

Machine Learning Inspired Energy-Efficient Hybrid Precoding for MmWave Massive MIMO Systems

Xinyu Gao*, Linglong Dai*, Ying Sun*, Shuangfeng Han[†], and Chih-Lin I[†]

*Tsinghua National Laboratory for Information Science and Technology (TNList),

Department of Electronic Engineering, Tsinghua University, Beijing, China

[†]Green Communication Research Center, China Mobile Research Institute, Beijing 100053, China

Abstract—Hybrid precoding is a promising technique for mmWave massive MIMO systems, as it can considerably reduce the number of required radio-frequency (RF) chains without obvious performance loss. However, most of the existing hybrid precoding schemes require a complicated phase shifter network, which still involves high energy consumption. In this paper, we propose an energy-efficient hybrid precoding architecture, where the analog part is realized by a small number of switches and inverters instead of a large number of high-resolution phase shifters. Our analysis proves that the performance gap between the proposed hybrid precoding architecture and the traditional one is small and keeps constant when the number of antennas goes to infinity. Then, inspired by the cross-entropy (CE) optimization developed in machine learning, we propose an adaptive CE (ACE)-based hybrid precoding scheme for this new architecture. It aims to adaptively update the probability distributions of the elements in hybrid precoder by minimizing the CE, which can generate a solution close to the optimal one with a sufficiently high probability. Simulation results verify that our scheme can achieve the near-optimal sum-rate performance and much higher energy efficiency than traditional schemes.

I. INTRODUCTION

Millimeter-wave (mmWave) massive multiple-input multiple-output (MIMO) has been considered as a promising technology for future 5G wireless communications [1], since it can provide wider bandwidth [2] and achieve higher spectral efficiency [3]. However, in MIMO systems, each antenna usually requires a dedicated radio-frequency (RF) chain (including high-resolution digital-to-analog converter, mixer, etc.) to realize the fully digital signal processing (e.g., precoding) [4]. For mmWave massive MIMO, this will result in unaffordable hardware complexity and energy consumption, as the number of antennas becomes huge and the energy consumption of RF chain is high [5]. To reduce the number of required RF chains, hybrid precoding has been recently proposed [6]. Its key idea is to decompose the fully digital precoder into a large-size analog beamformer (realized by the analog circuit) and a small-size digital precoder (requiring a small number of RF chains). Thanks to the low-rank characteristics of mmWave channels [2], a small-size digital precoder can achieve the spatial multiplexing gains, making hybrid precoding enjoy the near-optimal performance [5].

Nevertheless, most of the existing hybrid precoding schemes require a complicated phase shifter network, where each RF chain is connected to all antennas with high-resolution phase shifters [6], [7]. Although this architecture can provide high

design freedom to achieve the near-optimal performance, it requires hundreds or even thousands of high-resolution phase shifters with high hardware cost and energy consumption [5]. To solve this problem, two categories of schemes have been proposed very recently. The first category is to directly employ finite-resolution phase shifters instead of high-resolution phase shifters [8], [9]. It can reduce the energy consumption of phase shifter network without obvious performance loss, but it still requires a large number of phase shifters with considerable energy consumption. The second category is to utilize the switch network to replace the phase shifter network [10]–[12]. It can significantly reduce the hardware cost and energy consumption, but it suffers from an obvious performance loss.

In this paper, we propose a switch and inverter (SI)-based hybrid precoding architecture with considerably reduced hardware cost and energy consumption. Instead of using phase shifters, the analog part of the proposed architecture is realized by a small number of energy-efficient switches and inverters. Then, we provide the performance analysis to quantify the performance gap between the proposed hybrid precoding architecture and the traditional ones. After that, inspired by the cross-entropy (CE) optimization developed in machine learning [13], we propose an adaptive CE (ACE)-based hybrid precoding scheme for this new architecture. Specifically, according to the probability distributions of the elements in hybrid precoder, this scheme first randomly generates several candidate hybrid precoders. Then, it adaptively weights these candidate hybrid precoders according to their achievable sum-rates, and refines the probability distributions of elements in hybrid precoder by minimizing the CE. Repeating such procedure, we can finally generate a hybrid precoder close to the optimal one with a sufficiently high probability. Simulation results verify that our scheme can achieve the near-optimal sum-rate performance and much higher energy efficiency than traditional schemes.

Notation: Lower- and upper-case boldface letters denote vectors and matrices, respectively; $(\cdot)^T$, $(\cdot)^H$, $(\cdot)^{-1}$, $\text{tr}(\cdot)$, and $\|\cdot\|_F$ denote the transpose, conjugate transpose, inversion, trace, and Frobenius norm of a matrix, respectively; $|\cdot|$ denotes the absolute operator; $\mathbb{E}(\cdot)$ denotes the expectation; \otimes denotes the kronecker product; \mathbf{I}_N is the $N \times N$ identity matrix.

II. SYSTEM MODEL

In this paper, we consider a typical mmWave massive MIMO system, where the base station (BS) employs N antennas

and N_{RF} RF chains to simultaneously serve K single-antenna users (the extension to users with multiple-antennas is also possible as will be explained later). To fully achieve the multiplexing gains, we assume $N_{\text{RF}} = K$ [9]. For hybrid precoding systems as shown in Fig. 1, the $K \times 1$ received signal vector \mathbf{y} for all K users can be presented by

$$\mathbf{y} = \mathbf{H}\mathbf{F}_{\text{RF}}\mathbf{F}_{\text{BB}}\mathbf{s} + \mathbf{n}, \quad (1)$$

where $\mathbf{H} = [\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_K]^H$ is the channel matrix with \mathbf{h}_k presenting the $N \times 1$ channel vector between the BS and the k th user, \mathbf{s} is the $K \times 1$ transmitted signal vector for all K users satisfying $\mathbb{E}(\mathbf{s}\mathbf{s}^H) = \mathbf{I}_K$, \mathbf{F}_{RF} of size $N \times N_{\text{RF}}$ is the analog beamformer realized by analog circuit (different architectures incur different hardware constraints as will be discussed later), \mathbf{F}_{BB} is the baseband digital precoder of size $N_{\text{RF}} \times K$ satisfying the total transmit power constraint as $\|\mathbf{F}_{\text{RF}}\mathbf{F}_{\text{BB}}\|_F^2 = \rho$, where ρ is total transmit power. Finally, $\mathbf{n} \sim \mathcal{CN}(0, \sigma^2\mathbf{I}_K)$ of size $K \times 1$ is the additive white Gaussian noise (AWGN) vector, where σ^2 presents the noise power.

For the channel vector \mathbf{h}_k of the k th user, we adopt the geometric channel model to capture the characteristics of mmWave massive MIMO channels as [5]

$$\mathbf{h}_k = \sqrt{\frac{N}{L_k}} \sum_{l=1}^{L_k} \alpha_k^{(l)} \mathbf{a}(\varphi_k^{(l)}, \theta_k^{(l)}), \quad (2)$$

where L_k denotes the number of paths for user k , $\alpha_k^{(l)}$ and $\varphi_k^{(l)}$ ($\theta_k^{(l)}$) for $1 \leq l \leq L_k$ are the complex gain and azimuth (elevation) angle of departure (AoD) of the path l for user k , $\mathbf{a}(\varphi, \theta)$ presents the $N \times 1$ array steering vector. For the typical uniform planar array (UPA) with N_1 elements in horizon and N_2 elements in vertical ($N = N_1N_2$), we have [6]

$$\mathbf{a}(\varphi, \theta) = \mathbf{a}_{\text{az}}(\varphi) \otimes \mathbf{a}_{\text{el}}(\theta), \quad (3)$$

where $\mathbf{a}_{\text{az}}(\varphi) = \frac{1}{\sqrt{N_1}} [e^{j2\pi i(d_1/\lambda) \sin \varphi}]^T$ for $i \in \mathcal{I}(N_1)$, $\mathbf{a}_{\text{el}}(\theta) = \frac{1}{\sqrt{N_2}} [e^{j2\pi j(d_2/\lambda) \sin \theta}]^T$ for $j \in \mathcal{I}(N_2)$, $\mathcal{I}(n) = \{0, 1, \dots, n-1\}$, λ is the signal wavelength, and d_1 (d_2) is the horizontal (vertical) antenna spacing. At mmWave frequencies, we usually have $d_1 = d_2 = \lambda/2$ [14].

III. ENERGY EFFICIENT HYBRID PRECODING

In this section, we first describe the proposed SI-based hybrid precoding architecture. Then, we propose an ACE-based hybrid precoding scheme for this new architecture. Finally, the complexity analysis is provided.

A. The proposed SI-based hybrid precoding architecture

Fig. 1 (a) and (b) show the traditional precoding architectures, i.e., the one with finite-resolution phase shifters (PS-based architecture) [9] and the one with switches (SW-based architecture) [11], respectively.

As shown in Fig. 1 (a), the traditional PS-based architecture requires a complicated phase shifter network, and the corresponding energy consumption can be presented as

$$P_{\text{PS-architecture}} = \rho + N_{\text{RF}}P_{\text{RF}} + NN_{\text{RF}}P_{\text{PS}} + P_{\text{BB}}, \quad (4)$$

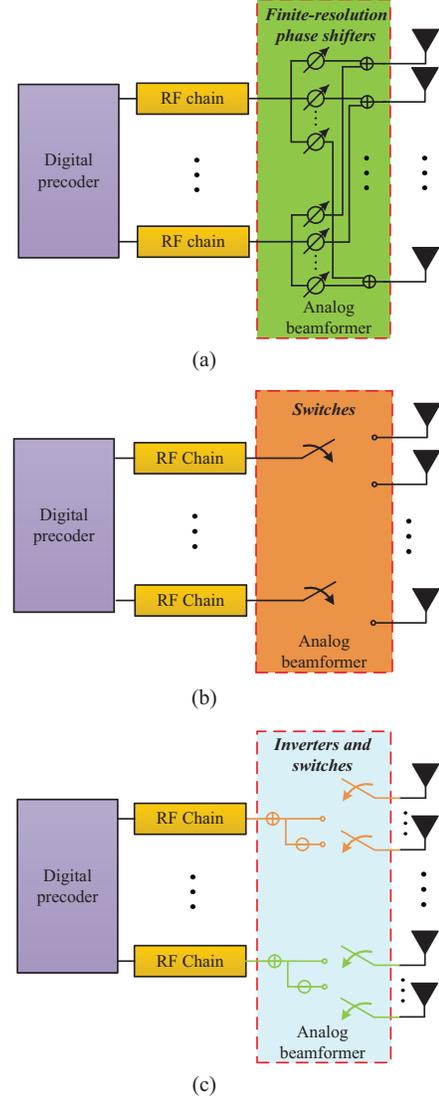


Fig. 1. Hybrid precoding: (a) traditional PS-based architecture; (b) traditional SW-based architecture; (c) proposed SI-based architecture.

where P_{RF} , P_{PS} , and P_{BB} are the energy consumption of RF chain, finite-resolution phase shifter, and baseband, respectively. Note that although the PS-based architecture enjoys high design freedom to achieve the near-optimal performance [9], it requires a large number (e.g., $N \times N_{\text{RF}} = 64 \times 16 = 1024$ [9]) of phase shifters. Moreover, the energy consumption of finite phase shifter is also considerable (e.g., $P_{\text{PS}} = 40\text{mW}$ for 4-bit phase shifter [11]). These make the traditional PS-based architecture still suffer from high energy consumption [14].

By contrast, as shown in Fig. 1 (b), SW-based architecture can efficiently relieve such problem by employing a small number (N_{RF} instead of $N \times N_{\text{RF}}$) of energy-efficient switches. The energy consumption of SW-based architecture can be presented as

$$P_{\text{SW-architecture}} = \rho + N_{\text{RF}}P_{\text{RF}} + N_{\text{RF}}P_{\text{SW}} + P_{\text{BB}}, \quad (5)$$

where P_{SW} is the energy consumption of switch, which is much lower than P_{PS} (e.g., $P_{\text{SW}} = 5\text{mW}$ [11]). Nevertheless,

since only N_{RF} antennas are active simultaneously, SW-based architecture cannot fully achieve the array gains of mmWave massive MIMO, leading to an obvious performance loss [15].

To overcome the problems faced by traditional architectures, we propose the SI-based architecture as shown in Fig. 1 (c), which can be considered as a better trade-off between the near-optimal PS-based architecture and the energy-efficient SW-based architecture. Specifically, in the proposed SI-based architecture, each RF chain is only connected to a sub antenna array with $M = N/N_{\text{RF}}$ (assumed to be an integer) antennas instead of all N antennas [16]. Moreover, each RF chain is connected to the sub antenna array via only one inverter and M switches instead of N phase shifters. The energy consumption of SI-based architecture can be presented by

$$P_{\text{SI-architecture}} = \rho + N_{\text{RF}} P_{\text{RF}} + N_{\text{RF}} P_{\text{IN}} + N P_{\text{SW}} + P_{\text{BB}}, \quad (6)$$

where P_{IN} is the energy consumption of inverter. It worth pointing out that the inverters can be realized by the digital chip with the energy consumption similar to switches (i.e., $P_{\text{IN}} \approx P_{\text{SW}}$) [15]. As a result, by comparing (4)-(6), we can conclude that the energy consumption of the proposed SI-based architecture is much lower than that of PS-based one. Furthermore, as all antennas are used, SI-based architecture can also achieve the potential array gains of mmWave massive MIMO, which will be further proved as follows.

To do this, we need to first explain the hardware constraints induced by the proposed SI-based architecture, which are different from those of the traditional ones. The first constraint is that the analog beamformer \mathbf{F}_{RF} should be a block diagonal matrix instead of a full matrix as

$$\mathbf{F}_{\text{RF}} = \begin{bmatrix} \bar{\mathbf{f}}_1^{\text{RF}} & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \bar{\mathbf{f}}_2^{\text{RF}} & & \mathbf{0} \\ \vdots & & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \dots & \bar{\mathbf{f}}_{N_{\text{RF}}}^{\text{RF}} \end{bmatrix}_{N \times N_{\text{RF}}}, \quad (7)$$

where $\bar{\mathbf{f}}_n^{\text{RF}}$ is the $M \times 1$ analog beamformer on the n th sub antenna array. The second constraint is that since only inverters and switches are used, all the N nonzero elements of \mathbf{F}_{RF} should belong to

$$\frac{1}{\sqrt{N}} \{-1, +1\}. \quad (8)$$

Based on these constraints, we have the following **Lemma 1**.

Lemma 1. *Assume that the channel \mathbf{h}_k of user k only has single path, i.e., $L_k = 1$ [12]. When $N \rightarrow \infty$ and $N/M = N_{\text{RF}}$, the ratio ζ between the array gains achieved by SI-based architecture and that achieved by PS-based architecture with sufficiently high-resolution phase shifters can be presented by*

$$\lim_{N \rightarrow \infty, \frac{N}{M} = N_{\text{RF}}} \zeta = \frac{4}{N_{\text{RF}} \pi^2}. \quad (9)$$

Proof: For the traditional PS-based architecture with sufficiently high-resolution phase shifters, the phases of the elements in the analog beamformer can be arbitrarily adjusted to capture the power of \mathbf{h}_k . Therefore, the array gains achieved

by PS-based architecture is $|\alpha_k^{(1)}|^2$. By contrast, the array gains achieved by SI-based architecture can be presented by

$$\begin{aligned} |\mathbf{h}_k^H \mathbf{f}_k^{\text{RF}}|^2 &= N |\alpha_k^{(1)}|^2 |\mathbf{a}^H(\varphi_k) \mathbf{f}_k^{\text{RF}}|^2 \\ &= \frac{1}{N} |\alpha_k^{(1)}|^2 \left| \sum_{m=1}^M e^{j\bar{\phi}_m} \right|^2 \\ &= \frac{1}{N} |\alpha_k^{(1)}|^2 \left(\left| \sum_{m=1}^M \cos(\bar{\phi}_m) \right|^2 + \left| \sum_{m=1}^M \sin(\bar{\phi}_m) \right|^2 \right), \end{aligned} \quad (10)$$

where \mathbf{f}_k^{RF} is the k th column of \mathbf{F}_{RF} including the zeros and $\bar{\phi}_m$ denotes the phase quantization error. Since the nonzero elements in \mathbf{f}_k^{RF} belong to $\frac{1}{\sqrt{N}} \{-1, +1\}$, $\bar{\phi}_m$ can be assumed to follow the uniform distribution $\mathcal{U}(-\pi/2, \pi/2)$ for $1 \leq m \leq M$ [10]. Then, we have

$$\begin{aligned} \lim_{\substack{N \rightarrow \infty \\ N/M = N_{\text{RF}}}} \zeta &= \frac{1}{N_{\text{RF}} M} \left\{ \left| \sum_{m=1}^M \cos(\bar{\phi}_m) \right|^2 + \left| \sum_{m=1}^M \sin(\bar{\phi}_m) \right|^2 \right\} \\ &= \frac{1}{N_{\text{RF}}} (\mathbb{E}[\cos(\bar{\phi}_m)])^2 + (\mathbb{E}[\sin(\bar{\phi}_m)])^2 \\ &= \frac{4}{N_{\text{RF}} \pi^2}, \end{aligned} \quad (11)$$

which verifies the conclusion in **Lemma 1**. \blacksquare

Lemma 1 indicates that although the proposed SI-based architecture suffers from some loss of array gains compared to the near-optimal PS-based architecture, the performance loss keeps constant and limited, which does not grow as the number of BS antennas goes to infinity. Recalling the low energy consumption of SI-based architecture, we can conclude that our scheme is a better trade-off between the traditional architectures, which will be also verified by simulation. Next, we will design a near-optimal hybrid precoding scheme for SI-based architecture with quite different hardware constraints.

B. ACE-based hybrid precoding scheme

We aim to design the analog beamformer \mathbf{F}_{RF} and the digital precoder \mathbf{F}_{BB} to maximize the achievable sum-rate R , which can be presented as

$$\begin{aligned} (\mathbf{F}_{\text{RF}}^{\text{opt}}, \mathbf{F}_{\text{BB}}^{\text{opt}}) &= \arg \max_{\mathbf{F}_{\text{RF}}, \mathbf{F}_{\text{BB}}} R, \\ \text{s.t. } & \mathbf{F}_{\text{RF}} \in \mathcal{F}, \\ & \|\mathbf{F}_{\text{RF}} \mathbf{F}_{\text{BB}}\|_F^2 = \rho, \end{aligned} \quad (12)$$

where \mathcal{F} denotes the set with all possible analog beamformers satisfying the two constraints (7) and (8) described above, and the achievable sum-rate R can be presented by

$$R = \sum_{k=1}^K \log_2(1 + \gamma_k), \quad (13)$$

where γ_k presents the signal-to-interference-plus-noise ratio (SINR) of the k th user as

$$\gamma_k = \frac{|\mathbf{h}_k^H \mathbf{F}_{\text{RF}} \mathbf{f}_k^{\text{BB}}|^2}{\sum_{k' \neq k} |\mathbf{h}_k^H \mathbf{F}_{\text{RF}} \mathbf{f}_{k'}^{\text{BB}}|^2 + \sigma^2}, \quad (14)$$

where \mathbf{f}_k^{BB} is the k th column of \mathbf{F}_{BB} .

It is worth pointing out that the constraints (7) and (8) on the analog beamformer \mathbf{F}_{RF} are non-convex. This makes (12) very difficult to be solved. Fortunately, as all the N nonzero elements of \mathbf{F}_{RF} belong to the set $\frac{1}{\sqrt{N}}\{-1, +1\}$, the number of possible \mathbf{F}_{RF} is finite. Therefore, (12) can be regarded as a non-coherent combining problem [13]. To solve it, we can first select a candidate \mathbf{F}_{RF} , and compute the optimal \mathbf{F}_{BB} according to the effective channel matrix $\mathbf{H}\mathbf{F}_{\text{RF}}$ without non-convex constraints. After all possible \mathbf{F}_{RF} 's have been searched, we can obtain the optimal analog beamformer $\mathbf{F}_{\text{RF}}^{\text{opt}}$ and digital precoder $\mathbf{F}_{\text{BB}}^{\text{opt}}$. However, such exhaustive search scheme requires to search 2^N possible \mathbf{F}_{RF} 's and \mathbf{F}_{BB} 's, which involves unaffordable complexity as N is usually large in mmWave massive MIMO systems (e.g., $N = 64$, $2^{64} \approx 1.84 \times 10^{19}$). To solve this problem, we propose an ACE algorithm, which can be considered as an improved version of the CE algorithm developed from machine learning [13].

At first, we would like to briefly introduce the conventional CE algorithm, which is a probabilistic model-based algorithm to solve the combining problem by an iterative procedure. In each iteration, the CE algorithm first generates S candidates (e.g., possible hybrid precoders in our problem) according to a probability distribution. Then, it computes the objective value (e.g., achievable sum-rate in our problem) of each candidate, and selects S_{elite} best candidates as "elite" [13]. Finally, based on the selected elites, the probability distribution is updated by minimizing the CE. Repeating such procedure, the probability distribution will be refined to generate a solution close to the optimal one with a sufficiently high probability. However, although the CE algorithm has been widely used in machine learning [13], it still has some disadvantages. One is that the contributions of all elites are treated as the same. Intuitively, the elite with better objective value should be more important when we update the probability distribution. Therefore, if we can adaptively weight the elites according to their objective values, better performance can be expected. Following this idea, we propose an ACE algorithm to solve (12).

The pseudo-code of the proposed ACE-based hybrid precoding scheme¹ is summarized in **Algorithm 1**, which can be explained as follows. At the beginning, we formulate the nonzero elements in \mathbf{F}_{RF} as $N \times 1$ vector $\mathbf{f} = [(\bar{\mathbf{f}}_1^{\text{RF}})^T, (\bar{\mathbf{f}}_2^{\text{RF}})^T, \dots, (\bar{\mathbf{f}}_{N_{\text{RF}}}^{\text{RF}})^T]^T$, and set the probability parameter $\mathbf{u} = [u_1, u_2, \dots, u_N]^T$ as an $N \times 1$ vector, where $0 \leq u_n \leq 1$ presents the probability that $f_n = 1/\sqrt{N}$, f_n is the n th element of \mathbf{f} . Then, by initializing $\mathbf{u}^{(0)} = \frac{1}{2} \times \mathbf{1}_{N \times 1}$ ($\mathbf{1}$ is the all-one vector), we assume that all the N nonzero elements of \mathbf{F}_{RF} belong to $\frac{1}{\sqrt{N}}\{-1, +1\}$ with equal probability, since no priori information is available. During the i th iteration, in step 1, we first generate S candidate analog beamformers $\{\mathbf{F}_{\text{RF}}^s\}_{s=1}^S$ based on the probability distribution $\Xi(\mathcal{F}; \mathbf{u}^{(i)})$ (i.e., generate $\{\mathbf{f}^s\}_{s=1}^S$ according to $\mathbf{u}^{(i)}$, and reshape them as matrices belong to \mathcal{F}). Then, in step 2, we

¹Note that the convergence of the proposed ACE-based hybrid precoding scheme can be proved by extending the Theorem 1 in [17].

Input: Channel matrix \mathbf{H} ; Number of iterations I ;
Number of candidates S ; Number of elites S_{elite} .
Initialization: $i = 0$; $\mathbf{u}^{(0)} = \frac{1}{2} \times \mathbf{1}_{N \times 1}$.
for $0 \leq i \leq I$
1. Randomly generate S candidate analog beamformers $\{\mathbf{F}_{\text{RF}}^s\}_{s=1}^S$ based on $\Xi(\mathcal{F}; \mathbf{u}^{(i)})$;
2. Compute S corresponding digital precoders $\{\mathbf{F}_{\text{BB}}^s\}_{s=1}^S$ based on the effective channel $\mathbf{H}_{\text{eq}}^s = \mathbf{H}\mathbf{F}_{\text{RF}}^s$;
3. Calculate the achievable sum-rate $\{R(\mathbf{F}_{\text{RF}}^s)\}_{s=1}^S$ (13);
4. Sort $\{R(\mathbf{F}_{\text{RF}}^s)\}_{s=1}^S$ in a descend order as
 $R(\mathbf{F}_{\text{RF}}^{[1]}) \geq R(\mathbf{F}_{\text{RF}}^{[2]}) \geq \dots \geq R(\mathbf{F}_{\text{RF}}^{[S]})$;
5. Select elites as $\mathbf{F}_{\text{RF}}^{[1]}, \mathbf{F}_{\text{RF}}^{[2]}, \dots, \mathbf{F}_{\text{RF}}^{[S_{\text{elite}}]}$;
6. Calculate weight w_s for each elite $\mathbf{F}_{\text{RF}}^{[s]}$, $1 \leq s \leq S_{\text{elite}}$;
7. Update $\mathbf{u}^{(i+1)}$ according to $\{w_s\}_{s=1}^{S_{\text{elite}}}$ and $\{\mathbf{F}_{\text{RF}}^{[s]}\}_{s=1}^{S_{\text{elite}}}$;
8. $i = i + 1$;
end for
Output: Analog beamformer $\mathbf{F}_{\text{RF}}^{[1]}$; Digital precoder $\mathbf{F}_{\text{BB}}^{[1]}$.

Algorithm 1: The proposed ACE-based hybrid precoding

calculate the corresponding digital precoder \mathbf{F}_{BB}^s according to the effective channel $\mathbf{H}_{\text{eq}}^s = \mathbf{H}\mathbf{F}_{\text{RF}}^s$ for $1 \leq s \leq S$. Note that there are lots of advanced digital precoder schemes [3]. In this paper, we adopt the classical ZF digital precoder with the near-optimal performance and low complexity as the example [3], and \mathbf{F}_{BB}^s can be computed as

$$\mathbf{G}^s = (\mathbf{H}_{\text{eq}}^s)^H (\mathbf{H}_{\text{eq}}^s (\mathbf{H}_{\text{eq}}^s)^H)^{-1}, \quad \mathbf{F}_{\text{BB}}^s = \beta^s \mathbf{G}^s, \quad (15)$$

where $\beta^s = \sqrt{\rho} / \|\mathbf{F}_{\text{RF}}^s \mathbf{G}^s\|_F$ is power normalized factor.

After that, in step 3, the achievable sum-rate $\{R(\mathbf{F}_{\text{RF}}^s)\}_{s=1}^S$ are calculated by substituting \mathbf{F}_{RF}^s and \mathbf{F}_{BB}^s (also a function of \mathbf{F}_{RF}^s) into (13). We sort $\{R(\mathbf{F}_{\text{RF}}^s)\}_{s=1}^S$ in a descend order in step 4. Then, the elites can be obtained in step 5. In the conventional CE algorithm, the next step is to using elites to update $\mathbf{u}^{(i+1)}$ by minimizing CE, which is equivalent to [13]

$$\mathbf{u}^{(i+1)} = \arg \max_{\mathbf{u}^{(i)}} \frac{1}{S} \sum_{s=1}^{S_{\text{elite}}} \ln \Xi(\mathbf{F}_{\text{RF}}^{[s]}; \mathbf{u}^{(i)}), \quad (16)$$

where $\Xi(\mathbf{F}_{\text{RF}}^{[s]}; \mathbf{u}^{(i)})$ denotes the probability to generate $\mathbf{F}_{\text{RF}}^{[s]}$. As mentioned above, in (16), the contributions of all elites are treated as the same, leading to performance degradation. To solve this problem, we propose to weight each elite adaptively based on its achievable sum-rate. Specifically, we first define an auxiliary parameter T presenting the average achievable sum-rate of all elites as

$$T = \frac{1}{S_{\text{elite}}} \sum_{s=1}^{S_{\text{elite}}} R(\mathbf{F}_{\text{RF}}^{[s]}). \quad (17)$$

Then, the weight w_s of the elite $\mathbf{F}_{\text{RF}}^{[s]}$ can be calculated in step 6 as $w_s = R(\mathbf{F}_{\text{RF}}^{[s]})/T$. Based on $\{w_s\}_{s=1}^{S_{\text{elite}}}$, (16) can be modified as

$$\mathbf{u}^{(i+1)} = \arg \max_{\mathbf{u}^{(i)}} \frac{1}{S} \sum_{s=1}^{S_{\text{elite}}} w_s \ln \Xi(\mathbf{F}_{\text{RF}}^{[s]}; \mathbf{u}^{(i)}). \quad (18)$$

Note that $\Xi(\mathbf{F}_{\text{RF}}^{[s]}; \mathbf{u}^{(i)}) = \Xi(\mathbf{f}^{[s]}; \mathbf{u}^{(i)})$, and the n th element $f_n^{[s]}$ of $\mathbf{f}^{[s]}$ is a Bernoulli random variable, where $f_n^{[s]} = 1/\sqrt{N}$ with probability $u_n^{(i)}$ and $f_n^{[s]} = -1/\sqrt{N}$ with probability $1 - u_n^{(i)}$. Therefore, we have

$$\Xi(\mathbf{F}_{\text{RF}}^{[s]}; \mathbf{u}^{(i)}) = \prod_{n=1}^N \left(u_n^{(i)}\right)^{\frac{1}{2}(1+\sqrt{N}f_n^{[s]})} \left(1-u_n^{(i)}\right)^{\frac{1}{2}(1-\sqrt{N}f_n^{[s]})}. \quad (19)$$

After substituting (19) into (18), the first derivative of the target in (18) with respect to $u_n^{(i)}$ can be derived as

$$\frac{1}{S} \sum_{s=1}^{S_{\text{elite}}} w_s \left(\frac{1 + \sqrt{N}f_n^{[s]}}{2u_n^{(i)}} - \frac{1 - \sqrt{N}f_n^{[s]}}{2(1 - u_n^{(i)})} \right). \quad (20)$$

Setting (20) to zero, $\mathbf{u}^{(i+1)}$ can be updated in step 7 as

$$u_n^{(i+1)} = \frac{\sum_{s=1}^{S_{\text{elite}}} w_s (\sqrt{N}f_n^{[s]} + 1)}{2 \sum_{s=1}^{S_{\text{elite}}} w_s}. \quad (21)$$

Such procedure above will be repeated ($i = i + 1$ in step 8) until the maximum number of iterations I is reached, and the analog beamformer and digital precoder will be selected as $\mathbf{F}_{\text{RF}}^{[1]}$ and $\mathbf{F}_{\text{BB}}^{[1]}$, respectively. Finally, it is worth pointing out that the proposed ACE-based hybrid precoding scheme can be also extended to the scenario where users employ multiple antennas. In this case, the analog beamformers at the BS and users should be jointly searched by the ACE algorithm.

C. Computational complexity analysis

In this subsection, the computational complexity of the proposed ACE-based hybrid precoding scheme is discussed.

From **Algorithm 1**, we can observe that the complexity of the ACE-based hybrid precoding scheme mainly comes from steps 2, 3, 6, and 7. In step 2, we need to compute S effective channel matrices $\{\mathbf{H}_{\text{eq}}^s\}_{s=1}^S$ and digital precoders $\{\mathbf{F}_{\text{BB}}^s\}_{s=1}^S$ according to (15). Therefore, the complexity of this part is $\mathcal{O}(SNK^2)$. In step 3, the achievable sum-rate of each candidate is computed. Since we employ the digital ZF precoder, the SINR γ_k^s of the k th user for the s th candidate is simplified to $\gamma_k^s = (\beta^s/\sigma)^2$. As a result, this part only involves the complexity $\mathcal{O}(S)$. In step 6, we calculate S_{elite} weights based on (17), which is quite simple with the complexity $\mathcal{O}(S_{\text{elite}})$. Finally, in step 7, the probability parameter $\mathbf{u}^{(i+1)}$ is updated according to (21) with the complexity $\mathcal{O}(NS_{\text{elite}})$.

In summary, after I iterations, the total computational complexity of the proposed ACE hybrid precoding scheme is $\mathcal{O}(ISNK^2)$. Since K is usually small, I and S also do not have to be very large as will be verified in the next section, we can conclude that the complexity of the proposed ACE-based hybrid precoding scheme is acceptable, which is comparable with the simple least squares (LS) algorithm.

IV. SIMULATION RESULTS

In this section, we provide the simulation results in terms of achievable sum-rate and energy-efficiency to evaluate the performance of the proposed ACE-based hybrid precoding

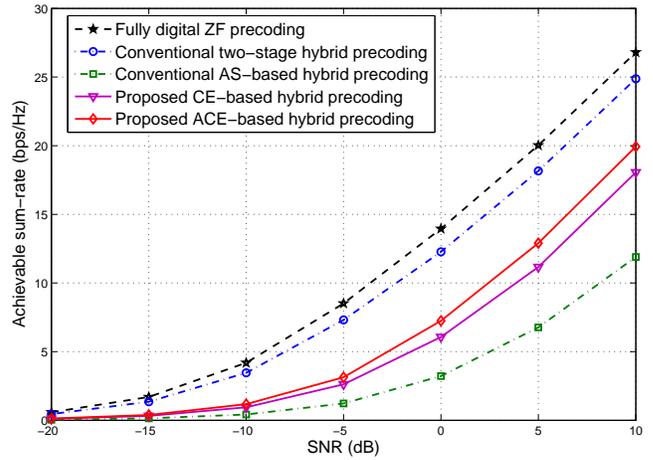


Fig. 2. Achievable sum-rate comparison.

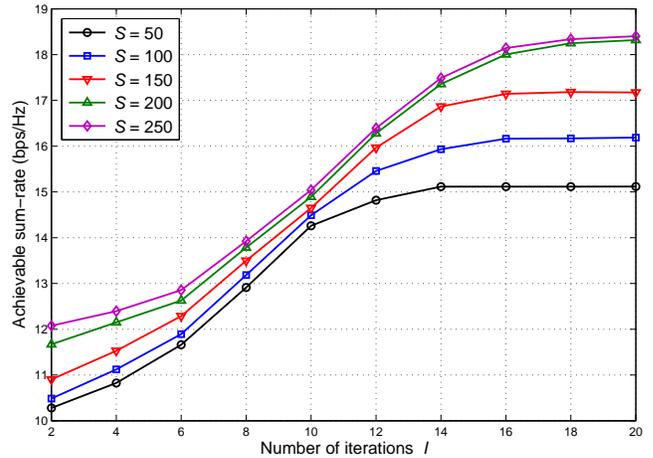


Fig. 3. Achievable sum-rate against S and I .

scheme. The simulation parameters are described as follows: We assume that the BS employs an UPA with antenna spacing $d_1 = d_2 = \lambda/2$. For the k th user, we generate the channel \mathbf{h}_k based on (2), where we assume: 1) $L_k = 3$; 2) $\alpha_k^{(l)} \sim \mathcal{CN}(0, 1)$ for $1 \leq l \leq L_k$; 3) $\varphi_k^{(l)}$ and $\theta_k^{(l)}$ follow the uniform distribution $\mathcal{U}(-\pi, \pi)$ for $1 \leq l \leq L_k$ [12]. Finally, the signal-to-noise ratio (SNR) is defined as ρ/σ^2 .

Fig. 2 shows the achievable sum-rate comparison in a typical mmWave massive MIMO system with $N = 64$, $N_{\text{RF}} = K = 4$. In Fig. 2, the proposed CE-based (i.e., using the conventional CE algorithm to solve (12)) and ACE-based hybrid precoding schemes are designed for SI-based architecture, where we set $S = 200$, $S_{\text{elite}} = 40$, and $I = 20$ for **Algorithm 1**, the conventional two-stage hybrid precoding scheme is designed for PS-based architecture with 4-bit phase shifters [9], and the conventional antenna selection (AS)-based hybrid precoding scheme is designed for SW-based architecture [11] with switches. Firstly, we can observe from Fig. 2 that the proposed ACE algorithm outperforms the traditional CE algorithm, where the SNR gap is about 1 dB. Note that the ACE algorithm only involves one additional step (i.e., step 6 in **Algorithm 1**) with negligible complexity. Therefore, the proposed ACE algorithm is more efficient. Moreover, Fig. 2 also shows that the proposed ACE-based hybrid precoding can

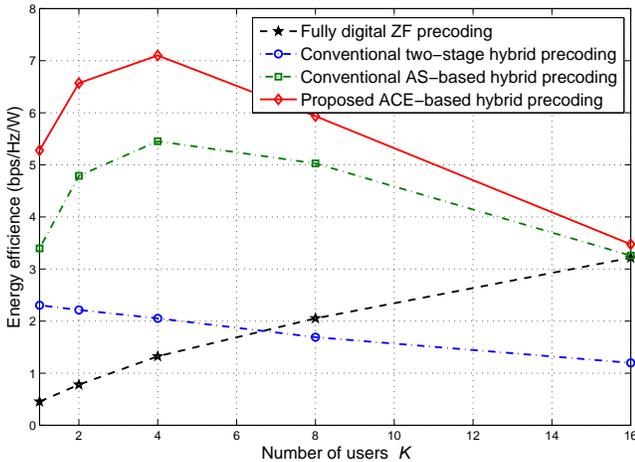


Fig. 4. Energy efficiency comparison.

achieve the sum-rate much higher than the conventional AS-based hybrid precoding, as it can achieve the potential array gains in mmWave massive MIMO systems. Finally, we can observe that the performance gap between ACE-based hybrid precoding and two-stage hybrid precoding is limited and keeps constant, which further verifies the conclusion in **Lemma 1**.

Fig. 3 shows the achievable sum-rate of the proposed ACE-based hybrid precoding against the number of candidates S and the number of iterations I , when $S_{\text{elite}}/S = 0.2$, $N = 64$, $N_{\text{RF}} = K = 4$, and $\text{SNR} = 10$ dB. From Fig. 3, we can observe that when S is small, increasing S will lead to an obvious improvement in the sum-rate performance. However, when S is sufficiently large, such trend is no longer obvious. This indicates that the number of candidates S does not have to be very large, e.g., $S = 200$ is enough. Furthermore, Fig. 3 also indicates that the proposed ACE-based hybrid precoding can converge with a small number of iterations, e.g., $I = 20$. These observations verify the rationality of the parameters for **Algorithm 1** we used in Fig. 2.

Fig. 4 shows the energy efficiency comparison when $N = 64$ is fixed and $N_{\text{RF}} = K$ varies from 1 to 16. The parameters for **Algorithm 1** are the same as Fig. 2. According to [15], [16], the energy efficiency can be defined as the ratio between the achievable sum-rate and the energy consumption, which should be (4), (5), and (6) for two-stage hybrid precoding, AS-based hybrid precoding, and ACE-based hybrid precoding, respectively. In this paper, we adopt the practical values $\rho = 30\text{mW}$ [16], $P_{\text{RF}} = 300\text{mW}$ [16], $P_{\text{BB}} = 200\text{mW}$ [15], $P_{\text{PS}} = 40\text{mW}$ (4-bit phase shifter) [15], and $P_{\text{SW}} = P_{\text{IN}} = 5\text{mW}$ [15]. From Fig. 4, we can observe that the proposed ACE-based hybrid precoding with SI-based architecture can achieve much higher energy efficiency than the others, especially when K is not very large (e.g., $K \leq 12$). Furthermore, it is interesting to observe that when $K \geq 8$, the energy efficiency of the two-stage hybrid precoding with PS-based architecture is even lower than that of the fully digital ZF precoding. This is due to the fact that as K grows, the number of phase shifters in PS-based architecture increases rapidly. As a result, the energy consumption of the phase shifter network will be huge, even higher than that of RF chains.

V. CONCLUSIONS

In this paper, we propose an energy-efficient SI-based hybrid precoding architecture, where the analog part is realized by a small number of switches and inverters. The performance analysis proves that the performance gap between the proposed SI-based architecture and the traditional near-optimal one keeps constant and limited. Then, by employing the idea of CE optimization in machine learning, we further propose an ACE-based hybrid precoding scheme with low complexity for SI-based architecture. Simulation results verify that our scheme can achieve the satisfying sum-rate performance and much higher energy efficiency than traditional schemes.

REFERENCES

- [1] A. L. Swindlehurst, E. Ayanoglu, P. Heydari, and F. Capolino, "Millimeter-wave massive MIMO: The next wireless revolution?" *IEEE Commun. Mag.*, vol. 52, no. 9, pp. 56–62, Sep. 2014.
- [2] T. S. Rappaport, S. Sun, R. Mayzus, H. Zhao, Y. Azar, K. Wang, G. N. Wong, J. K. Schulz, M. Samimi, and F. Gutierrez, "Millimeter wave mobile communications for 5G cellular: It will work!" *IEEE Access*, vol. 1, pp. 335–349, May 2013.
- [3] T. L. Marzetta, "Noncooperative cellular wireless with unlimited numbers of base station antennas," *IEEE Trans. Wireless Commun.*, vol. 9, no. 11, pp. 3590–3600, Nov. 2010.
- [4] H. Xie, F. Gao, S. Zhang, and S. Jin, "A unified transmission strategy for TDD/FDD massive MIMO systems with spatial basis expansion model," to appear in *IEEE Trans. Veh. Technol.*, 2017.
- [5] R. W. Heath, N. Gonzalez-Prelcic, S. Rangan, W. Roh, and A. Sayeed, "An overview of signal processing techniques for millimeter wave MIMO systems," *IEEE J. Sel. Top. Signal Process.*, vol. 10, no. 3, pp. 436–453, Apr. 2016.
- [6] O. El Ayach, S. Rajagopal, S. Abu-Surra, Z. Pi, and R. W. Heath, "Spatially sparse precoding in millimeter wave MIMO systems," *IEEE Trans. Wireless Commun.*, vol. 13, no. 3, pp. 1499–1513, Mar. 2014.
- [7] H. Xie, F. Gao, and S. Jin, "An overview of low-rank channel estimation for massive MIMO systems," *IEEE Access*, vol. 4, pp. 7313–7321, Nov. 2016.
- [8] F. Sohrabi and W. Yu, "Hybrid beamforming with finite-resolution phase shifters for large-scale MIMO systems," in *Proc. IEEE SPAWC Workshops*, Jul. 2015, pp. 136–140.
- [9] A. Alkhateeb, G. Leus, and R. W. Heath, "Limited feedback hybrid precoding for multi-user millimeter wave systems," *IEEE Trans. Wireless Commun.*, vol. 14, no. 11, pp. 6481–6494, Nov. 2015.
- [10] A. Alkhateeb, Y.-H. Nam, J. Zhang, and R. W. Heath, "Massive MIMO combining with switches," *IEEE Wireless Commun. Lett.*, vol. 5, no. 3, pp. 232–235, Jun. 2016.
- [11] R. Mendez-Rial, C. Rusu, A. Alkhateeb, N. González-Prelcic, and R. W. Heath, "Channel estimation and hybrid combining for mmWave: Phase shifters or switches?" in *Proc. ITA Workshops*, Feb. 2015, pp. 90–97.
- [12] A. Sayeed and J. Brady, "Beamspace MIMO for high-dimensional multiuser communication at millimeter-wave frequencies," in *Proc. IEEE GLOBECOM*, Dec. 2013, pp. 3679–3684.
- [13] R. Y. Rubinfeld and D. P. Kroese, *The cross-entropy method: A unified approach to combinatorial optimization, Monte-Carlo simulation and machine learning*. Springer Science & Business Media, 2013.
- [14] S. Han, C.-L. I, Z. Xu, and C. Rowell, "Large-scale antenna systems with hybrid precoding analog and digital beamforming for millimeter wave 5G," *IEEE Commun. Mag.*, vol. 53, no. 1, pp. 186–194, Jan. 2015.
- [15] R. Mendez-Rial, C. Rusu, N. González-Prelcic, A. Alkhateeb, and R. W. Heath, "Hybrid MIMO architectures for millimeter wave communications: Phase shifters or switches?" *IEEE Access*, vol. 4, pp. 247–267, Jan. 2016.
- [16] X. Gao, L. Dai, S. Han, C.-L. I, and R. W. Heath, "Energy-efficient hybrid analog and digital precoding for mmWave MIMO systems with large antenna arrays," *IEEE J. Sel. Areas Commun.*, vol. 34, no. 4, pp. 998–1009, Apr. 2016.
- [17] J.-C. Chen, C.-K. Wen, and K.-K. Wong, "An efficient sensor node selection algorithm for sidelobe control in collaborative beamforming," *IEEE Trans. Veh. Technol.*, vol. 65, no. 8, pp. 5984–5994, Aug. 2016.