONE-SHOT RECEIVER

In this section, we establish an achievable result for the capacity of the MIMO system depicted in Figure 1. The key idea in this achievability result is to configure \mathbf{A} such that the channel becomes a set of parallel SISO subchannels with phase quantization at the output. We then use the capacity-achieving input for this channel as our transmit strategy.

First, we apply the singular value decomposition (SVD) to the channel matrix \mathbf{H} to get the following matrix factorization:

$$\mathbf{H} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^H, \quad (5)$$

where $\mathbf{\Sigma} \in \mathbb{C}^{N_r \times N_t}$ is a diagonal matrix. The first $N_\sigma = \min\{N_t, N_r\}$ diagonal entries are the singular values $\{\sqrt{\lambda_i}\}_{i=1}^{N_\sigma}$ arranged in non-increasing order (that is $\lambda_i \geq \lambda_{i+1}$), and the remaining diagonal entries are zeros. The matrices $\mathbf{U} \in \mathbb{C}^{N_r \times N_r}$ and $\mathbf{V} \in \mathbb{C}^{N_t \times N_t}$ are unitary matrices. Suppose we define $\tilde{\mathbf{x}} = \mathbf{V}\mathbf{x}$ as the precoded transmitted symbols. Then, \mathbf{y} can be written as

$$\begin{aligned} \mathbf{y} &= \text{sign}(\Re\{\mathbf{A}(\mathbf{U}\mathbf{\Sigma}\mathbf{V}^H\tilde{\mathbf{x}}) + \mathbf{z}'\}) \\ &= \text{sign}(\Re\{\mathbf{A}(\mathbf{U}\mathbf{\Sigma}\mathbf{x}) + \mathbf{z}'\}). \end{aligned} \quad (6)$$

Without loss of generality, we can set the analog linear combiner \mathbf{A} as a product of two matrices $\mathbf{\Phi} \in \mathbb{C}^{N_q \times N_r}$ and $\mathbf{U}^H \in \mathbb{C}^{N_r \times N_r}$. Equation (6) then simplifies to

$$\mathbf{y} = \text{sign}(\Re\{\mathbf{\Phi}(\mathbf{\Sigma}\mathbf{x}) + \mathbf{z}'\}),$$

and the optimization problem (4) becomes

$$C = \max_{\mathbf{\Phi}} \max_{F_{\mathbf{x}}} I_{\mathbf{\Phi}}(\mathbf{x}; \mathbf{y}), \quad (7)$$

where we used the subscript $\mathbf{\Phi}$ to emphasize the dependence of the mutual information in the choice of $\mathbf{\Phi}$.

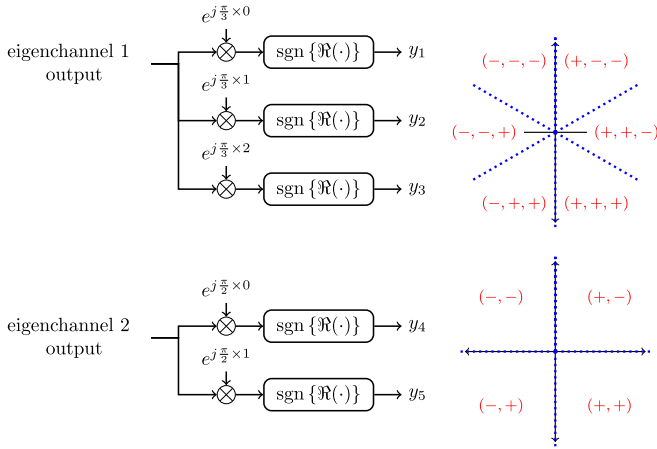


Fig. 2. Illustration of how Φ_{PH} is constructed for a MIMO channel with $N_\sigma = 2$, $N_q = 5$, and $\mathbf{s} = [3, 2]^T$.

A. Design Φ to Create Phase Quantizers

Prior to multiplying Φ , we have N_σ parallel complex AWGN channels. The matrix Φ can be used to create symmetric phase quantizers and connect these phase quantizers to each eigenchannel. First, we define an $N_r \times 1$ quantizer allocation vector $\mathbf{s} = [s_1, \dots, s_{N_\sigma}, \mathbf{0}_{1 \times N_r - N_\sigma}]^T$ such that $\sum_{i=1}^{N_\sigma} s_i = N_q$ and $s_i \geq 0$. The entries s_i correspond to the number of sign comparators used to discretize the output of the i -th channel. The set of s_i sign comparators creates a $2s_i$ -sector phase quantizer connected at the output of the i -th channel. An illustrative example is shown in Figure 2 where three (rotated) sign quantizers are connected at the output of eigenchannel 1 to realize a 6-sector phase-quantized channel and two (rotated) sign quantizers are connected at the output of eigenchannel 2 to realize a 4-sector phase-quantized channel. In this case, we have $\mathbf{s} = [3, 2]^T$ and Φ becomes

$$\Phi = [\phi_1, \phi_2], \quad \text{where} \quad \begin{cases} \phi_1 = [1, e^{j\frac{\pi}{3}}, e^{j\frac{2\pi}{3}}, 0, 0]^T \\ \phi_2 = [0, 0, 0, 1, e^{j\frac{\pi}{2}}]^T \end{cases}.$$

More generally, the Φ matrix is a horizontal stacking of $\phi_k \in \mathbb{C}^{N_q \times 1}$ for $k = 1, \dots, N_r$. The first N_σ column vectors can be expressed as

$$\phi_k = [\mathbf{0}_{1 \times \sum_{i=0}^{k-1} s_i} \quad \mathbf{e}_{1 \times s_k} \quad \mathbf{0}_{1 \times \sum_{i=k+1}^{N_q} s_i}]^T \quad (8)$$

for all $k \in \{1, \dots, N_\sigma\}$, where the l -th entry of $\mathbf{e}_{s_k \times 1}$, denoted as $\mathbf{e}_{s_k \times 1}^{(l)}$, is $e^{\frac{j\pi(l-1)}{s_k}}$. The remaining $N_r - N_\sigma$ column vectors of Φ are all-zero vectors. Effectively, we define $\mathcal{Q}_K^{\text{PH}}(\cdot) : \mathbb{C} \mapsto \mathbb{Z}$ as a function that performs a K -sector symmetric phase quantization to map a complex value to an integer value between 0 to $K-1$. The output of the i -th phase-quantized eigenchannel is then

$$y_q^{(i)} = \mathcal{Q}_{2s_i}^{\text{PH}} \left(\sqrt{\lambda_i} \mathbf{x}^{(i)} + \mathbf{z}^{(i)} \right), \quad (9)$$

where $\mathbf{x}^{(i)}$ and $\mathbf{z}^{(i)}$ are the i -th entry of \mathbf{x} and \mathbf{z} , respectively. We shall call this choice of Φ as Φ_{PH} and the corresponding analog combiner as $\mathbf{A}_{\text{PH}} = \Phi_{\text{PH}} \mathbf{U}^H$. Note that

$$C = \max_{\mathbf{A}} C(\mathbf{A}) \geq C(\mathbf{A}_{\text{PH}}) \quad (10)$$

with equality if \mathbf{A}_{PH} is the analog combiner configuration that maximizes $C(\mathbf{A})$. Thus, the quantity $C(\mathbf{A}_{\text{PH}})$ is the rate of our achievability scheme and is a lower bound for problem (4).

To solve $C(\mathbf{A}_{\text{PH}})$, we present two function definitions that correspond to the conditional probability and conditional entropy of a Gaussian channel with K -sector phase quantization at the output.

Definition 1: The phase quantization probability function, $W_y^{(K)}(\nu, \theta)$, is defined as

$$W_y^{(K)}(\nu, \theta) = \int_{\frac{2\pi}{K}y - \pi - \theta}^{\frac{2\pi}{K}(y+1) - \pi - \theta} f_{\Phi|N}(\phi|\nu) d\phi, \quad (11)$$

where $f_{\Phi|N}(\phi|\nu)$ is

$$= \frac{e^{-\nu}}{2\pi} + \frac{\sqrt{\nu} \cos(\phi) e^{-\nu \sin^2(\phi)} [1 - Q(\sqrt{2\nu} \cos(\phi))]}{\sqrt{\pi}}, \quad (12)$$

and $Q(x)$ is the Gaussian Q -function, $\theta \in [-\pi, \pi]$, and $\nu \geq 0$.

Definition 2: The phase quantization entropy, $w_K(\nu, \theta)$, is defined as¹

$$w_K(\nu, \theta) = - \sum_{y=0}^{K-1} W_y^{(K)}(\nu, \theta) \log W_y^{(K)}(\nu, \theta) \quad (13)$$

for any $\theta \in [-\pi, \pi]$, $\nu \geq 0$, and $K > 1$. The function $W_y^{(K)}(\nu, \theta)$ is given in Definition 1.

From [21, Theorem 1], the capacity-achieving input of a Gaussian channel with K -sector phase quantization at the output is a rotated K -phase shift keying (PSK) and the capacity can be computed numerically. Consequently, we formulate $C(\mathbf{A}_{\text{PH}})$ as

$$C(\mathbf{A}_{\text{PH}}) = \max_{s_i, \rho_i} \sum_{i=1}^{N_\sigma} C_\phi \left(s_i, \frac{\lambda_i \rho_i}{\sigma^2} \right) \quad (14a)$$

$$\text{s.t.} \quad \sum_{i=1}^{N_\sigma} s_i = N_q \quad (14b)$$

$$\sum_{i=1}^{N_\sigma} \rho_i \leq P \quad (14c)$$

$$\rho_i \geq 0, s_i \in \{0, 1, \dots, N_q\}, \quad (14d)$$

where $C_\phi(s, \nu)$ is the capacity of a scalar Gaussian channel with SNR = ν and $2s$ -sector phase quantization at the output. Mathematically,

$$C_\phi(s, \nu) = \begin{cases} \log 2 s - w_{2s} \left(\nu, \frac{\pi}{2s} \right), & s \geq 1 \\ 0, & \text{otherwise} \end{cases} \quad (15)$$

The objective function of (14) is the sum capacity of the N_σ parallel eigenchannels with phase quantization at the output. A $2s_i$ -PSK with amplitude $\sqrt{\rho_i}$ is transmitted over the i -th eigenchannel to attain the capacity. The optimization of $\{s_i\}_{i=1}^{N_\sigma}$ corresponds to the optimization of Φ_{PH} (consequently, \mathbf{A}_{PH}) whereas the optimization of $\{\rho_i\}_{i=1}^{N_\sigma}$ gives the optimal $F_{\mathbf{X}}$ that attains $C(\mathbf{A}_{\text{PH}})$.

¹All $\log(\cdot)$ terms in this paper are in base 2 unless specified otherwise.

B. Alternating Optimization Approach

One major issue with optimization problem (14) is its nonconvex structure due to the discrete parameters in the search space. In this subsection, we present a polynomial-time heuristic method that maximizes the objective function in (14). The approach is based on alternating optimization of the parameters.

First, we define the parameter N_s to be the number of active² eigenchannels that we intend to use for transmission. Note that the N_s active eigenchannels in the achievability scheme correspond to the eigenchannels having the N_s strongest singular values. The optimality of this choice is formalized in the following lemma.

Lemma 1: If the optimal strategy has N_s active eigenchannels, then the eigenvalues of those eigenchannels should be $\{\lambda_i\}_{i=1}^{N_s}$. In other words, these channels should have the strongest eigenvalues among N_σ eigenchannels.

Proof: See Appendix A. \square

By Lemma 1, we can rewrite problem (14) without loss of optimality as

$$C(\mathbf{A}_{\text{PH}}) = \max_{s_i, \rho_i, N_s} \sum_{i=1}^{N_s} \log 2 s_i - w_{2s_i} \left(\frac{\lambda_i \rho_i}{\sigma^2}, \frac{\pi}{2s_i} \right) \quad (16a)$$

$$\text{s.t.} \quad \sum_{i=1}^{N_s} s_i = N_q \quad (16b)$$

$$\sum_{i=1}^{N_s} \rho_i \leq P \quad (16c)$$

$$\rho_i \geq 0, s_i \in \{0, \dots, N_q\} \quad (16d)$$

$$N_s \in \{1, \dots, N_\sigma\}. \quad (16e)$$

Next, note that for fixed N_s and $\{s_i\}_{i=1}^{N_s}$, the optimization problem (16) can be simplified to

$$\min_{\rho_i} \sum_{i=1}^{N_s} w_{2s_i} \left(\frac{\lambda_i \rho_i}{\sigma^2}, \frac{\pi}{2s_i} \right) \quad (17a)$$

$$\text{s.t.} \quad \sum_{i=1}^{N_s} \rho_i \leq P \quad (17b)$$

$$\rho_i \geq 0 \quad \forall i \in \{1, 2, \dots, N_s\}, \quad (17c)$$

which has a convex structure. This is because $w_{2s_i}(\nu, \theta)$ is convex on ν [21, Proposition 2] and the summation of non-negative convex functions preserves convexity. Consequently, the optimal power allocation, denoted $\{\rho'_i\}_{i=1}^{N_s}$, can be solved efficiently using standard convex solvers.

For fixed N_s and $\{\rho'_i\}_{i=1}^{N_s}$, problem (16) becomes a discrete optimization over the parameters $\{s_i\}_{i=1}^{N_s}$. This optimal

allocation of the sign quantizers, denoted $\{s'_i\}_{i=1}^{N_s}$, can be solved using a dynamic programming approach. We define a state space $\mathcal{S}_{\text{state}}$ with each state being the tuple (i, n_q) , where $i \in \{0, \dots, N_s\}$ and $n_q \in \{0, \dots, N_q\}$. Define also the function $f(i, n_q)$ as the sum capacity of channels 1 to i when there are n_q sign quantizers that can be allocated to these i channels. This function $f(i, n_q)$ can be expressed using the recurrence relation in (18), as shown at the bottom of the page. The value of $f(N_s, N_q)$ gives the sum capacity for N_s active eigenchannels with N_q sign quantizers. We can also compute the optimum choice of k per state (i, n_q) as

$$s(i, n_q) = \arg \max_{k \in \{1, \dots, n_q\}} \left\{ f(i-1, n_q - k) + C_\phi \left(k, \frac{\lambda_i \rho_i}{\sigma^2} \right) \right\}.$$

The proof that the dynamic programming approach solves the optimal quantizer allocation for a fixed N_s and fixed power allocation is presented in Appendix B. We simply iterate over all $N_s \in \{1, \dots, N_\sigma\}$ and then for each value, alternate between the two optimization procedures until convergence.

The remaining computational bottleneck in solving the optimization problem is the evaluation of $C_\phi(s, \nu)$ in (15). This is because the phase quantization entropy contains integral terms that need to be computed numerically. To this end, we define a function $\tilde{C}_\phi(s, \nu)$ which closely approximates $C_\phi(s, \nu)$ as follows:

$$\tilde{C}_\phi(s, \nu) = \begin{cases} \log 2 s - \tilde{w}_{2s} \left(\nu, \frac{\pi}{2s} \right), & s \geq 1 \\ 0, & \text{otherwise} \end{cases}, \quad (19)$$

where

$$\tilde{w}_{2s} \left(\nu, \frac{\pi}{2s} \right) = - \sum_{y=0}^{2s-1} \tilde{W}_y^{(2s)} \left(\nu, \frac{\pi}{2s} \right) \log \tilde{W}_y^{(2s)} \left(\nu, \frac{\pi}{2s} \right), \quad (20)$$

$$\begin{aligned} & \tilde{W}_y^{(2s)} \left(\nu, \frac{\pi}{2s} \right) \\ &= \frac{1}{G} \cdot \frac{\pi}{sR} \sum_{r=0}^{R-1} f_{\Phi|N} \left(\frac{\pi}{s} \left(y + \frac{r}{R} \right) - \pi - \theta \middle| \nu \right), \end{aligned} \quad (21)$$

and $G = \sum_{y=0}^{2s-1} \tilde{W}_y^{(2s)} \left(\nu, \frac{\pi}{2s} \right)$. We approximate the integral terms using midpoint rule of definite integrals. R corresponds to the number of rectangles to be used in the approximation and $\frac{\pi}{sR}$ is the width of a rectangle. The factor $1/G$ ensures that the set $\left\{ \tilde{W}_y^{(2s)} \left(\nu, \frac{\pi}{2s} \right) \right\}_{y=0}^{2s-1}$ forms a probability simplex.

The plots of $\tilde{C}_\phi(s, \nu)$ and $C_\phi(s, \nu)$ for different s , and the squared approximation error are plotted in Figures 3a and 3b, respectively. Here, we fix $R = 9$ and observe different values of s . At this choice of R , we can see small approximation errors for all ν and s considered. Throughout this paper, the setting $R = 9$ is used for evaluating $\tilde{C}_\phi(s, \nu)$.

$$f(i, n_q) = \begin{cases} 0 & , i = 0 \text{ or } n_q = 0 \\ \max_{k \in \{1, \dots, n_q\}} \left\{ f(i-1, n_q - k) + \log 2k - w_{2k} \left(\frac{\lambda_i \rho_i}{\sigma^2}, \frac{\pi}{2k} \right) \right\} & , \text{ otherwise} \end{cases} \quad (18)$$

²An eigenchannel is active if it has a nonzero transmit power and a nonzero quantizer allocation. Otherwise, it is inactive.

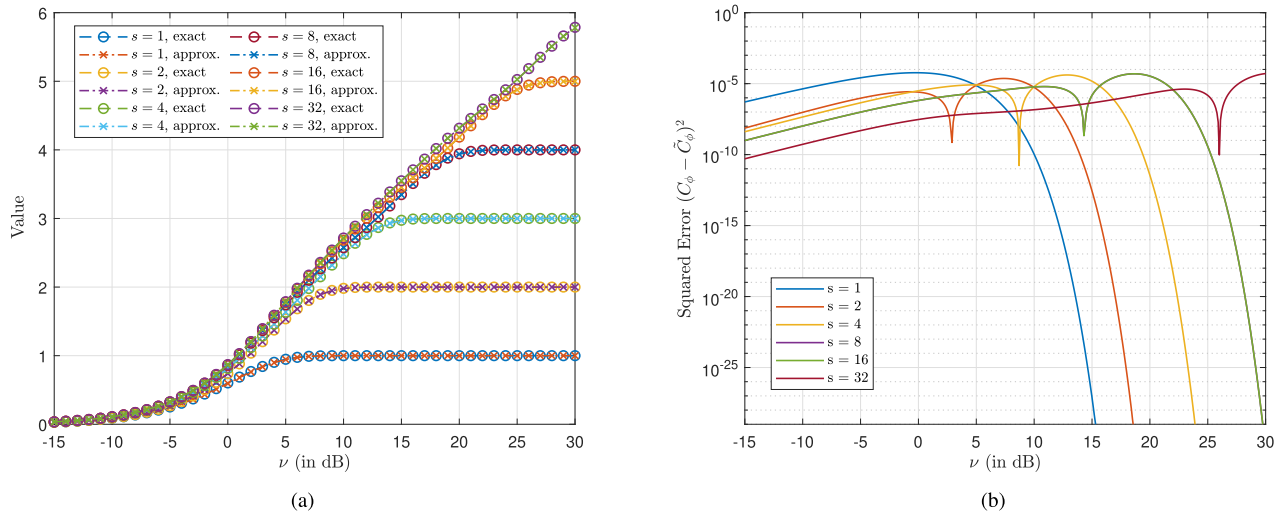


Fig. 3. (a) Plots of $C_\phi(s, \nu)$ and $\tilde{C}_\phi(s, \nu)$ for different s ; and (b) Squared approximation error of $C_\phi(s, SNR)$ and $\tilde{C}_\phi(s, \nu)$ vs ν for different values of s . We set $R = 9$.

Algorithm 1 outlines the alternating optimization approach. Here, the power allocation and quantizer allocation steps are referred to as optimization procedure #1 and #2, respectively. We use $\tilde{C}_\phi(s, \nu)$ instead of $C_\phi(s, \nu)$ in these optimization procedures. The optimized $\{\rho_i\}_{i=1}^{N_s}$, $\{s_i\}_{i=1}^{N_s}$, and N_s are then used to construct the approximate solution for the analog linear combiner \mathbf{A}_{PH} and the distribution $F_{\mathbf{X}}$ that attains $C(\mathbf{A}_{\text{PH}})$. We shall denote these approximate solutions produced by Algorithm 1 as $\hat{\mathbf{A}}_{\text{PH}}$ and $\hat{F}_{\mathbf{X}}$. While (19) is used to optimize \mathbf{A}_{PH} and $F_{\mathbf{X}}$, we still use (15) to evaluate C^* in Line 15. The computational complexity of computing $\hat{\mathbf{A}}_{\text{PH}}$ and $\hat{F}_{\mathbf{X}}$ using Algorithm 1 is

$$O\left(N_\sigma \{T_1 + T_2\} \log \frac{1}{\epsilon_1}\right), \quad (22)$$

where

$$\begin{aligned} T_1 &= O\left(\max\{N_\sigma^3 T_{\text{cap}}^3, \mathcal{F}\} \log \frac{N_\sigma}{\epsilon_2}\right) \\ T_2 &= O(N_\sigma N_q^2 T_{\text{cap}}) \\ T_{\text{cap}} &= O(N_q R). \end{aligned}$$

The quantities T_1 and T_2 account for the computational complexity of optimization procedures #1 and #2, respectively. For optimization procedure #1, interior point method is used to solve the convex power allocation problem. The expression $\max\{N_\sigma^3 T_{\text{cap}}^3, \mathcal{F}\}$ in T_1 is the number of operations per iteration [46, Section 1.3.1] and $\log\left(\frac{N_\sigma}{\epsilon_2}\right)$ accounts for the number of iterations [46, Section 11.3.3]. The parameter \mathcal{F} is the cost of evaluating the first and second derivatives of the objective function and constraints. The parameters ϵ_1 and ϵ_2 set the convergence criterion for the alternating optimization scheme and the optimization procedure #1, respectively. Finally, T_{cap} is the computational complexity of evaluating $\tilde{C}_\phi(s, \nu)$.

We note that Algorithm 1 only guarantees convergence to a local optimal solution of problem (14). This is because $C(\mathbf{A}_{\text{PH}})$ is not jointly convex over the parameters $\{\rho_i\}$ and $\{s_i\}$. Furthermore, $C(\mathbf{A}_{\text{PH}})$ is just a lower bound on C . Thus,

we have

$$C(\hat{\mathbf{A}}_{\text{PH}}) \leq C(\mathbf{A}_{\text{PH}}) \leq C.$$

It is therefore necessary to establish an upper bound on (4) to gauge how far, at worst, is $C(\hat{\mathbf{A}}_{\text{PH}})$ to the exact capacity C .

IV. CAPACITY UPPER BOUND FOR THE HYBRID ONE-SHOT RECEIVER

A. Capacity Upper Bound 1

We now derive an upper bound for the capacity expression in (4). First, we use the high SNR capacity results in [41] and apply it to \mathbb{C}^{N_σ} to establish C in the infinite SNR regime, which we denote as C_∞ . That is, we maximize the number of regions created by arranging the hyperplanes in N_σ complex dimensions. By doing this, we maximize the output entropy. For the special case where all the hyperplanes intersect the origin, [47] proved that the number of regions created by N_q hyperplanes in $2N_\sigma$ real dimensions is

$$\mathcal{M}_\infty = 2 \sum_{i=0}^{2N_\sigma-1} \binom{N_q-1}{i}, \quad (23)$$

where $\binom{n}{m} = 0$ if $m > n$. Consequently, $C_\infty = \log \mathcal{M}_\infty$ by choosing an $F_{\mathbf{X}}$ that induces a uniform output distribution.

B. Capacity Upper Bound 2

The capacity upper bound established in the previous subsection is quite loose in the finite SNR case. To establish a tighter capacity upper bound in the finite SNR regime, we note that the sign quantizer only requires the phase information of its input in order to output -1 or $+1$. Suppose we define a function $\Theta_k(\cdot)$ that extracts the phase information of the k -th output of the analog combiner, and another function that maps the phase detector output to $+1$ when it is inside the region $[-\frac{\pi}{2}, \frac{\pi}{2}]$ and is mapped to -1 otherwise. Then, the system model in Figure 4 is the same as the system model depicted

Algorithm 1 Alternating Optimization to Get Local Optimal Solution to Problem (14)

Input: $N_\sigma, \{\lambda_i\}_{i=1}^{N_\sigma}, \sigma^2, N_q, P$
Output: $\hat{\mathbf{A}}_{\text{PH}}, \hat{F}_{\mathbf{X}}, C^*$

```

1  $C^* = 0, N_s^* = 0, N_s = 1$  // Initialize
2 Set  $\epsilon_1$  // Set convergence condition
3 for  $N_s = 1$  to  $N_\sigma$  do
4   /* Uniform allocation of  $s_i$  */
5    $s_i = \lfloor \frac{N_q}{N_s} \rfloor \forall i \in \{1, \dots, N_q \bmod N_s\}$ ,
6    $s_i = \lfloor \frac{N_q}{N_s} \rfloor \forall i \in \{N_q \bmod N_s + 1, \dots, N_s\}$ 
7   repeat
8     /* Optimization Procedure #1 (optimize power allocation) */
9     Solve  $\{\rho_i^*\}_{i=1}^{N_s}$  using the formulation in (17).
10     $C_1 = \sum_{i=1}^{N_s} \tilde{C}_\phi(s_i, \frac{\lambda_i \rho_i^*}{\sigma^2})$ 
11     $\rho_i = \rho_i^* \forall i \in \{1, \dots, N_s\}$ 
12    /* Optimization Procedure #2 (optimize quantizer allocation) */
13    Solve  $\{s_i^*\}_{i=1}^{N_s}$  using dynamic programming.
14     $C_2 = \sum_{i=1}^{N_s} \tilde{C}_\phi(s_i^*, \frac{\lambda_i \rho_i^*}{\sigma^2})$ 
15     $s_i = s_i^* \forall i \in \{1, \dots, N_s\}$ 
16  until  $C_2 - C_1 < \epsilon_1$ 
17  if  $C_2 > C^*$  and  $s_i > 0 \forall i \in \{1, \dots, N_s\}$  // Update optimal values
18  then
19     $C^* = \sum_{i=1}^{N_s} \log 2 s_i - w_{2s_i}(\frac{\lambda_i \rho_i^*}{\sigma^2}, \frac{\pi}{2s_i})$ 
20     $N_s^* = N_s, s'_i = s_i, \rho'_i = \rho_i \forall i \in \{1, \dots, N_s^*\}$ 
21   $\Phi_{\text{PH}} = [\phi_1, \dots, \phi_{N_s^*}, \dots, 0]$ , where  $\phi_k$  is from (8).
22   $\hat{\mathbf{A}}_{\text{PH}} = \Phi_{\text{PH}} \mathbf{U}^H$ 
23   $\hat{F}_{\mathbf{X}_i} = 2s'_i$ -PSK with amplitude  $\sqrt{\rho'_i}$ 

```

in Figure 1. It then follows that the capacity of the two system models are equal. By the DPI, we have

$$I_{\mathbf{A}}(\mathbf{x}; \mathbf{y}) \leq I_{\mathbf{A}}(\mathbf{x}; \tilde{\mathbf{y}}),$$

and an upper bound on the channel capacity can be established by maximizing the mutual information between the input \mathbf{x} and the outputs of the phase detectors. Define $\Theta_{i:j}$ to be the phase detector outputs from index i to index j , and $\mathbf{z}'_{i:j}$ to be the noise components at the output of the analog combiner from index i to index j . The following lemma further simplifies the capacity maximization problem.

Lemma 2: Suppose we have $N_\sigma \leq N_q$. Suppose further that there exists an analog linear combiner $\mathbf{A} = \begin{bmatrix} \mathbf{A}_1 \\ \mathbf{A}_2 \end{bmatrix}$ such that $\mathbf{A}\mathbf{H} = \begin{bmatrix} \mathbf{A}_1\mathbf{H} \\ \mathbf{A}_2\mathbf{H} \end{bmatrix} = \begin{bmatrix} \mathbf{B}_1 \\ \mathbf{B}_2 \end{bmatrix}$, where $\mathbf{B}_1 \in \mathbb{C}^{N_\sigma \times N_t}$ and $\mathbf{B}_2 \in \mathbb{C}^{(N_q - N_\sigma) \times N_t}$. If \mathbf{B}_1 is full rank and $\tilde{\mathbf{y}} = [\tilde{\mathbf{y}}^{(1)} \tilde{\mathbf{y}}^{(2)}]^H$, where $\tilde{\mathbf{y}}^{(1)} = \Theta_{1:N_\sigma}(\mathbf{B}_1\mathbf{x} + \mathbf{z}'_{1:N_\sigma})$ and $\tilde{\mathbf{y}}^{(2)} = \Theta_{N_\sigma+1:N_q}(\mathbf{B}_2\mathbf{x} + \mathbf{z}'_{N_\sigma+1:N_q})$, then

$$I_{\mathbf{A}}(\mathbf{x}; \tilde{\mathbf{y}}) = I_{\mathbf{A}_1}(\mathbf{x}; \tilde{\mathbf{y}}^{(1)}).$$

Moreover, if $I_{\mathbf{A}}(\mathbf{x}; \tilde{\mathbf{y}})$ is achieved by a \mathbf{B}_1 that is not full rank, then there exists an $\mathbf{A}'\mathbf{H} = \begin{bmatrix} \mathbf{B}'_1 \\ \mathbf{B}'_2 \end{bmatrix}$ (where \mathbf{B}'_1 being full rank) and a distribution $F_{\mathbf{X}'}$ such that $I_{\mathbf{A}}(\mathbf{x}; \tilde{\mathbf{y}}) = I_{\mathbf{A}'}(\mathbf{x}'; \tilde{\mathbf{y}})$.

Proof: See Appendix C. \square

To put it simply, Lemma 2 shows that we only need to consider N_σ phase detector outputs since considering more phase detector outputs than N_σ does not increase the mutual information. Thus, without loss of generality, we can consider the maximization of $I_{\mathbf{A}_1}(\mathbf{x}; \tilde{\mathbf{y}}^{(1)})$ with $\mathbf{B}_1 = \mathbf{A}_1\mathbf{H}$ being full rank. In the remainder, we will show that C in problem (4) can be upper bounded by the capacity of a Gaussian MIMO channel with phase detectors at the output. The capacity-achieving input of this channel is also characterized.

Suppose we denote $\tilde{y}_i^{(1)}$ as the i -th element of $\tilde{\mathbf{y}}^{(1)}$, $\tilde{\mathbf{y}}_{i:j}^{(1)}$ as the elements of $\tilde{\mathbf{y}}^{(1)}$ from index i to index j , and $h_{\mathbf{A}_1}(\cdot)$ as the differential entropy function induced by \mathbf{A}_1 . Then, we get the following upper bound on $\max_{\mathbf{A}_1} \max_{F_{\mathbf{X}}} I_{\mathbf{A}_1}(\mathbf{x}; \tilde{\mathbf{y}}^{(1)})$:

$$\begin{aligned}
&= \max_{\mathbf{A}_1} \max_{F_{\mathbf{X}}} h_{\mathbf{A}_1}(\tilde{\mathbf{y}}^{(1)}) - h_{\mathbf{A}_1}(\tilde{\mathbf{y}}^{(1)}, \mathbf{x}) \\
&= \max_{\mathbf{A}_1} \max_{F_{\mathbf{X}}} \sum_{i=1}^{N_\sigma} h_{\mathbf{A}_1}(\tilde{y}_i^{(1)} | \tilde{\mathbf{y}}_{1:i-1}^{(1)}) \\
&\quad - \sum_{i=1}^{N_\sigma} h_{\mathbf{A}_1}(\tilde{y}_i^{(1)} | \tilde{\mathbf{y}}_{1:i-1}^{(1)}, \mathbf{x}) \\
&\stackrel{(a)}{\leq} \max_{\mathbf{A}_1} \max_{F_{\mathbf{X}}} \sum_{i=1}^{N_\sigma} h_{\mathbf{A}_1}(\tilde{y}_i^{(1)}) - \sum_{i=1}^{N_\sigma} h_{\mathbf{A}_1}(\tilde{y}_i^{(1)} | \tilde{\mathbf{y}}_{1:i-1}^{(1)}, \mathbf{x}) \\
&\stackrel{(b)}{\leq} \max_{\mathbf{A}_1} \max_{F_{\mathbf{X}}} \sum_{i=1}^{N_\sigma} h_{\mathbf{A}_1}(\tilde{y}_i^{(1)}) \\
&\quad - \sum_{i=1}^{N_\sigma} h_{\mathbf{A}_1}(\tilde{y}_i^{(1)} | \tilde{\mathbf{y}}_{1:i-1}^{(1)}, \tilde{\mathbf{y}}_{i+1:N_\sigma}^{(1)}, \mathbf{x}) \\
&= \max_{\tilde{\mathbf{A}}_1 \text{ s.t. } \mathbf{B}_1 \text{ has orthogonal rows}} \max_{F_{\mathbf{X}}} \sum_{i=1}^{N_\sigma} \left[h_{\tilde{\mathbf{A}}_1}(\tilde{y}_i^{(1)}) \right. \\
&\quad \left. - h_{\tilde{\mathbf{A}}_1}(\tilde{y}_i^{(1)} | \mathbf{x}) \right] \\
&= \max_{\tilde{\mathbf{A}}_1 \text{ s.t. } \mathbf{B}_1 \text{ has orthogonal rows}} \max_{F_{\mathbf{X}}} \sum_{i=1}^{N_\sigma} I_{\tilde{\mathbf{A}}_1}(\mathbf{x}; \tilde{y}_i^{(1)}) \\
&= \max_{F_{\mathbf{X}}} \sum_{i=1}^{N_\sigma} I_{\mathbf{U}_1^H}(\mathbf{x}; \tilde{y}_i^{(1)}).
\end{aligned}$$

The first two lines follow from the definition of mutual information and the use of the chain rule. The inequalities (a) and (b) follow from the fact that conditioning reduces entropy. Note that equality in both (a) and (b) can be achieved if and only if $\tilde{y}_i^{(1)}$ is independent of $\tilde{\mathbf{y}}_{1:i-1}^{(1)}$ and $\tilde{\mathbf{y}}_{i+1:N_\sigma}^{(1)}$. Since we are able to choose \mathbf{A}_1 , an appropriate choice of \mathbf{A}_1 that achieves equality in the third and fourth line should make \mathbf{B}_1 have mutually orthogonal rows. Existence of such \mathbf{B}_1 is guaranteed since $N_\sigma \leq N_t$. Doing so gives us the fifth line. We call this choice $\tilde{\mathbf{A}}_1$. The sixth line follows from the definition of mutual information. Finally, we obtain the last line using the fact that, out of all the orthonormal bases for the range of \mathbf{H} , SVD produces the orthonormal basis for which the total channel gain along each direction is maximized.³ Thus, $\tilde{\mathbf{A}}_1$ should be \mathbf{U}_1^H , where \mathbf{U}_1 contains the columns of \mathbf{U}

³Recall the Maximum Variance Formulation of Principal Component Analysis (PCA) and its connection to SVD.

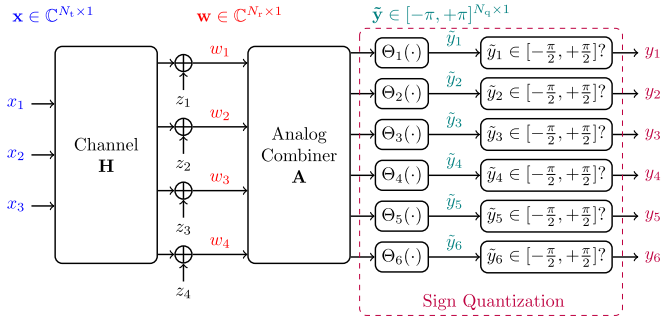


Fig. 4. Modified System Model with Sign Quantization broken down to two-stages: (1) Phase Detection and (2) Phase-to-bit.

corresponding to the N_σ nonzero eigenvalues. Consequently, $\mathbf{B}_1 = \Sigma_{N_\sigma}$ (the first N_σ rows of Σ).

The above result shows that $I_{\mathbf{A}_1}(\mathbf{x}; \tilde{\mathbf{y}}^{(1)})$ can be upper bounded by the input-output mutual information of a Gaussian product channel with N_σ eigenchannels; each eigenchannel has a phase detector at the output. Using this, we can now establish the capacity upper bound in the finite SNR.

Proposition 1: The solution to (4) can be upper bounded by (24), as shown at the bottom of the page.

Proof: By considering the limiting case of [21, Theorem 1] with $b \rightarrow \infty$, the capacity-achieving input of the i -th Gaussian channel with phase detector at the output is an ∞ -PSK (i.e. a circle with radius $\sqrt{\rho_i}$) for all SNR. Optimal power allocation is applied to the eigenchannels to get the desired capacity upper bound result. \square

The first and second terms of (24) are the differential output entropy and conditional differential entropy, respectively. Problem (24) is a convex optimization problem. The convexity of the objective function in (24) comes from [21, Proposition 2]. Finally, by combining (24) and C_∞ , we get the upper bound

$$C_{\text{UB}} = \min \{C_{\phi\text{-detector}}, C_\infty\}. \quad (25)$$

For the zero threshold ADC case, this capacity upper bound is tighter than the TSC bound established in [37] since it takes into account the fact that sign quantization throws away the amplitude information. Currently, the analysis framework we developed only applies to the $\mathbf{t} = \mathbf{0}$ case. Extension of the analysis framework to general \mathbf{t} is a potential subject for future work and we are exploring proof techniques that can be used to extend this current framework to general \mathbf{t} . In the next section, we compare (25) and the TSC bound.

V. NUMERICAL EVALUATION OF THE ACHIEVABILITY SCHEME AND UPPER BOUND

In this section, we investigate the performance of our achievability scheme and how close it is to our established

capacity upper bound in (25). We consider a fixed 4×4 channel \mathbf{H} with eigenvalues $\lambda_1 = 1.6, \lambda_2 = 1.2, \lambda_3 = 0.8$, and $\lambda_4 = 0.4$. We also set $P = 1$ and vary the SNR by changing the noise variance σ^2 . We fix the number of sign quantizers to $N_q = 12$. Even though a 4×4 MIMO setup is considered to generate the numerical results, we note that the insights obtained in this small setup are applicable to larger MIMO settings.

The achievable rate of Algorithm 1, denoted $C(\hat{\mathbf{A}}_{\text{PH}})$, is depicted in Figure 5a. The individual rates of each eigenchannel are also given in Figure 5a to see how the rate at each eigenchannel changes with SNR. For comparison, we also superimpose the capacity upper bound given in (25) and the TSC bound. Note that the gap between the TSC bound and our capacity upper bound corresponds to the rate loss incurred when the amplitude information of the received signal is thrown away in a MIMO Gaussian channel. The values of $\{\rho'_i\}_{i=1}^{N_s}$ and $\{s'_i\}_{i=1}^{N_s}$ computed by Algorithm 1 are shown in Figures 6a and 6b, respectively. These parameters are used to construct $\hat{\mathbf{A}}_{\text{PH}}$ and $\hat{F}_{\mathbf{X}}$ according to Lines 17-19 of Algorithm 1.

It can be observed that $C(\hat{\mathbf{A}}_{\text{PH}})$ is tight in the low SNR regime and outperforms the naive approach of simply setting the ρ_i 's and s_i 's to be equal. In this regime, the optimal strategy is for the transmitter to use $2N_q$ -PSK signaling and send the symbols over the strongest eigenchannel. Simultaneously, the receiver configures \mathbf{A}_{PH} in such a way that all sign quantizers are connected to the output of the strongest eigenchannel. With $N_q \rightarrow \infty$, the transmission strategy produced by the achievability scheme is an ∞ -PSK sent over the strongest eigenchannel. Meanwhile, the analog linear combiner is configured to form an ∞ -bit phase quantizer (or a phase detector). Thus, the gap between the achievable rate and capacity upper bound in the low SNR regime vanishes as N_q grows unbounded. There also exists SNR thresholds, above which we activate the strongest inactive eigenchannel for transmission. The individual rates in Figure 5a have nonmonotonic behavior and sharp transitions within the SNR range considered. This can be attributed to the discrete nature of optimizing $\{s_i\}$. Nonetheless, the achievable rate remains smooth.

At this point, one might expect that uniform allocation would have the same performance as Algorithm 1 in all SNR regimes if $\mathbf{H} = \mathbf{I}_{N_\sigma \times N_\sigma}$ (i.e. eigenvalues are equal). In the classical waterfilling scheme for Gaussian channels, there is no loss of optimality if power is uniformly allocated among subchannels with identical eigenvalues. However, as depicted in Figure 5a, there is a gap between $C(\mathbf{A}_{\text{PH}})$ and the achievable rate of the hybrid one-shot receiver under uniform allocation of $\{\rho_i\}$ and $\{s_i\}$ in the low SNR regime. Thus, there is still some benefit, albeit small, in optimizing $\{\rho_i\}$ and $\{s_i\}$ when $\mathbf{H} = \mathbf{I}_{N_\sigma \times N_\sigma}$.

$$C_{\phi\text{-detector}} = N_\sigma \log(2\pi) - \min_{\rho_i: \sum_i \rho_i = P} \left\{ \sum_{i=1}^{N_\sigma} \int_{-\pi}^{\pi} f_{\Phi|N} \left(\phi \left| \frac{\lambda_i \rho_i}{\sigma^2} \right. \right) \log \frac{1}{f_{\Phi|N} \left(\phi \left| \frac{\lambda_i \rho_i}{\sigma^2} \right. \right)} d\phi \right\} \quad (24)$$

where $f_{\Phi|N}(\phi|\nu)$ is given in (12)

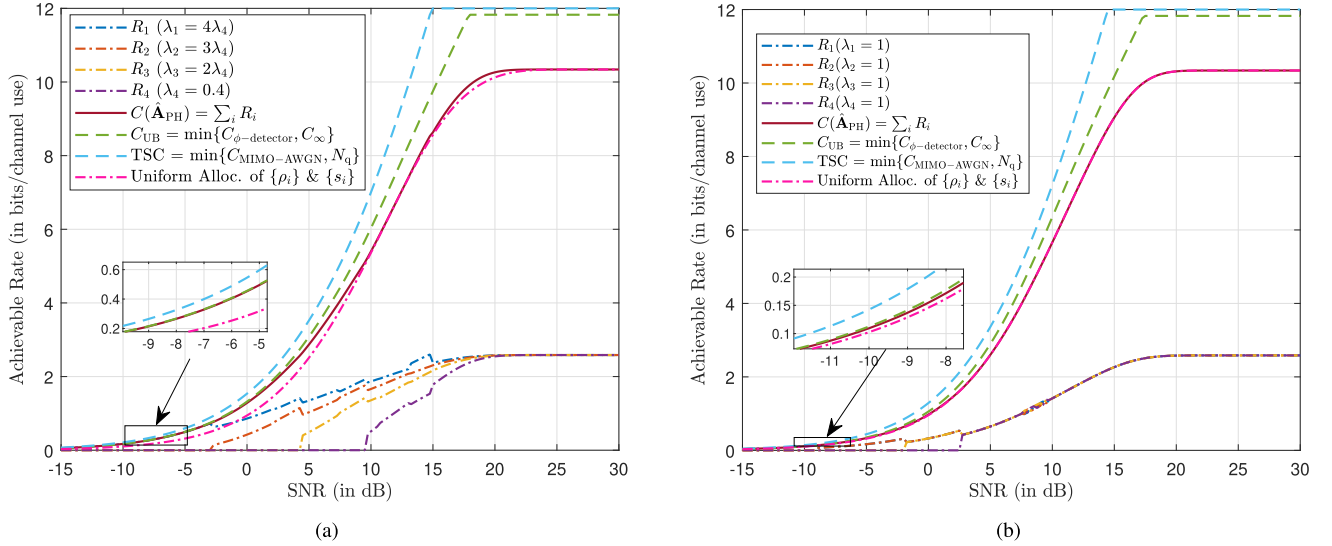


Fig. 5. Rate vs. SNR of the Achievability Scheme when (a) $\{\lambda_i\}_{i=1}^4 = \{1.6, 1.2, 0.8, 0.4\}$ and when (b) $\{\lambda_i\}_{i=1}^4 = \{1.0, 1.0, 1.0, 1.0\}$. Also superimposed are the individual rates, TSC bound, proposed upper bound, and the achievable rate under uniform allocation of $\{\rho_i\}$ and $\{s_i\}$.

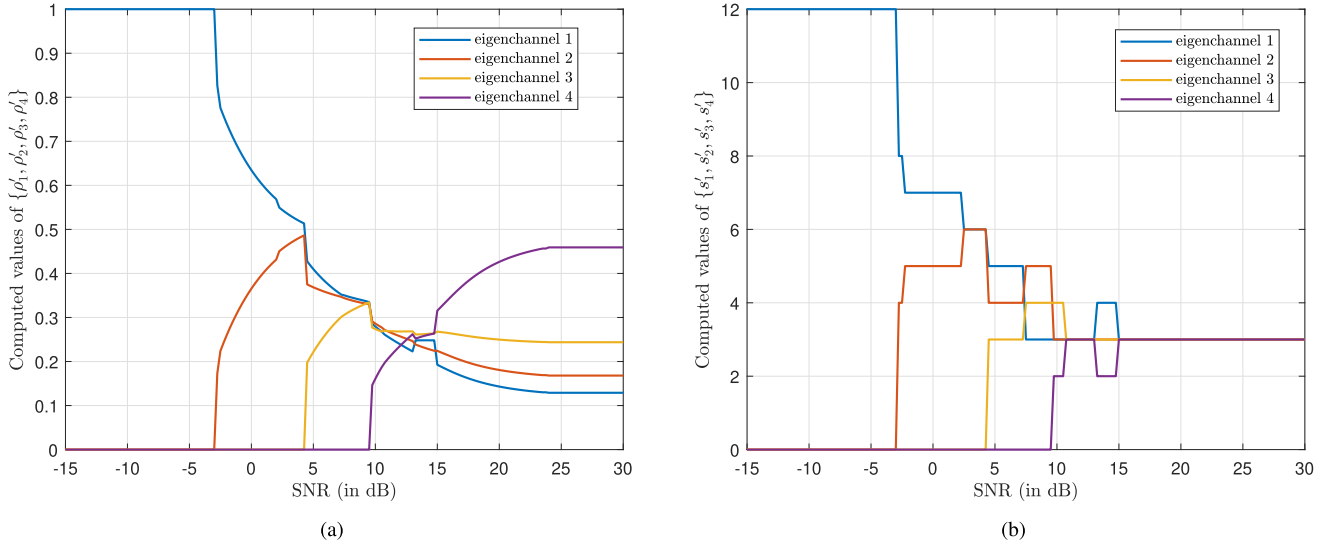


Fig. 6. Computed values of (a) $\{\rho'_i\}_{i=1}^4$ and (b) $\{s'_i\}_{i=1}^4$ as a function of SNR for Figure 5a.

Another intriguing observation in our numerical results is that the values of $\{\rho'_i\}_{i=1}^{N_\sigma}$ and $\{s'_i\}_{i=1}^{N_\sigma}$ given in Figures 6a and 6b do not always favor the eigenchannel with the largest eigenvalue. For instance, the strongest eigenchannel does not get the largest share in the available quantizers at SNR = 9 dB and SNR = 10 dB. To validate the result of Algorithm 1, we perform exhaustive search on the optimal $\{\rho_i\}_{i=1}^{N_\sigma}$ and $\{s_i\}_{i=1}^{N_\sigma}$ to get $C(\mathbf{A}_{\text{PH}})$. That is, we solve the convex power allocation strategy in (17) for all possible configurations of $\{s_i\}_{i=1}^{N_\sigma}$. This guarantees finding the global optimal solution in problem (14). There are a total of $\binom{N_q + N_\sigma - 1}{N_\sigma - 1}$ configurations that satisfy $\sum_{i=1}^{N_\sigma} s_i = N_q$. We tabulate the top 5 solutions for the joint optimization of $\{\rho_i\}_{i=1}^{N_\sigma}$ and $\{s_i\}_{i=1}^{N_\sigma}$ for SNR = 9 dB and SNR = 10 dB in Tables Ia and Ib, respectively. It can be seen that the optimal $\{\rho_i\}_{i=1}^{N_\sigma}$ and $\{s_i\}_{i=1}^{N_\sigma}$

(marked in blue) obtained by exhaustive search match those produced by Algorithm 1.

Despite the good agreement between our achievability scheme and the established capacity upper bound in the low SNR regime, the performance gap between the two widens as the SNR is increased in Figure 5a. We shall refer to the rate of our scheme in the infinite SNR regime as $R_\infty^{(\text{scheme})}$. $R_\infty^{(\text{scheme})}$ is maximized if the quantizers are uniformly allocated to the N_σ eigenchannels. In case N_q is not divisible by N_σ , we simply distribute the excess sign quantizers equally among the $N_q \bmod N_\sigma$ strongest eigenchannels. The number of possible outputs is

$$\mathcal{M}_{\text{scheme}} = \begin{cases} \left\{ 2 \left\lfloor \frac{N_q}{N_\sigma} \right\rfloor \right\}^u \times \left\{ 2 \left\lceil \frac{N_q}{N_\sigma} \right\rceil \right\}^{N_\sigma - u}, & \left\lfloor \frac{N_q}{N_\sigma} \right\rfloor > 0 \\ \left\{ 2 \left\lfloor \frac{N_q}{N_\sigma} \right\rfloor \right\}^u, & \text{otherwise} \end{cases}$$

TABLE I
 TOP 5 (OUT OF 455) JOINT OPTIMIZATION OF $\{\rho_i\}_{i=1}^{N_\sigma}$ AND $\{s_i\}_{i=1}^{N_\sigma}$ PRODUCED BY EXHAUSTIVE SEARCH FOR
 (a) SNR = 9 dB AND (b) SNR = 10 dB. HERE, $N_\sigma = 4$, $N_q = 12$, AND $\{\lambda_i\}_{i=1}^4 = \{1.6, 1.2, 0.8, 0.4\}$

(a)									
Rank #	R (in bpcu)	s_1	s_2	s_3	s_4	ρ_1	ρ_2	ρ_3	ρ_4
1	4.8349	3	5	4	0	0.3388	0.3328	0.3284	1.0614e-32
2	4.8299	3	4	5	0	0.33622	0.32926	0.33452	2.3924e-32
3	4.8245	5	3	4	0	0.32269	0.34696	0.33035	0
4	4.8237	3	6	3	0	0.33583	0.33724	0.32692	1.4565e-33
5	4.8217	4	3	5	0	0.32537	0.34113	0.3335	1.254e-32

(b)									
Rank #	R (in bpcu)	s_1	s_2	s_3	s_4	ρ_1	ρ_2	ρ_3	ρ_4
1	5.3992	3	3	4	2	0.27939	0.28615	0.27362	0.16083
2	5.3952	3	3	3	3	0.27632	0.28255	0.27284	0.16828
3	5.3927	3	3	5	1	0.29554	0.30523	0.30048	0.098749
4	5.39	3	4	3	2	0.28069	0.27659	0.2789	0.16382
5	5.3882	3	5	2	2	0.28219	0.27762	0.27293	0.16725

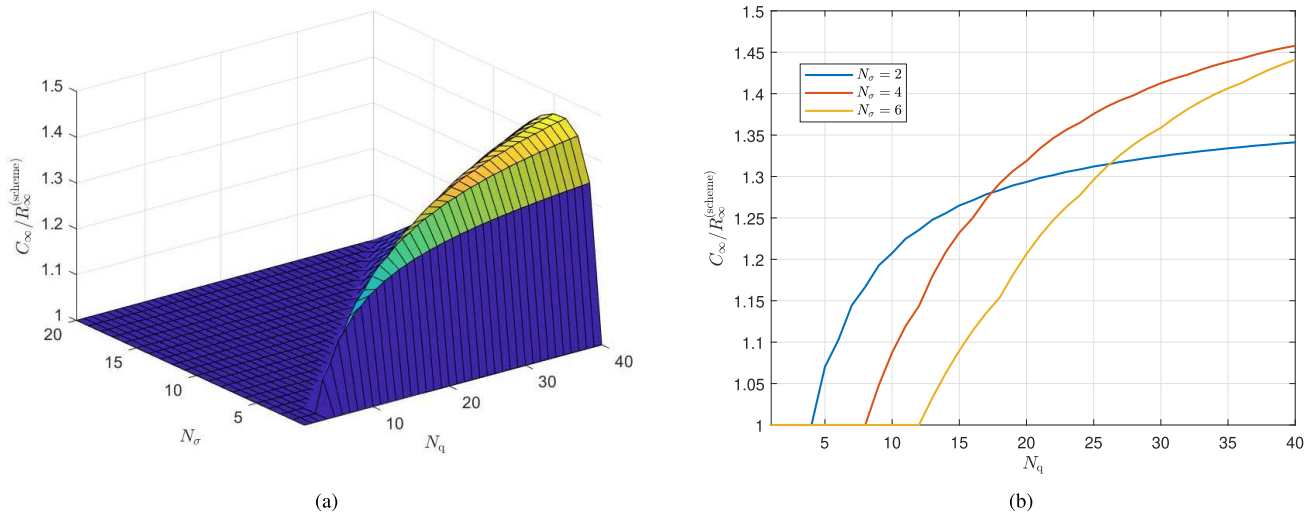


Fig. 7. (a) Ratio of $C_\infty/R_\infty^{(\text{scheme})}$ as a function N_σ and N_q , and (b) plot of $C_\infty/R_\infty^{(\text{scheme})}$ for $N_\sigma = 2, 4, 6$.

where $u = N_q \bmod N_\sigma$. The rate of our achievability scheme in the infinite SNR regime is $R_\infty^{(\text{scheme})} = \log \mathcal{M}_{\text{scheme}}$. To quantify the gap between C_∞ and $R_\infty^{(\text{scheme})}$ under different N_q and N_σ , we plot their ratio as a function of N_σ and N_q in Figure 7a. It can be observed that $R_\infty^{(\text{scheme})}$ coincides with C_∞ when $2N_q \leq N_\sigma$. This is the case in which we assign at most one sign quantizer to each real dimension. When $2N_q > N_\sigma$, a logarithmic increase in $C_\infty/R_\infty^{(\text{scheme})}$ is observed as N_q is increased. This is depicted in Figure 7b. The suboptimality of our achievability scheme in the infinite SNR regime comes from the restriction imposed on the matrix Φ . In our scheme, we simply design Φ to connect a sign quantizer to a single eigenchannel and then choose the input distribution for each eigenchannel independently. On the other hand, C_∞ is derived based on the intuition that a sign quantizer can be connected to multiple eigenchannels; thereby creating a hyperplane that passes through more than 2 real dimensions. In effect, the number of quantization regions created by N_q intersecting hyperplanes can exceed that of our achievability scheme.

VI. MIMO RECEIVER WITH PIPELINED PHASE ADC

To overcome the rate loss in the high SNR regime, we deviate our attention away from the one-shot receiver model and instead incorporate analog temporal and spatial processing techniques in the receiver design. To this end, we present a new MIMO receiver that utilizes an analog linear combiner and a more complex form of ADC structure. We call this ADC structure the *pipelined phase ADC*, since the key idea is borrowed from the pipelined ADC topology [48]. In the subsequent discussion, we shall elaborate on the operation of this pipelined phase ADC. We then give a formal description of the new MIMO receiver in Section VI-B and show how it can achieve the high SNR capacity of N_q bits/channel use.

A. The Pipelined Phase ADC

Figure 8a depicts a block diagram of the pipelined phase ADC. This ADC is composed of $L - 1$ pipeline stages, where L depends on the number of 1-bit ADCs. Each pipeline stage

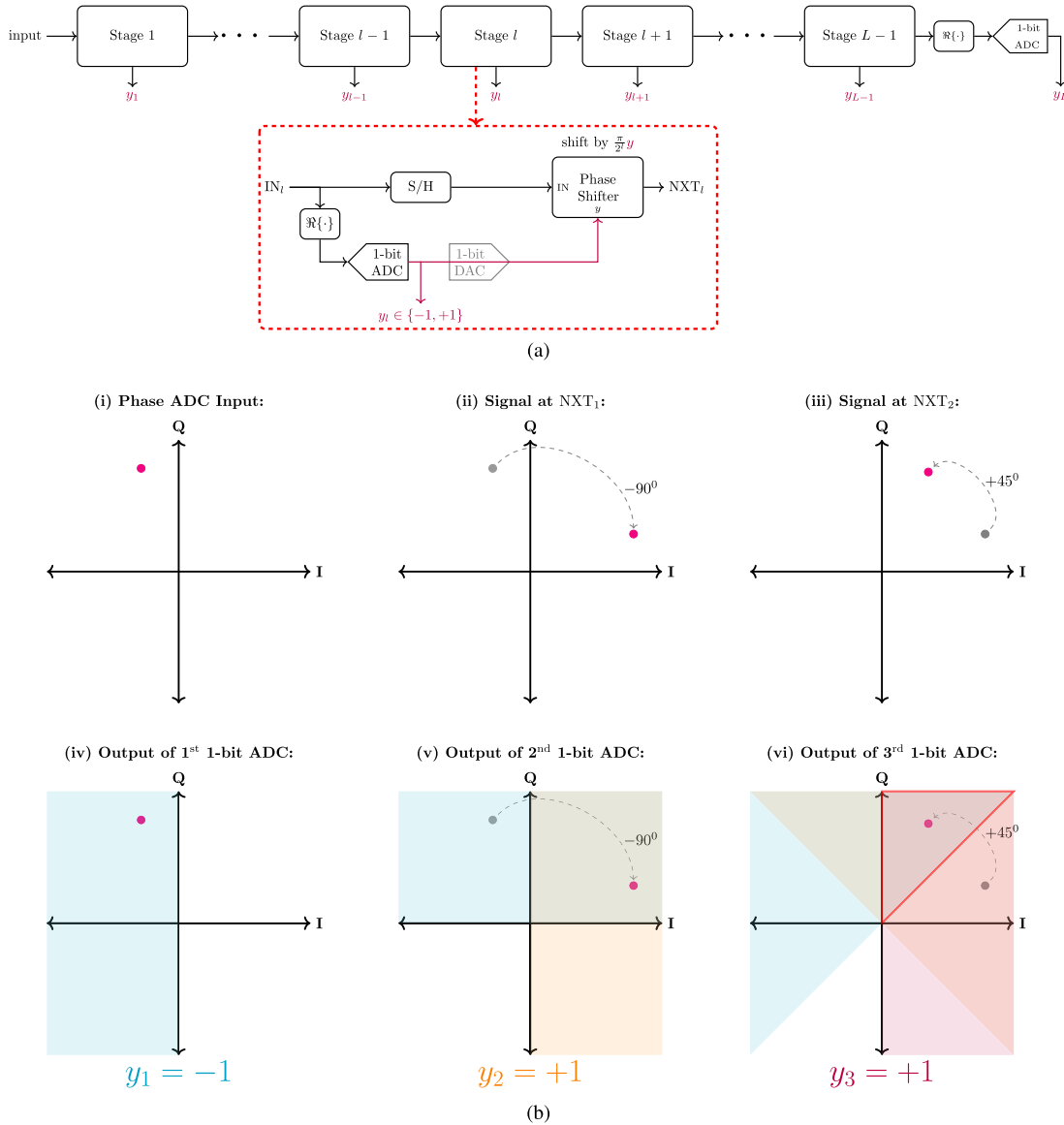


Fig. 8. (a) Block diagram of an $(L-1)$ -stage Pipelined Phase ADC and (b) an illustrative example (for $L=3$) to demonstrate the quantization process per stage.

consists of an analog delay element, in the form of a sample-and-hold (S/H) block, and a phase shifter to perform analog temporal processing. At the l -th pipeline stage, we apply 1-bit quantization to the input to get y_l . The 1-bit output y_l is then used by the phase shifter to apply an appropriate rotation to the S/H output.⁴ This phase shifted signal is then fed to the next pipeline stage for further processing.

We give a more concrete example of the quantization mechanism through an illustrative example in Figure 8b. Here, we assume $L=3$ so there are two pipeline stages and three 1-bit ADCs. The input signal is given by the magenta dot (shown in Figure 8b.i). The real component of the signal is fed to a 1-bit ADC of the 1st pipelined stage to produce y_1 (Figure 8b.iv). Since it falls at the LHS of the y -axis

⁴We note that a 1-bit digital-to-analog converter (DAC) might be required to interface the ADC output to the phase shifter. However, the 1-bit DAC can be eliminated if the analog phase shifter is digitally-controlled.

(cyan region), the 1-bit ADC outputs $y_1 = -1$. Consequently, this implies a clockwise phase shift of $\frac{\pi}{2} = 90^\circ$ (Figure 8b.ii); which will be fed to the next pipeline stage. This signal falls at the RHS of the y -axis (orange region); thus producing $y_2 = +1$ (Figure 8b.v). Notice that the intersection of the cyan and orange regions forms a quantization region of a 2-bit phase quantizer. With $y_2 = +1$, the signal at NXT_2 is phase shifted by $\frac{\pi}{4} = 45^\circ$ counter clockwise. The last ADC outputs $y_3 = +1$ since the resulting signal falls in the RHS of y -axis (purple region) (Figure 8b.vi). The intersection of the cyan, orange, and purple regions is a quantization region of a 3-bit phase quantizer.

In general, the pipelined phase ADC enables us to create an L -bit phase quantizer with length- $L-1$ delay using L 1-bit ADCs. Note that a flash ADC structure would require 2^L comparators to construct an L -bit phase quantizer. Moreover, because the analog pipelining structure enables each 1-bit

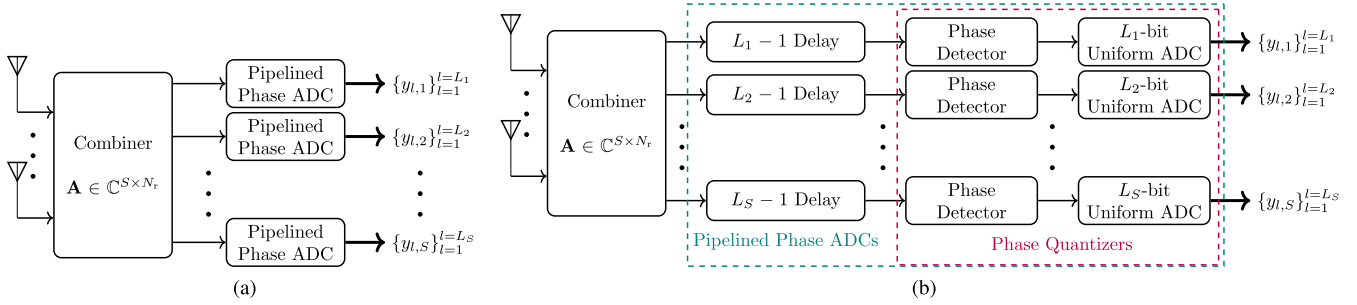


Fig. 9. (a) System Model of MIMO Receiver with Pipelined Phase ADCs, and its (b) equivalent model.

ADC to extract 1-bit of information at each channel use, the maximum rate of L bits/channel use is achievable. Note that the first $L - 1$ channel uses are strictly less than this rate since some analog delay elements do not contain signals initially. However, the definition of channel capacity applies for asymptotically large block lengths. Thus, this finite length delay is negligible in the asymptotic regime.

The proposed pipelined phase ADC has some resemblance with the ADC mechanism used in the adaptive threshold receiver recently proposed in [37]. While both ADC topologies exploit analog domain pipelining, the latter adaptively chooses the locations of the 1-bit ADC thresholds in the current channel use based on the previous channel uses. The former applies an appropriate phase shift to the input of the next pipeline stage depending on the 1-bit ADC output in the current pipeline stage.

B. Proposed Receiver

The block diagram for the proposed MIMO receiver employing pipelined phase ADCs is shown in Figure 9a. An analog combiner is used to perform analog spatial processing of the received signals. The S data streams at the analog combiner output are each fed to an L_S -bit pipelined phase ADC to produce $N_q = \sum_{i=1}^S L_i$ bits every channel use.

The achievability scheme presented in Section III can be extended to this receiver structure. We consider the SVD of the channel, which gives a precoded transmit strategy $\tilde{\mathbf{x}} = \mathbf{V}\mathbf{x}$ and an analog linear combiner $\mathbf{A} = \mathbf{U}_1^H$. Thus, there are $S = N_\sigma$ parallel data streams. Note that an L -bit pipelined phase ADC is equivalent to an L -bit phase quantizer (with some finite delay due to pipelining). As a result, the same reasoning in Section III can be used to adapt (14) to this receiver architecture. The resulting optimization problem is

$$\max_{L_i, \rho_i} \sum_{i=1}^{N_\sigma} L_i - w_{2L_i} \left(\frac{\lambda_i \rho_i}{\sigma^2}, \frac{\pi}{2L_i} \right) \quad (26a)$$

$$\text{s.t.} \quad \sum_{i=1}^{N_\sigma} L_i = N_q \quad (26b)$$

$$\sum_{i=1}^{N_\sigma} \rho_i \leq P \quad (26c)$$

$$\rho_i \geq 0, L_i \in \{1, \dots, N_q\} \quad (26d)$$

Consequently, Algorithm 1 can also be adapted to produce a heuristic solution to (26) by modifying Lines 7 and 10 accordingly.

As SNR grows unbounded, $w_{2L_i}(\cdot, \frac{\pi}{2L_i})$ vanishes and (26a) approaches N_q bits/channel use. This is the maximum rate that any channel with N_q -bit output quantization can achieve. Note that this rate can be larger than C_∞ established in Section IV-A. This is because the inclusion of analog temporal processing allows an analog sample to be quantized multiple times; thus the combinatorial geometry approach used in Section IV-A should be modified accordingly. To this end, we simply use the trivial upper bound N_q bit/channel use.

For the capacity upper bound in the finite SNR regime, we can use the DPI argument in Section IV-B to show that $C_{\phi\text{-detector}}$ upper bounds the capacity of our proposed MIMO receiver. The following corollary of Proposition 1 extends this result.

Corollary 1: The capacity of the Gaussian channel employing the MIMO receiver with pipelined phase ADCs can be upper bounded by (24).

Proof: To prove the claim, we consider the equivalent receiver model in Figure 9b. By DPI, the capacity is bounded by the maximum mutual information between the transmitted symbols and the output of the phase detector. Moreover, the finite length delay prior to the phase detector does not change the capacity. \square

We also point out that the use of sample-and-hold blocks, 1-bit digital-to-analog converters (DACs), and analog phase shifters in the pipelined phase ADCs entails additional power consumption to the MIMO receiver. In this work, we simply focus at how analog spatial and temporal processing can be used to maximize the achievable rate of a MIMO receiver for a given number of 1-bit ADCs. Determining the best architecture from an energy efficiency standpoint can be a future direction of this study.

VII. NUMERICAL RESULTS FOR THE MIMO RECEIVER WITH PIPELINED PHASE ADCS

In this section, we examine the achievable rate of the MIMO receiver with pipelined phase ADCs. A 4×4 MIMO Gaussian channel with $N_q = 12$ available 1-bit ADCs is considered. We set $P = 1$ and vary the SNR by changing the noise variance σ^2 . Moreover, we look at two sets of eigenvalues for the experiment setup: (a) $\lambda_1 = 1.6, \lambda_2 = 1.2, \lambda_3 = 0.8,$

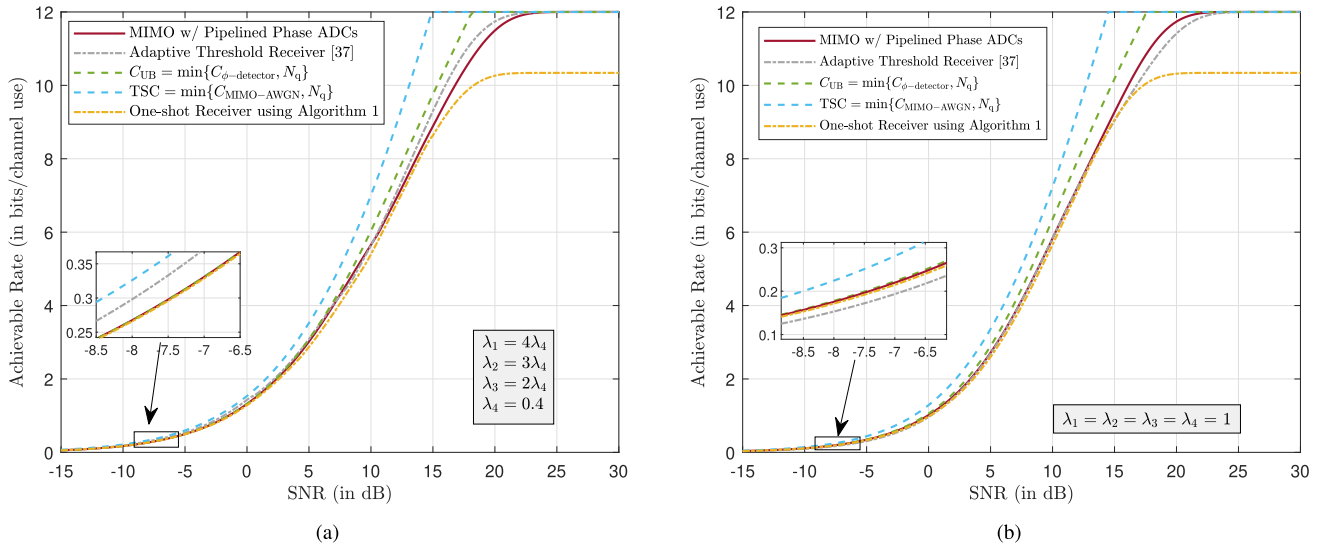


Fig. 10. Achievable Rate vs. SNR of the Proposed Receiver for (a) $\lambda_1 = 1.6, \lambda_2 = 1.2, \lambda_3 = 0.8,$ and $\lambda_4 = 0.4$; and (b) $\lambda_1 = \lambda_2 = \lambda_3 = \lambda_4 = 1$ ($N_q = 12$). The rate is shot compared to that of the adaptive threshold receiver.

and $\lambda_4 = 0.4$; and (b) $\lambda_1 = \lambda_2 = \lambda_3 = \lambda_4 = 1$. We shall refer to these channel setups as setup A and setup B.

The achievable rates of the MIMO receiver employing pipelined phase ADCs for the two for setup A and setup B are depicted in Figure 10a and 10b, respectively. To compute the achievable rate, we modify Algorithm 1 as described in Section VI. We also superimpose the achievable rate of the hybrid one-shot receiver with zero threshold ADCs, the capacity upper bound described in Section VI, and the TSC bound. It can be observed that the achievability scheme for the MIMO receiver with pipelined phase ADCs is tight with our established upper bound in the low SNR regime and also attains the high SNR capacity of N_q bits/channel. While the gap between the achievable rate of MIMO receiver with pipelined phase ADCs and that of the hybrid one-shot receiver in Section III is small in the low SNR regime, this gap gradually increases with SNR. This demonstrates that the rate increase provided by incorporating analog temporal processing in our receiver design is more pronounced in the high SNR regime.

We compare the performance of our proposed receiver to that of the adaptive threshold receiver in [37]. Since the adaptive threshold receiver is designed for real MIMO channels, we made some modifications in the channel setup for fair comparison. We considered an 8×8 real MIMO Gaussian channel with eigenvalues $\lambda'_{2i-1} = \lambda'_{2i} = \lambda_i$ for $i = 1, 2, 3, 4$. Furthermore, we set $\text{SNR} = 2P/\sigma^2$ instead of $\text{SNR} = P/\sigma^2$. The adaptive threshold receiver effectively creates parallel real eigenchannels with uniform quantization at the output. The transmit strategy for the achievability scheme described in [37] is to send equiprobable pulse amplitude modulation (PAM) over each eigenchannel. The power allocation per PAM strategy is obtained using the conventional waterfilling algorithm for the unquantized AWGN channel. Using this power allocation scheme, an exhaustive search procedure is performed to allocate the 1-bit ADCs. The achievable rate

results in Figure 10a and 10b show that the achievable rate of our proposed receiver may have inferior or superior performance than the adaptive threshold receiver depending on the SNR and channel eigenvalues.

One potential reason why our proposed receiver works better in some cases is because the transmit power and the 1-bit ADC allocation are jointly optimized. This is in contrast to the adaptive threshold receiver which performs separate optimization of the transmit power and ADC allocation. On the other hand, the adaptive threshold receiver extracts the amplitude information, which is neglected by our proposed receiver. This may explain why the adaptive threshold receiver outperformed our proposed receiver in Figure 10a. Nonetheless, we point out that the adaptive threshold receiver requires AGCs to adjust the dynamic range of the received signal.

VIII. CONCLUSION

In this work, we analyzed the capacity of a point-to-point Gaussian MIMO channel in which the receiver is equipped with N_q 1-bit ADCs and an analog linear combiner prior to quantization. In particular, we focused on the zero-threshold ADC case. Our first contribution is an achievability scheme in which the analog combiner is configured to create parallel Gaussian channels with phase quantization at the output. The achievable rate of this constructed channel is evaluated using an alternating optimization approach. We then established a new capacity upper bound that is tighter than the TSC bound when the ADCs are restricted to have zero threshold. This upper bound serves as a measure of the worst-case gap between the rate of our achievability scheme and the capacity of the channel. Our numerical results showed that the rate of our achievability scheme is tight in the low SNR regime. However, a performance gap exists in the high SNR regime whenever $N_\sigma \leq 2N_q$. More precisely, when this condition is satisfied, we observed that the ratio of the channel capacity and the rate of our achievability scheme in the infinite SNR

regime grows logarithmically with the number of 1-bit ADCs. To overcome this, a new receiver is proposed that implements joint analog spatial and temporal processing through the use of an analog combiner and pipelined phase ADCs. We showed that the proposed receiver achieves the high SNR capacity of N_q bits/channel use and outperforms the adaptive threshold receiver [37] when the channel eigenvalues are equal. Further research needs to be conducted to be able to generalize these results to multi-user setting and different fading environments. As mentioned in Section VI, investigation of the best architecture from an energy efficiency viewpoint is another interesting research direction.

APPENDIX A PROOF OF LEMMA 1

Suppose the optimal strategy \mathcal{O}' uses the ordered set of eigenchannels with eigenvalues

$$\mathcal{S}' = \{\lambda_1, \dots, \lambda_{N_s+1}\} \setminus \lambda_i$$

for some integer $i \in \{1, \dots, N_s\}$. The power and sign quantizers allocated to the k -th eigenchannel in this optimal strategy \mathcal{O}' are denoted as ρ'_k and s'_k , respectively. Define another strategy \mathcal{O}^* which uses the ordered set of eigenchannels with eigenvalues $\mathcal{S}^* = \{\lambda_1, \dots, \lambda_{N_s}\}$. Let ρ_k^* and s_k^* be the power and sign quantizer allocation in the k -th eigenchannel when the strategy \mathcal{O}^* is used. If we set $\rho_k^* = \rho'_k$ and $s_k^* = s'_k$ $\forall k = 1, \dots, i-1, i+1, \dots, N_s$ and let $\rho_i^* = \rho'_{N_s+1}$ and $s_i^* = s'_{N_s+1}$, then the difference between the rate of \mathcal{O}^* and \mathcal{O}' is

$$R_{\mathcal{O}^*} - R_{\mathcal{O}'} = w_{2s_i} \left(\frac{\lambda_{N_s+1} \rho_i^*}{\sigma^2}, \frac{\pi}{2s_i} \right) - w_{2s_i} \left(\frac{\lambda_i \rho_i^*}{\sigma^2}, \frac{\pi}{2s_i} \right) \geq 0.$$

The inequality comes from the monotonic decreasing property of phase quantization entropy with respect to ν [21, Proposition 1] and $\lambda_{N_s+1} \leq \lambda_i$. This contradicts the assumption that strategy \mathcal{O}' is optimal.

APPENDIX B PROOF OF CORRECTNESS OF THE DYNAMIC PROGRAMMING APPROACH

The key technique in showing the correctness is through strong induction. First, the base case $i = 0$ or $n_q = 0$ is true since if there is no channel to send information or there is no sign quantizer that can output produce the output \mathbf{y} , then $I(\mathbf{x}; \mathbf{y})$ should be 0. Next, we prove the inductive step. Assume $f(i', n'_q)$ be the optimal solution for all $i' < i$. We need to show $f(i, n_q)$ is the optimal solution for the state (i, n_q) . Note that the sum capacity of channels 1 to i with n_q available sign quantizers can be expressed as

$$\begin{aligned} & \sum_{j=1}^i \log 2 s_j - w_{2s_j} \left(\frac{\lambda_j \rho_j}{\sigma^2}, \frac{\pi}{2s_j} \right) \\ &= \sum_{j=1}^{i-1} \log 2 s_j - w_{2s_j} \left(\frac{\lambda_j \rho_j}{\sigma^2}, \frac{\pi}{2s_j} \right) \end{aligned}$$

$$\begin{aligned} & + \log 2 s_i - w_{2s_i} \left(\frac{\lambda_i \rho_i}{\sigma^2}, \frac{\pi}{2s_i} \right) \\ & \leq f(i-1, n_q - s_i) + \log 2 s_i - w_{2s_i} \left(\frac{\lambda_i \rho_i}{\sigma^2}, \frac{\pi}{2s_i} \right). \end{aligned}$$

The first line follows from isolating the capacity of the i -th channel from channels i to $i-1$. The inequality in the second line follows from the optimality of $f(i', n'_q)$ and equality is achieved by choosing the optimal s_i . Hence, the problem has an optimal substructure property. The algorithm considers all possible choices of s_i and compares their values. Thus, optimality of $f(i, n_q)$ is guaranteed.

APPENDIX C PROOF OF LEMMA 2

By the chain rule of mutual information, we have

$$\begin{aligned} I_{\mathbf{A}}(\mathbf{x}; \tilde{\mathbf{y}}) &= I_{\mathbf{A}}(\mathbf{x}; \tilde{\mathbf{y}}^{(1)}, \tilde{\mathbf{y}}^{(2)}) \\ &= I_{\mathbf{A}_1}(\mathbf{x}; \tilde{\mathbf{y}}^{(1)}) + I_{\mathbf{A}_2}(\mathbf{x}; \tilde{\mathbf{y}}^{(2)} | \tilde{\mathbf{y}}^{(1)}). \end{aligned}$$

The claim is proven if we can show that $I_{\mathbf{A}_2}(\mathbf{x}; \tilde{\mathbf{y}}^{(2)} | \tilde{\mathbf{y}}^{(1)}) = 0$. In other words, $\mathbf{x} \rightarrow \tilde{\mathbf{y}}^{(1)} \rightarrow \tilde{\mathbf{y}}^{(2)}$ should form a Markov chain. The term $I_{\mathbf{A}_2}(\mathbf{x}; \tilde{\mathbf{y}}^{(2)} | \tilde{\mathbf{y}}^{(1)})$ can be expressed as

$$\begin{aligned} &= I_{\mathbf{A}_2}(\mathbf{x}; e^{j\tilde{\mathbf{y}}^{(2)}} | \tilde{\mathbf{y}}^{(1)}) \\ &= I_{\mathbf{A}_2}(\mathbf{x}; e^{j\{\mathbf{z}'_{N_\sigma+1:N_q} - \mathbf{B}^{(2)}\{\mathbf{B}^{(1)}\}^\dagger \mathbf{z}'_{1:N_\sigma}\}} | \tilde{\mathbf{y}}^{(1)}) \\ &= I_{\mathbf{A}_2}(\mathbf{x}; \mathbf{z}'_{N_\sigma+1:N_q} - \mathbf{B}^{(2)}\{\mathbf{B}^{(1)}\}^\dagger \mathbf{z}'_{1:N_\sigma} | \tilde{\mathbf{y}}^{(1)}) \\ &= 0. \end{aligned}$$

The equality in the first line follows from the fact that $e^{j(\cdot)}$ is bijective. The second equality is obtained by noting that

$$\begin{aligned} & \exp(j\tilde{\mathbf{y}}^{(2)}) \cdot \exp(-j\mathbf{B}^{(2)}\{\mathbf{B}^{(1)}\}^\dagger \tilde{\mathbf{y}}^{(1)}) \\ &= \exp(j(\tilde{\mathbf{y}}^{(2)} - \mathbf{B}^{(2)}\{\mathbf{B}^{(1)}\}^\dagger \tilde{\mathbf{y}}^{(1)})) \\ &= \exp(j(\mathbf{z}'_{N_\sigma+1:N_q} - \mathbf{B}^{(2)}\{\mathbf{B}^{(1)}\}^\dagger \mathbf{z}'_{1:N_\sigma})), \end{aligned}$$

where $\{\cdot\}^\dagger$ is the Moore-Penrose inverse operator. Since the transmitted symbol and the phase of the additive noise components are independent, $I_{\mathbf{A}_2}(\mathbf{x}; \tilde{\mathbf{y}}^{(2)} | \tilde{\mathbf{y}}^{(1)}) = 0$.

To prove the second claim, we note that we have full control over \mathbf{A} . If \mathbf{B}_1 is not full rank, we can create \mathbf{A}' by permuting the rows of \mathbf{A} to make \mathbf{B}'_1 full rank. This, in effect, reorders the elements of $\tilde{\mathbf{y}}$ but does not change the mutual information since reordering is a bijective mapping. Thus, $I_{\mathbf{A}}(\mathbf{x}; \tilde{\mathbf{y}}) = I_{\mathbf{A}'}(\mathbf{x}; \tilde{\mathbf{y}})$. If no such row permutation of rows can make \mathbf{B}_1 full rank, then this implies that $\text{rank}\{\mathbf{A}\mathbf{H}\} < N_\sigma$. We can simply use an analog combiner \mathbf{A}' with $\text{rank}\{\mathbf{A}'\mathbf{H}\} = N_\sigma$ and employ a transmit strategy $F_{\mathbf{X}'}$ that only uses the $\text{rank}\{\mathbf{A}\mathbf{H}\}$ out of the N_σ eigenchannels so that $I_{\mathbf{A}}(\mathbf{x}; \tilde{\mathbf{y}}) = I_{\mathbf{A}'}(\mathbf{x}'; \tilde{\mathbf{y}})$.

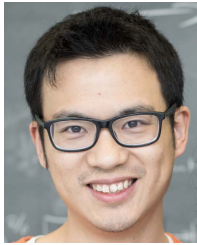
REFERENCES

- [1] C. Risi, D. Persson, and E. G. Larsson, "Massive MIMO with 1-bit ADC," 2014, *arXiv:1404.7736*.
- [2] T. S. Rappaport *et al.*, "Millimeter wave mobile communications for 5G cellular: It will work!" *IEEE Access*, vol. 1, pp. 335–349, 2013.

- [3] S. Rangan, T. S. Rappaport, and E. Erkip, "Millimeter-wave cellular wireless networks: Potentials and challenges," *Proc. IEEE*, vol. 102, no. 3, pp. 366–385, Feb. 2014.
- [4] H. Halbauer and T. Wild, "Towards power efficient 6G sub-THz transmission," in *Proc. Joint Eur. Conf. Netw. Commun. 6G Summit (EuCNC/6G Summit)*, Jun. 2021, pp. 25–30.
- [5] R. H. Walden, "Analog-to-digital converter survey and analysis," *IEEE J. Sel. Areas Commun.*, vol. 17, no. 4, pp. 539–550, Apr. 1999.
- [6] E. Björnson, M. Matthaiou, and M. Debbah, "Massive MIMO with non-ideal arbitrary arrays: Hardware scaling laws and circuit-aware design," *IEEE Trans. Wireless Commun.*, vol. 14, no. 8, pp. 4353–4368, Aug. 2015.
- [7] S. Jacobsson, G. Durisi, M. Coldrey, U. Gustavsson, and C. Studer, "One-bit massive MIMO: Channel estimation and high-order modulations," in *Proc. IEEE Int. Conf. Commun. Workshop (ICCW)*, Jun. 2015, pp. 1304–1309.
- [8] Y. Li, C. Tao, G. Seco-Granados, A. Mezghani, A. L. Swindlehurst, and L. Liu, "Channel estimation and performance analysis of one-bit massive MIMO systems," *IEEE Trans. Signal Process.*, vol. 65, no. 15, pp. 4075–4089, Apr. 2017.
- [9] N. I. Bernardo, J. Zhu, and J. Evans, "On minimizing symbol error rate over fading channels with low-resolution quantization," *IEEE Trans. Commun.*, vol. 69, no. 11, pp. 7205–7221, Nov. 2021.
- [10] S. Gayan, R. Senanayake, H. Inaltekin, and J. Evans, "Reliability characterization for SIMO communication systems with low-resolution phase quantization under Rayleigh fading," *IEEE Open J. Commun. Soc.*, vol. 2, pp. 2660–2679, 2021.
- [11] F. Sun, J. Singh, and U. Madhow, "Automatic gain control for ADC-limited communication," in *Proc. IEEE Global Telecommun. Conf. GLOBECOM*, Dec. 2010, pp. 1–5.
- [12] O. Dabeer and U. Madhow, "Channel estimation with low-precision analog-to-digital conversion," in *Proc. IEEE Int. Conf. Commun.*, May 2010, pp. 1–6.
- [13] L. Jun, L. Zhongqiang, and X. Xingzhong, "Low-complexity synchronization scheme with low-resolution ADCs," *Information*, vol. 9, no. 12, p. 313, Dec. 2018.
- [14] M. Schluter, M. Dörpinghaus, and G. P. Fettweis, "Bounds on phase, frequency, and timing synchronization in fully digital receivers with 1-bit quantization and oversampling," *IEEE Trans. Commun.*, vol. 68, no. 10, pp. 6499–6513, Oct. 2020.
- [15] J. Singh, O. Dabeer, and U. Madhow, "On the limits of communication with low-precision analog-to-digital conversion at the receiver," *IEEE Trans. Commun.*, vol. 57, no. 12, pp. 3629–3639, Dec. 2009.
- [16] A. Mezghani and J. A. Nossek, "On ultra-wideband MIMO systems with 1-bit quantized outputs: Performance analysis and input optimization," in *Proc. Int. Symp. Inf. Theory (ISIT)*, Jun. 2007, pp. 1286–1289.
- [17] S. Krone and G. Fettweis, "Fading channels with 1-bit output quantization: Optimal modulation, ergodic capacity and outage probability," in *Proc. IEEE Inf. Theory Workshop*, Aug. 2010, pp. 1–5.
- [18] T. Koch and A. Lapidoth, "At low SNR, asymmetric quantizers are better," *IEEE Trans. Inf. Theory*, vol. 59, no. 9, pp. 5421–5445, Sep. 2013.
- [19] M. N. Vu, N. H. Tran, D. G. Wijeratne, K. Pham, K.-S. Lee, and D. H. N. Nguyen, "Optimal signaling schemes and capacity of non-coherent Rician fading channels with low-resolution output quantization," *IEEE Trans. Wireless Commun.*, vol. 18, no. 6, pp. 2989–3004, Jun. 2019.
- [20] M. Rahman, M. Ranjbar, and N. Tran, "On the capacity-achieving scheme and capacity of 1-bit ADC Gaussian-mixture channels," *EAI Endorsed Trans. Ind. Netw. Intell. Syst.*, vol. 7, no. 22, Jan. 2020, Art. no. 162830.
- [21] N. I. Bernardo, J. Zhu, and J. Evans, "On the capacity-achieving input of channels with phase quantization," *IEEE Trans. Inf. Theory*, vol. 68, no. 9, pp. 5866–5888, Sep. 2022.
- [22] N. I. Bernardo, J. Zhu, and J. Evans, "On the capacity-achieving input of the Gaussian channel with polar quantization," *IEEE Trans. Commun.*, vol. 70, no. 9, pp. 5912–5928, Sep. 2022.
- [23] A. Mezghani and J. A. Nossek, "Capacity lower bound of MIMO channels with output quantization and correlated noise," in *Proc. IEEE Int. Symp. Inf. Theory*, Jul. 2012, pp. 1–5.
- [24] J. Mo and R. W. Heath, Jr., "Capacity analysis of one-bit quantized MIMO systems with transmitter channel state information," *IEEE Trans. Signal Process.*, vol. 63, no. 20, pp. 5498–5512, Oct. 2015.
- [25] O. Orhan, E. Erkip, and S. Rangan, "Low power analog-to-digital conversion in millimeter wave systems: Impact of resolution and bandwidth on performance," in *Proc. Inf. Theory Appl. Workshop (ITA)*, Feb. 2015, pp. 191–198.
- [26] O. T. Demir and E. Björnson, "The Bussgang decomposition of nonlinear systems: Basic theory and MIMO extensions [lecture notes]," *IEEE Signal Process. Mag.*, vol. 38, no. 1, pp. 131–136, Jan. 2021.
- [27] J. Zhang, L. Dai, X. Li, Y. Liu, and L. Hanzo, "On low-resolution ADCs in practical 5G millimeter-wave massive MIMO systems," *IEEE Commun. Mag.*, vol. 56, no. 7, pp. 205–211, Jul. 2018.
- [28] J. Liu, Z. Luo, and X. Xiong, "Low-resolution ADCs for wireless communication: A comprehensive survey," *IEEE Access*, vol. 7, pp. 91291–91324, 2019.
- [29] J. Choi, G. Lee, A. Alkhateeb, A. Gatherer, N. Al-Dhahir, and B. L. Evans, "Advanced receiver architectures for millimeter-wave communications with low-resolution ADCs," *IEEE Commun. Mag.*, vol. 58, no. 8, pp. 42–48, Aug. 2020.
- [30] N. Liang and W. Zhang, "Mixed-ADC massive MIMO uplink in frequency-selective channels," *IEEE Trans. Commun.*, vol. 64, no. 11, pp. 4652–4666, Nov. 2016.
- [31] J. Mo, A. Alkhateeb, S. Abu-Surra, and R. W. Heath, Jr., "Hybrid architectures with few-bit ADC receivers: Achievable rates and energy-rate tradeoffs," *IEEE Trans. Wireless Commun.*, vol. 16, no. 4, pp. 2274–2287, Apr. 2017.
- [32] K. Roth and J. A. Nossek, "Achievable rate and energy efficiency of hybrid and digital beamforming receivers with low resolution ADC," *IEEE J. Sel. Areas Commun.*, vol. 35, no. 9, pp. 2056–2068, Sep. 2017.
- [33] N. Shlezinger, Y. C. Eldar, and M. R. D. Rodrigues, "Hardware-limited task-based quantization," *IEEE Trans. Signal Process.*, vol. 67, no. 20, pp. 5223–5238, Oct. 2019.
- [34] N. Shlezinger and Y. C. Eldar, "Task-based quantization with application to MIMO receivers," *Commun. Inf. Syst.*, vol. 20, no. 2, pp. 131–162, 2020.
- [35] Y.-S. Jeon, S.-N. Hong, and N. Lee, "Supervised-learning-aided communication framework for MIMO systems with low-resolution ADCs," *IEEE Trans. Veh. Technol.*, vol. 67, no. 8, pp. 7299–7313, Aug. 2018.
- [36] Y. Zhang, M. Alrabeiah, and A. Alkhateeb, "Deep learning for massive MIMO with 1-bit ADCs: When more antennas need fewer pilots," *IEEE Wireless Commun. Lett.*, vol. 9, no. 8, pp. 1273–1277, Aug. 2020.
- [37] A. Khalili, F. Shirani, E. Erkip, and Y. C. Eldar, "MIMO networks with one-bit ADCs: Receiver design and communication strategies," *IEEE Trans. Commun.*, vol. 70, no. 3, pp. 1580–1594, Mar. 2022.
- [38] S. Rini, L. Barletta, Y. C. Eldar, and E. Erkip, "A general framework for MIMO receivers with low-resolution quantization," in *Proc. IEEE Inf. Theory Workshop (ITW)*, Nov. 2017, pp. 599–603.
- [39] A. Khalili, F. Shirani, E. Erkip, and Y. C. Eldar, "Tradeoff between delay and high SNR capacity in quantized MIMO systems," in *Proc. IEEE Int. Symp. Inf. Theory (ISIT)*, Jul. 2019, pp. 597–601.
- [40] S. Jacobsson, G. Durisi, M. Coldrey, U. Gustavsson, and C. Studer, "Throughput analysis of massive MIMO uplink with low-resolution ADCs," *IEEE Trans. Wireless Commun.*, vol. 16, no. 6, pp. 4038–4051, Jun. 2017.
- [41] A. Khalili, S. Rini, L. Barletta, E. Erkip, and Y. C. Eldar, "On MIMO channel capacity with output quantization constraints," in *Proc. IEEE Int. Symp. Inf. Theory (ISIT)*, Jun. 2018, pp. 1355–1359.
- [42] A. Dytso, M. Goldenbaum, H. V. Poor, and S. S. Shitz, "When are discrete channel inputs optimal?— optimization techniques and some new results," in *Proc. 52nd Annu. Conf. Inf. Sci. Syst. (CISS)*, Mar. 2018, pp. 1–6.
- [43] R. E. Blahut, "Computation of channel capacity and rate-distortion functions," *IEEE Trans. Inf. Theory*, vol. IT-18, no. 4, pp. 460–473, Jul. 1972.
- [44] J. Huang and S. P. Meyn, "Characterization and computation of optimal distributions for channel coding," *IEEE Trans. Inf. Theory*, vol. 51, no. 7, pp. 2336–2351, Jul. 2005.
- [45] A. Khalili, S. Shahsavari, F. Shirani, E. Erkip, and Y. C. Eldar, "On throughput of millimeter wave MIMO systems with low resolution ADCs," in *Proc. IEEE Int. Conf. Acoust., Speech Signal Process. (ICASSP)*, May 2020, pp. 5255–5259.
- [46] S. Boyd and L. Vandenberghe, *Convex Optimization*. Cambridge, U.K.: Cambridge Univ. Press, Mar. 2004.
- [47] C. Ho and S. Zimmerman, "On the number of regions in an m -dimensional space cut by n hyperplanes," *Gazette Austral. Math. Soc.*, vol. 33, no. 4, pp. 260–264, 2006.
- [48] T. B. Cho and P. R. Gray, "A 10 b, 20 Msample/s, 35 mW pipeline A/D converter," *IEEE J. Solid-State Circuits*, vol. 30, no. 3, pp. 166–172, Mar. 1995.



Neil Irwin Bernardo (Graduate Student Member, IEEE) received the B.S. degree in electronics and communications engineering and the M.S. degree in electrical engineering from the University of the Philippines Diliman in 2014 and 2016, respectively. He is currently pursuing the Ph.D. degree in engineering with The University of Melbourne, Australia. He has been a Faculty Member of the University of the Philippines Diliman since 2014. His research interests include wireless communications, signal processing, and information theory.



Jingge Zhu (Member, IEEE) received the B.S. and M.S. degrees in electrical engineering from Shanghai Jiao Tong University, Shanghai, China, in 2008 and 2011, respectively, the Dipl.-Ing. degree in technische Informatik from Technische Universität Berlin, Berlin, Germany, in 2011, and the Doctorat ès Sciences degree from the Ecole Polytechnique Fédérale (EPFL), Lausanne, Switzerland, in 2016.

He was a Post-Doctoral Researcher at the University of California at Berkeley, Berkeley, from 2016 to 2018. He is currently a Lecturer at The

University of Melbourne, Australia. His research interests include information theory with applications in communication systems and machine learning. He received the Discovery Early Career Research Award (DECRA) from the Australian Research Council in 2021, the IEEE Heinrich Hertz Award for Best Communications Letters in 2013, the Early Post-Doctoral Mobility Fellowship from Swiss National Science Foundation in 2015, and the Chinese Government Award for Outstanding Students Abroad in 2016.



Yonina C. Eldar (Fellow, IEEE) received the B.Sc. degree in physics and the B.Sc. degree in electrical engineering from Tel-Aviv University (TAU), Tel-Aviv, Israel, in 1995 and 1996, respectively, and the Ph.D. degree in electrical engineering and computer science from the Massachusetts Institute of Technology (MIT), Cambridge, in 2002.

She is currently a Professor with the Department of Mathematics and Computer Science, Weizmann Institute of Science, Rehovot, Israel. Previously, she was a Professor with the Department of Electrical

Engineering, Technion. She is also a Visiting Professor at MIT, a Visiting Scientist at the Broad Institute, and an Adjunct Professor at Duke University. She was a Visiting Professor at Stanford. She is the author of the book *Sampling Theory: Beyond Bandlimited Systems* and the coauthor of five other books published by Cambridge University Press. Her research interests are in the broad areas of statistical signal processing, sampling theory and compressed sensing, learning and optimization methods, and their applications to biology, medical imaging, and optics.

Dr. Eldar has received many awards for excellence in research and teaching, including the IEEE Signal Processing Society Technical Achievement Award in 2013, the IEEE/AESS Fred Nathanson Memorial Radar Award in 2014, and the IEEE Kiyo Tomiyasu Award in 2016. She received the Michael Bruno Memorial Award from the Rothschild Foundation, the Weizmann Prize for Exact Sciences, the Wolf Foundation Krill Prize for Excellence in Scientific Research, the Henry Taub Prize for Excellence in Research (twice), the Hershel Rich Innovation Award (three times), the Award for Women with Distinguished Contributions, the Andre and Bella Meyer Lectureship, the Career Development Chair at the Technion, the Muriel & David Jacknow Award for Excellence in Teaching, and the Technion's Award for Excellence in Teaching (two times). She received several best paper awards and best demo awards together with her research students and colleagues, including the SIAM outstanding Paper Prize, the UFFC Outstanding Paper Award, the Signal Processing Society Best Paper Award, and the IET Circuits, Devices and Systems Premium Award. She was selected as one of the 50 most influential women in Israel and Asia. She was a Horev Fellow of the Leaders in Science and Technology Program at the Technion and an Alon Fellow. She is also a Highly Cited Researcher. She is a member of the Israel Academy of Sciences and Humanities (elected 2017) and a EURASIP Fellow. She was a member of the Young Israel Academy of Science and Humanities and the Israel Committee for Higher Education. She is the Editor-in-Chief of *Foundations and Trends in Signal Processing*, a member of the IEEE Sensor Array and Multichannel Technical Committee, and serves on several other IEEE committees. In the past, she was a Signal Processing Society Distinguished Lecturer, a member of the IEEE Signal Processing Theory and Methods and Bio Imaging Signal Processing technical committees, and served as an Associate Editor for the IEEE TRANSACTIONS ON SIGNAL PROCESSING, the *EURASIP Journal of Signal Processing*, the *SIAM Journal on Matrix Analysis and Applications*, and the *SIAM Journal on Imaging Sciences*. She was the co-chair and the technical co-chair of several international conferences and workshops.



Jamie Evans (Senior Member, IEEE) was born in Newcastle, Australia, in 1970. He received the B.S. degree in physics and the B.E. degree in computer engineering from The University of Newcastle, in 1992 and 1993, respectively, and the M.S. and Ph.D. degrees in electrical engineering from The University of Melbourne, Australia, in 1996 and 1998, respectively. From March 1998 to June 1999, he was a Visiting Researcher with the Department of Electrical Engineering and Computer Science, University of California at Berkeley, Berkeley. Since

returning to Australia in July 1999, he has held an academic positions at The University of Sydney, The University of Melbourne, and Monash University. He is currently a Professor in electrical and electronic engineering and the Pro Vice-Chancellor (Education) at The University of Melbourne. His research interests are in communications theory, information theory, and statistical signal processing, with a focus on wireless communications networks. He received the University Medal upon graduation from The University of Newcastle. He was also awarded the Chancellor's Prize for excellence for his Ph.D. thesis.