

Fully- / Partially-Connected Hybrid Beamforming Architectures for mmWave MU-MIMO

Xiaoshen Song, *Student Member, IEEE*, Thomas Kühne, and Giuseppe Caire
Fellow, IEEE

Abstract

Hybrid digital analog (HDA) beamforming has attracted considerable attention in practical implementation of millimeter wave (mmWave) multiuser multiple-input multiple-output (MU-MIMO) systems due to the low power consumption with respect to its fully digital baseband counterpart. The implementation cost, performance, and power efficiency of HDA beamforming depends on the level of connectivity and reconfigurability of the analog beamforming network. In this paper, we investigate the performance of two typical architectures that can be regarded as extreme cases, namely, the fully-connected (FC) and the one-stream-per-subarray (OSPS) architectures. In the FC architecture each RF antenna port is connected to all antenna elements of the array, while in the OSPS architecture the RF antenna ports are connected to disjoint subarrays. We jointly consider the initial beam acquisition and data communication phases, such that the latter takes place by using the beam direction information obtained by the former. We use the state-of-the-art beam alignment (BA) scheme previously proposed by the authors and consider a family of MU-MIMO precoding schemes well adapted to the beam information extracted from the BA phase. We also evaluate the power efficiency of the two HDA architectures taking into account the power dissipation at different hardware components as well as the power backoff under typical power amplifier constraints. Numerical results show that the two architectures achieve similar sum spectral efficiency, while the OSPS architecture is advantageous with respect to the FC case in terms of hardware complexity and power efficiency, at the sole cost of a slightly longer BA time-to-acquisition due to its reduced beam angle resolution.

Index Terms

X. Song is sponsored by the China Scholarship Council (201604910530). This work was funded by the European Union's Horizon 2020 research and innovation programme under grant agreement No. 779305 (SERENA).

Millimeter Waves, MU-MIMO, HDA Beamforming, Beam Acquisition, Spectral Efficiency, Power Efficiency.

I. INTRODUCTION

Millimeter wave (mmWave) multiuser multiple-input multiple-output (MU-MIMO) communications have emerged as one of the most promising techniques for the second phase of 5G wireless systems, aimed at achieving broadband data communications at unprecedented high rates (≥ 1 Gb/s per user), in very dense urban small-cell environments. The relatively underutilized mmWave spectrum (30-300 GHz) allows to achieve a target ~ 1 Gb/s per data stream with ~ 1 GHz signal bandwidth, provided that the system can support a spectral efficiency of about 1 bit/s/Hz. Such relatively low spectral efficiency per stream can be achieved with rather standard modulation and coding techniques (e.g., binary codes of rate 1/2 mapped onto a QPSK constellation), when that the signal to interference plus noise ratio (SINR) at the receiver is between 0 and 3 dB (depending on the gap to capacity of the underlying code).¹

Due to the severe isotropic pathloss incurred by mmWave frequencies, large antenna gains are required both at the base station (BS) side and the user equipment (UE) side. Fortunately, the small carrier wavelength associated with mmWave frequencies enables large antenna arrays to be packed in a small form factor, such that the required large antenna gain can be obtained using beamforming. For example, in a single-user scenario where the signal-to-noise ratio (SNR) at the receiver in isotropic propagation conditions² is between -30 and -20 dB (a quite realistic situation for outdoor mmWave channels), a combined Tx and Rx beamforming gain of 30 dB is needed such that, when the Tx and the Rx beams are well aligned, the resulting SNR *after beamforming* reaches the desired target (a bit above 0 dB, as argued before).

Realizing fast and accurate digitally steerable beamforming at mmWave, however, is not a trivial task. One main challenge is that the conventional full digital transceiver architecture (with one radio frequency (RF) chain per antenna element) is infeasible due to hardware cost, power

¹With ideal single-user capacity achieving codes for the Gaussian channel, we have that $\log(1 + \text{SINR}) = 1$ bit/s/Hz is achieved for $\text{SINR} = 1$ (i.e., 0 dB). In practice, gaps of a fraction of a dB to 3-4 dB are obtained by actual coding schemes adopted in current standards.

² Here the isotropic propagation conditions correspond to one active antenna at the transmitter (Tx) and one active antenna at the receiver (Rx), respectively.

consumption, and above all power dissipation in the small integrated array form factor. Each RF chain consists of (roughly speaking) analog-to-digital/ digital-to-analog (A/D, D/A) converters, up/down-conversion mixers, filters, power amplifiers (PAs), and low-noise amplifiers (LNAs). It follows that a design goal for mmWave transceivers is to reduce the number of RF chains to be significantly smaller than the number of antenna array elements.

For this reason, the concatenation of digital and analog beamforming, known as hybrid digital analog (HDA) beamforming architecture, has been widely considered. In such a context, the limited number of RF chains are used to enable the multistream baseband processing, while an analog processing is used to realize the antenna beamforming gain. A primary objective of HDA beamforming is to maximize the multiuser sum rate, while keeping the hardware costs, complexity, and power efficiency, within some desirable targets.

A. Related Work

A large number of works have addressed HDA beamforming for mmWave communication systems. Rather than giving a complete account of such considerable body of literature (out of scope of the present non-tutorial paper), we consider a few significant representatives and examine their proposed approaches in a critical manner. A common assumption in most of existing works is that the analog part of the HDA precoder can only utilize phase control. This phase control can be realized through either phase shifters [1–6] or lenses [7, 8]. Consequently, the problem of finding the (sub-) optimal analog and digital precoding matrices is transformed into a series of relatively complicated decomposition steps [2–6], since the underlying optimization problem is non-convex. This phase-only control assumption may somewhat reduce the hardware complexity. However, the signaling freedom is also drastically reduced and the corresponding optimization computational complexity is typically prohibitive for practical real-time implementations. These drawbacks motivate the exploration of an analog precoding architecture with both phase and amplitude controls [9, 10]. In fact, it has been demonstrated in practice that simultaneous phase and amplitude control is fully feasible at mmWaves with good accuracy, low complexity, and low cost [11, 12].

Another severe limitation appearing in several HDA beamforming works is the assumption of invariant instantaneous channel coefficients over a large time duration [1, 13, 14]. It is known that, in order to overcome the heavy signal attenuation, communication at mmWaves requires an initial beam acquisition (which we refer as *beam alignment (BA)*) [7, 15, 16]. The goal

of BA is to find a pair of narrow beams connecting each UE with the BS.³ Thus, the nearly invariant channel assumption only makes sense *after BA is achieved*, since once the beams are aligned, the communication occurs only through a single narrow path with small effective angular spread, whose delay and Doppler shift can be easily compensated using standard synchronization techniques [17–19]. However, before BA is achieved, the channel delay spread and time-variation can be large due to the presence of several multipath components, each with its own delay and Doppler shift. In this case, the instantaneous channel coefficients change very fast. Any BA algorithms relying on an invariant instantaneous channel assumption are no-longer feasible, since for example, even a small motion of a few centimeters traverses several wavelengths, potentially producing multiple deep fades [20, 21].

In addition, a large number of works on HDA architectures investigated only the data communication phase and assume full channel state information (CSI) [2–6, 10, 22, 23], i.e., that the vectors of baseband complex channel coefficients at each array element are known. These works focus on the optimization of the HDA precoder using the full CSI knowledge. Unfortunately, this assumption is obviously not feasible in a realistic system. In order to acquire such coefficients, one should be able to sample each antenna element, i.e., one would need an RF chain per antenna element or exhaustively measure all elements successively. Hence, if full CSI knowledge was possible, no HDA beamforming would be needed, since we could simply implement baseband digital beamforming/multiuser precoding, which is performance-wise more efficient. As a matter of fact, it makes sense to study HDA architectures under the assumption that only a low-dimensional projection of the channel vectors can be measured by the limited number of RF chains. To this end, a hybrid precoding scheme exploiting implicit CSI (i.e., the couplings of all possible pairs of analog beamforming vectors) was proposed in [24]. However, the work in [24] (as well as in [4, 6, 10, 22, 23]) is limited to a single-user configuration and does not treat the MU-MIMO case.

It is known that MU-MIMO is superior to single-user beamforming from a network spectral efficiency perspective even under HDA, provided that the user density is rich enough such that the BS can schedule subsets of UEs to be served by spatial multiplexing with sufficient angular separation [25, 26]. Hence, this motivates us to consider the implementation of MU-MIMO schemes under realistic HDA architecture constraints. Two “extreme” HDA architectures are

³E.g., in line-of-sight (LOS) propagation, the aligned directions typically coincide with the AoA and AoD of the LoS path.

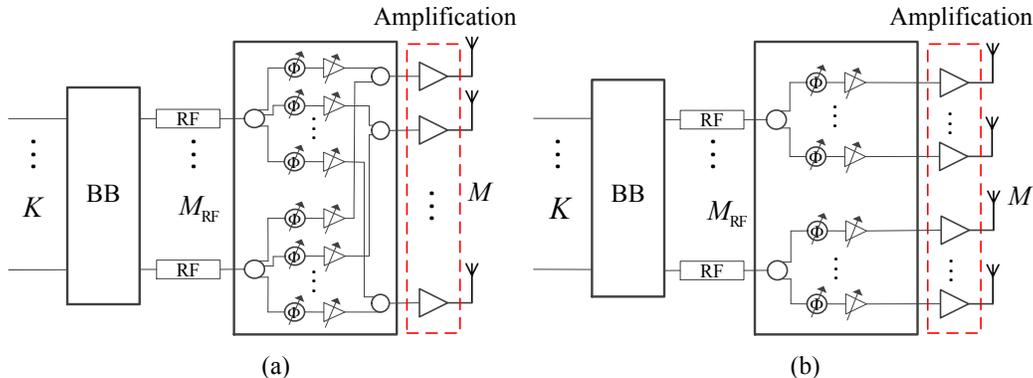


Fig. 1: Hybrid digital analog (HDA) transmitter architectures: (a) fully-connected (FC), (b) partially-connected with one-stream-per-subarray (OSPS). The “BB” block denotes digital baseband beamforming, K is the number of data streams, M_{RF} is the number of RF chains, and M is the number of antennas.

depicted in Fig. 1 [27]. Fig. 1 (a) shows a fully-connected (FC) architecture, where each RF antenna port is connected to all antenna elements of the array. At the other extreme, Fig. 1 (b) shows what we refer to as the one-stream-per-subarray (OSPS) architecture, where each RF antenna port is connected to a disjoint subarray. A common theme that underlies most of the HDA works is that the FC architecture outperforms the OSPS architecture only at the cost of higher hardware complexity. However, many reference works [3, 8, 10, 22, 23] ignore hardware impairments [6], such as the power dissipation and the PA nonlinear distortion. In particular, the nonlinear PAs employed at the BS can drastically distort the transmit signal when operated close to saturation [28]. To this end, a certain power backoff from the saturation power of a PA should be considered accordingly for different signaling schemes and transceiver architectures, such that the PAs can always work in their linear operating region.

B. Contributions

In this paper we overcome the shortcomings of the present literature outlined before, and comprehensively evaluate the performance of HDA architectures (in particular, as shown in Fig. 1), where we assume both amplitude and phase control for each analog path. Our main focus is on the MU-MIMO downlink, but similar and symmetric conclusions can be reached for the uplink as well. Our main contributions are summarized as follows:

1) *More general and realistic mmWave channel model.* We consider a quite general mmWave wireless channel model, taking into account the fundamental features of mmWave channels such as fast time-variation due to Doppler, frequency-selectivity, and the AoA-AoD sparsity

[20, 21, 29]. The numerical results based on our proposed channel model are further verified on the 3D geometry based channel generator QuaDRiGa [30], which has become a standard tool in industrial R&D as well as in 3GPP standardization.

2) *More practical hardware impairments and power efficiency analysis.* When comparing the HDA beamforming performance of different transmitter architectures, we take into account the practical hardware impairments, particularly, the potential power dissipation of the underlying analog network components, as well as the unavoidable power backoff for the nonlinear PAs. While the former is not difficult to be compensated, the latter is highly dependent on the peak-to-average power ratio (PAPR) of the input signal, which (as illustrated in Section V) should be carefully investigated in terms of different signaling and modulating schemes. On top of the potential hardware impairments, we also evaluate the power efficiency of the most power consuming PAs with respect to different transmitter architectures. Numerical results show that the OSPA architecture with single-carrier (SC) modulation achieves the highest power efficiency.

3) *A joint evaluation of initial BA and data communication.* As mentioned before, a main limitation in most hybrid beamforming works is that they only focus on the data communication and assume full CSI. To address this issue, we consider both initial BA and consecutive data communication in this paper. We assume that the precoder in the data communication phase can only exploit a limited amount of CSI, which is obtained along the beams acquired in the BA phase. Hence, the signaling and communication procedure in our paper captures the fundamental features of practical mmWave communication.

4) *Low-complexity data transmit precoding.* In the BA phase, we use our previously proposed BA scheme [16, 18, 19], after which each UE obtains a sparse estimate of the channel gains associated to all pairs of AoA-AoD on a finely spaced discrete grid, corresponding to the Tx and Rx beamforming codebooks. For the data communication phase, we consider three alternative precoding options on top of the effective channel after the BA phase. These are referred to as beam steering (BST), analog maximum ratio transmission (MRT), and joint analog maximum ratio and baseband zeroforcing (MR-ZF), respectively. The proposed schemes are very suitable for practical implementations due to the low-time-overhead and low-complexity. In particular, the MR-ZF precoding scheme proposed in this paper outperforms the state-of-the-art counterparts in the literature.

Notation: We denote vectors, matrices, and scalars by \mathbf{a} , \mathbf{A} , and a (A), respectively. For an integer $K \in \mathbb{Z}$, $[K]$ denotes the index set $\{1, \dots, K\}$. We represent sets by calligraphic \mathcal{A} and

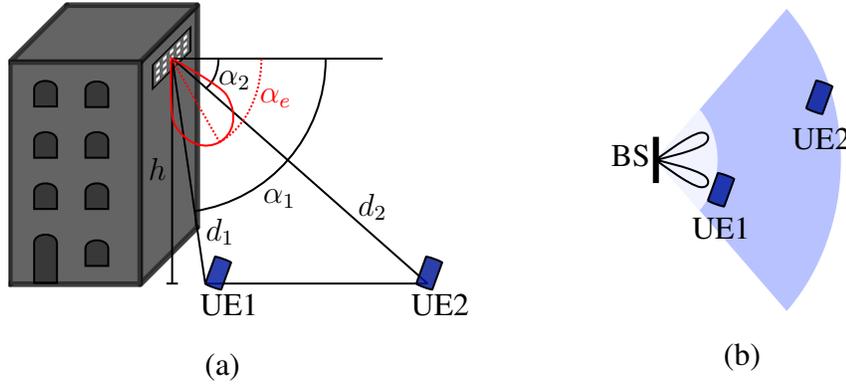


Fig. 2: Illustration of a small cell scenario with (a) 3D side view and (b) 2D top view. In this paper, the initial beam alignment refers to the beam training/searching in the azimuth plane as shown in (b).

their cardinality with $|\mathcal{A}|$. We use $\mathbb{E}[\cdot]$ for the expectation, $\|\cdot\|$ for l_2 -norm, \otimes for continuous-time convolution, \otimes for the Kronecker product, \odot for Hadamard product.

II. CHANNEL AND SIGNAL MODELS

A. Channel Model

One of the main new features of 5G wireless networks is the densely spread small cell layer [31]. In small cell configurations as illustrated in Fig. 2 (a), the BS creates a fixed arc-like sectorized beam in the elevation direction. The orientation of the BS beam center in the elevation direction tends to be fixed with an elevation angle α_e [32]. It follows that the probing area in the range direction is restricted and the intensive initial beam searching takes place mainly in the azimuth direction. For notation simplicity, in this paper we only focus on the 2D azimuth plane. Extension to the 3D geometry is conceptually straightforward although may lead to a rather high dimensional search for the initial beam acquisition phase. In the small cell scenario as illustrated in Fig. 2, where the beam shape in the elevation direction is fixed a priori in order to define the cell footprint area, the 2D azimuth geometry is fully justified. We assume that the BS serving simultaneously K UEs. The BS is equipped with a uniform linear array (ULA) of M antennas and M_{RF} RF chains, where $K \leq M_{\text{RF}} \ll M$. Each UE is equipped with a ULA of N antennas and $N_{\text{RF}} \ll N$ RF chains. Since the focus of this paper is the BS architecture, we consider the case of $N_{\text{RF}} = 1$, where the extension to $N_{\text{RF}} > 1$ is straightforward and was considered in our work on BA [16, 18, 19]. The propagation channel between the BS and the k -th UE, $k \in [K]$, consists of $L_k \ll \max\{M, N\}$ *significant* multipath components. As a result,

the $N \times M$ baseband equivalent impulse response of the channel at time slot s can be written as

$$\begin{aligned} \mathbf{H}_{s,k}(t, \tau) &= \sum_{l=1}^{L_k} \rho_{s,k,l} e^{j2\pi\nu_{k,l}t} \mathbf{a}_R(\phi_{k,l}) \mathbf{a}_T(\theta_{k,l})^H \delta(\tau - \tau_{k,l}) \\ &= \sum_{l=1}^{L_k} \mathbf{H}_{s,k,l}(t) \delta(\tau - \tau_{k,l}), \end{aligned} \quad (1)$$

where $\mathbf{H}_{s,k,l}(t) := \rho_{s,k,l} e^{j2\pi\nu_{k,l}t} \mathbf{a}_R(\phi_{k,l}) \mathbf{a}_T(\theta_{k,l})^H$ and $\delta(\cdot)$ denotes the Dirac delta function. Each l -th multipath component is identified by the tuple $(\phi_{k,l}, \theta_{k,l}, \tau_{k,l}, \nu_{k,l})$ of angle of arrival (AoA), angle of departure (AoD), delay, and Doppler shift, respectively. The vectors $\mathbf{a}_T(\theta_{k,l}) \in \mathbb{C}^M$ and $\mathbf{a}_R(\phi_{k,l}) \in \mathbb{C}^N$ are the array response vectors of the BS and the k -th UE at the AoD $\theta_{k,l}$ and the AoA $\phi_{k,l}$, respectively. With the ULA configuration and the assumption that the spacing of the ULA antennas in each array (subarray) equals to a half-wavelength $\lambda/2$, the elements of $\mathbf{a}_T(\theta_{k,l})$ and $\mathbf{a}_R(\phi_{k,l})$ are given by

$$[\mathbf{a}_T(\theta)]_{(i'-1)\cdot\hat{M}+d} = e^{j(d-1)\pi \sin(\theta)} \cdot e^{j\Psi(i',\theta)}, d \in [\hat{M}] \quad (2a)$$

$$[\mathbf{a}_R(\phi)]_n = e^{j(n-1)\pi \sin(\phi)}, n \in [N], \quad (2b)$$

where in (2a) we assume that $(i' \equiv 1, \hat{M} = M)$ for the FC architecture as shown in Fig. 1(a), and $(i' \in [M_{\text{RF}}], \hat{M} = \frac{M}{M_{\text{RF}}})$ for the OSPS architecture as shown in Fig. 1(b). The additional term $\Psi(i', \theta)$ in (2a) takes into account the phase shifts among different subarrays, given by

$$\Psi(i', \theta) = \frac{2\pi}{\lambda} (i' - 1) \cdot D_x \cdot \sin(\theta), \quad (3)$$

where i' indicates the index of the subarrays and $D_x \geq 0$ denotes the subarray center-to-center spacing in the scan direction. Hence, in the special case with $D_x = 0$, all the subarrays are co-located⁴; while with $D_x = \frac{M}{M_{\text{RF}}} \cdot \frac{\lambda}{2}$, the antenna element layout in the scan direction for the OSPS architecture is exactly the same as for the FC architecture.

For the sake of modeling simplicity, we assume in (1) that each multipath component has a very narrow footprint over the AoA-AoD-delay domain. The extension to more widely spread multipath clusters is straightforward and will be applied in the numerical simulations. We

⁴ In this paper, we consider a 2D geometry w.r.t. the azimuth plane as illustrated in Fig. 2(b). In practice, the co-located layout can be obtained by stacking the arrays on top of each other in the vertical dimension. Strictly speaking this yields a rectangular array configuration, but since each row forms an individually driven array, adaptive beamforming in the elevation direction is not possible, therefore the beamforming geometry is still two-dimensional.

adopt a block fading model, where the coefficient of the l -th multipath component $\rho_{s,k,l}$ is constant over a short interval (within one slot) and changes from slot to slot according to a wide-sense stationary process statistics characterized by its power spectral density (Doppler spectrum) [33]. When the channel *coherence time* (related to the inverse of the bandwidth of the Doppler spectrum, see [33]) is significantly larger than the slot duration but equal or smaller than the (non-consecutive) slot separation in time, a convenient model is to consider the coefficients as i.i.d. across different slots. Moreover, the Doppler shift $\nu_{k,l}$ as defined in (1) introduces a continuous phase rotation for each channel sample. Each multipath component (channel tap coefficient) is formed by the superposition of a large number of micro-scattering components (e.g., due to rough surfaces) having (approximately) the same AoA-AoD and delay. By the central limit theorem, it is customary to model the superposition of these many small effects as Gaussian [34, 35]. Hence, the multipath component coefficients can be modeled as Rice fading given by

$$\rho_{s,k,l} \sim \sqrt{\gamma_{k,l}} \left(\sqrt{\frac{\eta_{k,l}}{1 + \eta_{k,l}}} + \frac{1}{\sqrt{1 + \eta_{k,l}}} \check{\rho}_{s,k,l} \right), \quad (4)$$

where $\gamma_{k,l}$ denotes the overall multipath component strength, $\eta_{k,l} \in [0, \infty)$ indicates the strength ratio between the specular reflection (or LOS) and the scattered components, and $\check{\rho}_{s,k,l} \sim \mathcal{CN}(0, 1)$ is a zero-mean unit-variance complex Gaussian random variable whose value changes in an i.i.d. fashion across different slots. In particular, $\eta_{k,l} \rightarrow \infty$ indicates a pure LOS path while $\eta_{k,l} = 0$ indicates a pure scattered path, affected by Rayleigh fading.

The AoA-AoDs $(\phi_{k,l}, \theta_{k,l})$ in (1) can take on arbitrary values in the continuous AoA-AoD domain. Following the widely used approach of [36], known as *beam-domain representation*, we obtain a finite-dimensional representation of the channel response (1). More precisely, we consider the discrete set of AoA-AoDs

$$\Phi := \left\{ \check{\phi} : (1 + \sin(\check{\phi}))/2 = \frac{n-1}{N}, n \in [N] \right\}, \quad (5a)$$

$$\Theta := \left\{ \check{\theta} : (1 + \sin(\check{\theta}))/2 = \frac{m-1}{M}, m \in [M] \right\}. \quad (5b)$$

It follows that the corresponding sets $\mathcal{A}_R := \{\mathbf{a}_R(\check{\phi}) : \check{\phi} \in \Phi\}$ and $\mathcal{A}_T := \{\mathbf{a}_T(\check{\theta}) : \check{\theta} \in \Theta\}$ form discrete dictionaries to represent the channel response. For the ULAs considered in this

paper, the dictionaries \mathcal{A}_R and \mathcal{A}_T , after suitable normalization, reduce to the columns of unitary *Discrete Fourier Transform* (DFT) matrices $\mathbf{F}_N \in \mathbb{C}^{N \times N}$ and $\mathbf{F}_M \in \mathbb{C}^{M \times M}$, with elements

$$[\mathbf{F}_N]_{n,n'} = \frac{1}{\sqrt{N}} e^{j2\pi(n-1)(\frac{n'-1}{N}-\frac{1}{2})}, n, n' \in [N], \quad (6a)$$

$$[\mathbf{F}_M]_{m,m'} = \frac{1}{\sqrt{M}} e^{j2\pi(m-1)(\frac{m'-1}{M}-\frac{1}{2})}, m, m' \in [M]. \quad (6b)$$

Consequently, based on a subarray basis indexed by i' , the beam-domain representation of the channel response (1) is given by [7, 36]

$$\check{H}_{s,k}^{i'}(t, \tau) = \mathbf{F}_N^H \mathbf{H}_{s,k}(t, \tau) \cdot \left(\mathbf{F}_M \odot \mathbf{1}_{\{(i'-1)\hat{M}+1:i'\hat{M}, 1:M\}} \right) = \sum_{l=1}^{L_k} \check{H}_{s,k,l}^{i'}(t) \delta(\tau - \tau_l), \quad (7)$$

where ($i' \equiv 1, \hat{M} = M$) for the FC architecture, and ($i' \in [M_{\text{RF}}], \hat{M} = \frac{M}{M_{\text{RF}}}$) for the OSPA architecture. Here we define $\check{H}_{s,k,l}^{i'}(t) := \mathbf{F}_N^H \mathbf{H}_{s,k,l}(t) \cdot \left(\mathbf{F}_M \odot \mathbf{1}_{\{(i'-1)\hat{M}+1:i'\hat{M}, 1:M\}} \right)$ as the beam-domain l -th multipath component between the k -th UE and the BS, where $\mathbf{1}_{\{a_1:a_2, b_1:b_2\}} \in \mathbb{C}^{M \times M}$ is an indicator matrix, with 1 at the components indexed by rows from a_1 to a_2 and by columns from b_1 to b_2 , otherwise zero. The indicator matrix takes into account the fact that the number of antenna elements for each subarray in the OSPA architecture is M_{RF} times less than that in the FC architecture.

As shown in our earlier work [16] (and the references therein), for the FC architecture, as the number of antennas M at the BS and N at the UE increases, the DFT basis provides a good sparsification of the propagation channel. As a result, $\check{H}_{s,k}^{i'}(t, \tau)$ can be approximated as a sparse matrix, with non-zero elements in the locations corresponding to small clusters of discrete AoA-AoD pairs. For the OSPA architecture, note that the indices of non-zero elements in $\check{H}_{s,k}^{i'}(t, \tau)$ are identical for all $i' \in [M_{\text{RF}}]$. However, the channel sparsity depends on the number of antennas in each subarray. In both cases, we may encounter a grid error in (7) since the AoAs-AoDs do not necessarily fall into the uniform grid $\Phi \times \Theta$. Nevertheless, as shown in [16], the grid error becomes negligible by increasing the number of (subarray) antennas (i.e., the grid resolution). In our simulations, we do not constrain the AoA-AoD pairs of the physical channel to take on values on the discrete grid; therefore, the grid discretization effect is fully taken into account in our numerical results.

B. Signaling Model

Because of space limitation, in this paper we focus on SC signaling. Similar conclusions can be reached for OFDM, although the latter is generally more fragile to frame synchronization

errors, large PAPR, and, before BA is achieved, to inter-carrier interference due to the fact that the Doppler spread between the several multipath components may be large [19, 37]. Let $\mathbf{x}_s(t) = [x_{s,1}(t), x_{s,2}(t), \dots, x_{s,K}(t)]^\top$ denote the continuous-time baseband equivalent signal (either pilot or data signal), transmitted over the s -th slot. With HDA beamforming, the beamformed signal at the output of the transmitter over the s -th slot is generally given by

$$\hat{\mathbf{x}}_s(t) = \sqrt{E_0} \cdot \mathbf{U}_s^{\text{RF}} \cdot \mathbf{W}_s^{\text{BB}} \cdot \mathbf{x}_s(t), \quad (8)$$

where for simplicity of exposition we restrict to the case of uniform power allocation, with $E_0 = \frac{P_{\text{tot}} T_c}{K}$ indicating the per-chip energy of each signal stream, where P_{tot} denotes the total radiated power at the BS and $T_c = \frac{1}{B}$ denotes the chip duration with B indicating the signaling bandwidth. In (8), we define $\mathbf{W}_s^{\text{BB}} \in \mathbb{C}^{M_{\text{RF}} \times K}$ and $\mathbf{U}_s^{\text{RF}} \in \mathbb{C}^{M \times M_{\text{RF}}}$ as the baseband (digital) and the RF analog beamforming matrices, respectively. Note that, depending on the transmitter architecture, the analog beamforming matrix \mathbf{U}_s^{RF} takes on the form

$$[\tilde{\mathbf{u}}_{s,1}, \tilde{\mathbf{u}}_{s,2}, \dots, \tilde{\mathbf{u}}_{s,M_{\text{RF}}}] \quad \text{and} \quad \begin{bmatrix} \tilde{\mathbf{u}}_{s,1} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \tilde{\mathbf{u}}_{s,2} & \cdots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \tilde{\mathbf{u}}_{s,M_{\text{RF}}} \end{bmatrix} \quad (9)$$

for the FC (left) and the OSPA (right) architectures, respectively, where $\tilde{\mathbf{u}}_{s,i} \in \mathbb{C}^{\hat{M}}$, $i \in [M_{\text{RF}}]$, with $\hat{M} = M$ for the FC architecture and $\hat{M} = \frac{M}{M_{\text{RF}}}$ for the OSPA architecture. Hence, in both cases \mathbf{U}_s^{RF} has dimension $M \times M_{\text{RF}}$, but FC has a full matrix, while OSPA has a block-diagonal matrix, due to the constrained connectivity. Without loss of generality, the beamforming vectors are normalized as $\sum_{i=1}^{M_{\text{RF}}} \|\mathbf{u}_{s,i}\|^2 = M_{\text{RF}}$.

The beamformed signal (8) goes through the channel as defined in (1). At the UE side, because of the HDA architecture, the UE does not have direct access to each antenna element. Instead, at each slot s , the UE obtains only a projection of the received signal by applying some beamforming vector in the analog domain. We consider a single RF chain at each UE as mentioned before. Thus, the received signal at the k -th UE side is given by

$$\begin{aligned} \hat{y}_{s,k}(t) &= \mathbf{v}_{s,k}^H \mathbf{H}_{s,k}(t, \tau) \circledast \hat{\mathbf{x}}_s(t) + z_{s,k}(t) \\ &= \sqrt{E_0} \mathbf{v}_{s,k}^H \mathbf{H}_{s,k}(t, \tau) \circledast (\mathbf{U}_s^{\text{RF}} \cdot \mathbf{W}_s^{\text{BB}} \cdot \mathbf{x}_s(t)) + z_{s,k}(t), \end{aligned} \quad (10)$$

where $\mathbf{v}_{s,k} \in \mathbb{C}^N$ denotes the normalized beamforming vector with $\|\mathbf{v}_{s,k}\| = 1$ at the k -th UE, and $z_{s,k}(t)$ is the continuous-time complex *Additive White Gaussian Noise* (AWGN) at the output of the UE RF chain, with a *Power Spectral Density* (PSD) of N_0 Watt/Hz.

In the following, we will evaluate the performance of different transmitter architectures as shown in Fig. 1. For this purpose, it is useful to first define the channel SNR before beamforming (BBF) $\text{SNR}_{\text{BBF},k}$, given by

$$\text{SNR}_{\text{BBF},k} = \frac{P_{\text{tot}} \sum_{l=1}^{L_k} \gamma_{k,l}}{N_0 B}. \quad (11)$$

where k is the index of the UE and $\gamma_{k,l}$ denotes the strength of the l -th multipath component. The SNR in (11) indicates the ratio of the total received signal power (summing over all the multipath components) over the total noise power at the receiver baseband processor input, assuming that the signal is isotropically transmitted by the BS and isotropically received at the k -th UE over the total bandwidth B . As mentioned before, one of the challenges of mmWaves communication is that the SNR before beamforming SNR_{BBF} in (11) may be very low.

III. BEAM ACQUISITION AND DATA TRANSMISSION

We evaluate the performance of the FC and OSPS architectures including both the BA phase and the consequent data transmission phase, where the latter uses the beam information obtained by the former. For the BA phase we use the scheme proposed in our previous work [19], that compares favorably with respect to several competing schemes proposed in the literature. For the sake of space limitation, we provide here only a high-level summary of the scheme and invite the reader to consider [19] for the full details. Fig. 3 (a) illustrates the considered frame structure, which consists of three parts: the beacon slot, the random access control channel (RACCH) slot, and the data slot. As shown in Fig. 3 (b), the BS broadcasts its pilot signals periodically over the beacon slots. The measurements are collected at each UE locally and independently of other UEs. Based on measurements accumulated over a sequence of several beacon slots, each UE can estimate a set of strongly coupled AoA-AoD pairs, corresponding to the directions of strong propagation paths between the UE and the BS arrays. These determine the beamforming direction for possible data transmission. During the RACCH slot, the BS stays in listening mode and the UEs send beamformed uplink packets. These packets contain basic information such as the UE ID and the beam indices of the selected BS beam directions. The BS responds with

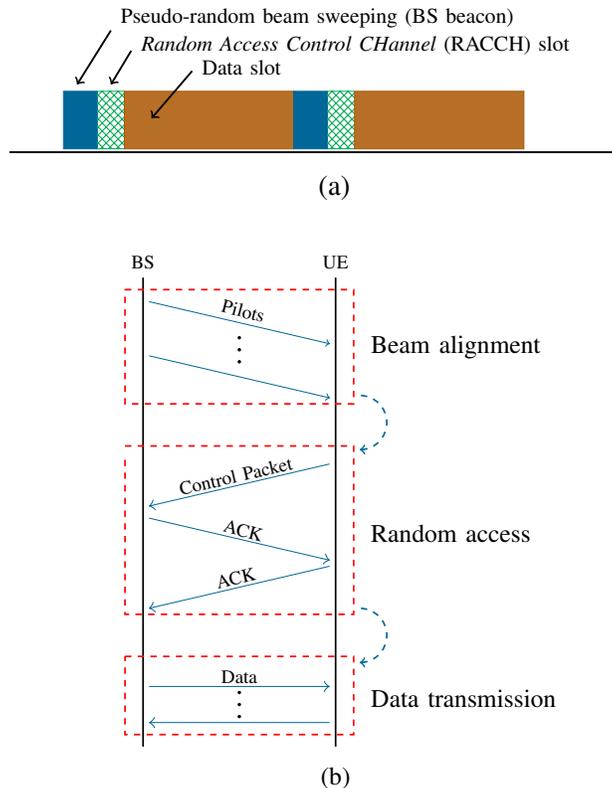


Fig. 3: Illustration of (a) the frame structure in the underlying system and (b) the communication process between the BS and a generic UE. The initial beam alignment phase is periodically done over beacon slots, followed by a random access stage to build up the connection between the BS and the active UEs, and consecutively the data communication.

an acknowledgment (ACK) data packet in the data subslot of a next frame, using the indicated beam indices for transmission. From this moment on, the BS and the UE are connected in the sense that, if the procedure is successful, they can communicate by aligning their beams along a small number of multipath components with strong average power transmission.

As explained in details in [19], the BS beacon signals are formed by M_{RF} different PN sequences, each of which undergoes a “multifinger” beam pattern obtained by selecting a subset of the columns (or masked DFT columns as in the case of OSPS). The beamforming patterns send the signal energy uniformly distributed along subsets of the BS AoD grid. The beamforming patterns follow a pre-determined pseudo-random sequence, similar in the spirit to the primary synchronization code of a W-CDMA 3G system for BS identification. During the beacon slot, each UE k receives using its own pseudo-random sequence of multifinger beam patterns, and integrates the received signal energy over the multiple time segments within a beacon slot in order

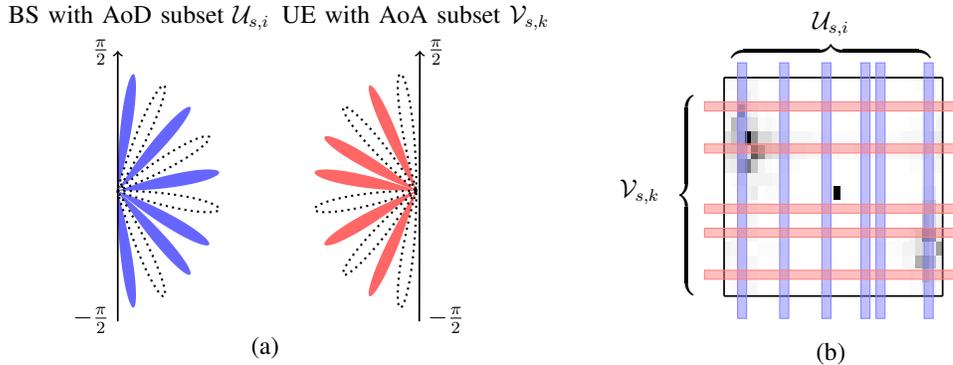


Fig. 4: (a) Illustration of the subset of AoA-AoDs at time slot s probed by the i -th beacon stream transmitted by the BS and received by the k -th UE, for $\hat{M} = N = 10$. The AoD subset is given by $\mathcal{U}_{s,i} = \{1, 3, 4, 6, 8, 10\}$ (numbered counterclockwise) with beamforming vector $\check{\mathbf{u}}_{s,i} = \frac{1}{\sqrt{6}}[1, 0, 1, 0, 1, 0, 1, 1, 0, 1]^T$. The AoA subset is given by $\mathcal{V}_{s,k} = \{2, 4, 5, 7, 9\}$ (numbered counterclockwise) with receive beamforming vector $\check{\mathbf{v}}_{s,k} = \frac{1}{\sqrt{5}}[0, 1, 0, 1, 1, 0, 1, 0, 1, 0]^T$. (b) The beam-domain channel gain matrix (with one LOS component and two scattered multipath components indicated by the dark spots, generated by the QuaDRiGa simulator) measured along $\mathcal{V}_{s,k} \times \mathcal{U}_{s,i}$.

to obtain an estimate of the average received energy. As a result, this fully non-coherent energy measurement yields (approximately) the average energy sum of several multipath components. These multipath components corresponds to the AoA-AoD pairs in the grid for which the BS transmit directions and the UE receive directions meet. Fig. 4 (a) shows an example of transmit and receive multifinger beam patterns and Fig. 4 (b) shows the corresponding masks of crossing AoA-AoD directions, superimposed with the second moments (channel gain) of the beam-domain channel matrix generated by the QuaDRiGa simulator. The goal of the BA algorithm run at the UE side is to identify the position of the strong components, i.e., the small dark spots in the plot of Fig. 4(b). It turns out that this problem can be cast as the reconstruction of a sparse non-negative vector from noisy linear measurements, which can be efficiently obtained by solving a non-negative least-squares (NNLS) problem. It can be shown that NNLS naturally induce sparsity in the solution, and it is very efficient to solve by a plethora of well-known algorithms (e.g., projected gradient). The full details of the BA scheme, as well as extensive comparison with other competing state-of-the-art schemes, are provided in [19].

We denote by $\Gamma_k \in \mathbb{C}^{N \times M}$ as the matrix of second moments of the beam-domain channel coefficients between the BS array and the k -th UE array. An example of Γ_k is illustrated in Fig. 5(a). Also, Fig. 5(b) shows the estimate Γ_k^* of Γ_k provided by the NNLS estimation at UE k . Once the BA algorithm yields Γ_k^* , the k -th UE will send a beamformed control packet

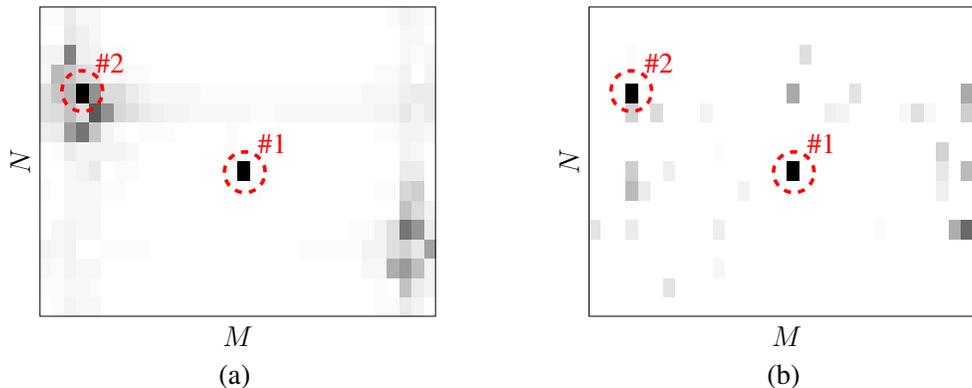


Fig. 5: Illustration of the second moments of the beam-domain channel matrix Γ_k : (a) the actual QuaDRiGa generated Γ_k , (b) the NNLS estimated Γ_k^* . The dashed circles indicate the top $p = 2$ strongest components in Γ_k and Γ_k^* , respectively. We announce a success in the BA phase if the locations of the strongest component in Γ_k and in Γ_k^* are consistent.

to the BS in the RACCH. The UE chooses the beamforming direction corresponding to the strongest AoA direction obtained from Γ_k^* , meanwhile the BS stays in listening mode during the RACCH, using a sectored beamforming configuration. In this way, full beamforming gain at the UE transmit side and a limited sector beamforming gain at the BS receive side can be achieved. Once the RACCH packet is received, the BS can use the transmit beam indicated by UE k to communicate data. In the next section, we focus on the data communication phase assuming that the RACCH has been correctly received, therefore, both the BS and the UE know the indices of the strong components in Γ_k^* . Notice that if the NNLS estimation fails, it is likely that the RACCH will not be received or will be received in error, because the beamforming gain at the UE side will be poor. In this case, the UE will not receive a data packet and after a given time-out will try the BA procedure again. Also in the (very unlikely) case of a collision in the RACCH, the same time-out procedure can be exploited. Therefore, data communication effectively takes place only when a) the strong multipath components in Γ_k are correctly estimated and b) when the RACCH decoding is successful. In [19] we have already argued that the probability that the BA procedure fails is dominated by the error probability in the estimation of the strong components of Γ_k . Hence, a sensible system design approach consists of allowing a sufficient number of beacon slots such that the probability of success in identifying the strong components of Γ_k is close to 1, and designing the HDA beamforming scheme in the assumption that the estimation of Γ_k is correct. As a result, we shall compare the FC and OSPA architectures in terms of number of beacon slots needed to achieve a BA success probability near 1, and their achieved spectral

efficiency under such condition. In any case, the designed HDA precoders in our simulations are always obtained from the true NNLS estimation $\mathbf{\Gamma}_k^*$, and not by the genie-aided exact knowledge of $\mathbf{\Gamma}_k$.

A. Data Communication Phase

We assume that the BS simultaneously schedules $K = M_{\text{RF}}$ UEs. With the small cell configuration as illustrated in Fig. 2, the distance differences between each UE and the BS are very small, implying that the received power w.r.t. the LOS path for each UE within the BS coverage are similar. Although schedulers such as random or proportionate fair scheduler are commonly used in sub 6 GHz, the directionality of the mmWave channel instead calls for schedulers that select groups of users with good angular separation (directional scheduler) [26]. More precisely, we assume that the selected K UEs have similar received power in terms of the strongest path, and their strongest AoDs in the downlink are at least $\Delta\theta_{\min}$ away from each other.

Let $\mathbf{x}^d(t) = [x_1^d(t), x_2^d(t), \dots, x_K^d(t)]^T$ denote the complex baseband data signal,⁵ with $x_k^d(t)$, $k \in [K]$, corresponding to the k -th UE, given by

$$x_k^d(t) = \sum_{n=1}^{N_d} d_{k,n} p_r(t - nT_c), \quad (12)$$

where $p_r(t)$ is the unit-energy square-root Nyquist pulse shaping filter, $\{d_{k,n}\}$ denote the unit-energy modulation symbols belonging to a suitable modulation constellation [33], and N_d indicates the number of the transmit symbols. Accordingly, the received data signal at the k -th UE is given by (refer to (10))

$$\begin{aligned} \hat{y}_k(t) &= \sqrt{E_0} \mathbf{v}_k^H \mathbf{H}_k(t, \tau) \otimes (\mathbf{U}^{\text{RF}} \cdot \mathbf{W}^{\text{BB}} \cdot \mathbf{x}^d(t)) + z_k(t) \\ &= \sum_{n=1}^{N_d} \sum_{k'=1}^K \sum_{l=1}^{L_k} d_{k',n} \sqrt{E_0} \mathbf{v}_k^H \mathbf{H}_{k,l}(t) \mathbf{U}^{\text{RF}} \mathbf{w}_{k'} p_r(t - \tau_{k',l} - nT_c) + z_k(t) \\ &= \sum_{n=1}^{N_d} \sum_{k'=1}^K \sum_{l=1}^{L_k} C_{k,k',l,n} e^{j\Delta_{k,n,l}} p_r(t - \tau_{k',l} - nT_c) + z_k(t) \end{aligned} \quad (13)$$

⁵From now on, we ignore the slot index s for notation simplicity, also because once a successful BA is achieved, the channel statistical property, the precoding vector at the BS, and combining vector at each UE are invariant within many slots. However, note that this invariance holds only until a new updated BA takes place, implying that the underlying channel may encounter large mobility, blockage, etc.

where $\mathbf{w}_{k'}$ denotes the k' -th column of \mathbf{W}^{BB} , $\Delta_{k,n,l} = 2\pi(\check{\nu}_{k,l} + \nu_{k,l}nT_c)$, and $C_{k,k',l,n} := \rho_{k,l}d_{k',n}\sqrt{E_0}(\mathbf{v}_k^{\text{H}}\mathbf{a}_{\text{R}}(\phi_{k,l})\mathbf{a}_{\text{T}}(\theta_{k,l})^{\text{H}}\mathbf{U}^{\text{RF}}\mathbf{w}_{k'})$. We assume that each UE uses standard timing synchronization with respect to its strongest multipath component indexed by l^1 , which is selected by its initial BA. To decode the data signal, each UE performs matched filtering with respect to the symbol pulse $p_r(t)$, sampling at epochs $t = \hat{n}T_c + \tau_{k,l^1}$. It follows that the discrete-time baseband signal received at the k -th UE receiver takes on the form

$$\begin{aligned} y_k[\hat{n}] &= y_k(t)|_{t=\hat{n}T_c+\tau_{k,l^1}} = \hat{y}_k(t) \otimes p_r^*(-t)|_{t=\hat{n}T_c+\tau_{k,l^1}} \\ &= \sum_{n=1}^{N_d} \sum_{k'=1}^K \sum_{l=1}^{L_k} C_{k,k',l,n} e^{j\Delta_{k,n,l}} \varphi_r[\hat{n}_{k,k',\hat{n},n,l}^{\Delta}] + \sum_{n=1}^{N_d} z_k^c[\hat{n}] \\ &= \sum_{n=1}^{N_d} \left(\sum_{l=1}^{L_k} C_{k,k,l,n} e^{j\Delta_{k,n,l}} \varphi_r[\hat{n}_{k,k,\hat{n},n,l}^{\Delta}] + \sum_{\substack{k' \neq k \\ l=1}}^{L_k} C_{k,k',l,n} e^{j\Delta_{k,n,l}} \varphi_r[\hat{n}_{k,k',\hat{n},n,l}^{\Delta}] + z_k^c[\hat{n}] \right), \quad (14) \end{aligned}$$

where $\hat{n}_{k,k',\hat{n},n,l}^{\Delta} := (\hat{n} - n)T_c + \tau_{k,l^1} - \tau_{k',l}$, $\varphi_r[t^{\Delta}] = \varphi_r(t)|_{t=t^{\Delta}} := \int p_r(\tau)p_r^*(\tau - t^{\Delta})d\tau$, and $z_k^c[\hat{n}]$ denotes the noise at the output of the matched filter with variance $N_0 \cdot \int |p_r(t)|^2 dt = N_0$. As we can see, the first term in (14) corresponds to the desired data symbol $d_{k,n}$ multiplied by a different complex coefficient over each path l .⁶ Whereas, the last two terms in (14) correspond to the multiuser interference and noise, respectively. By treating the multiuser interference as noise, the asymptotic ergodic spectral efficiency of the k -th UE is given by

$$R_k = \log_2 \left(1 + \frac{\mathbb{E} \left[\left| \sum_{l=1}^{L_k} C_{k,k,l,n} e^{j\Delta_{k,n,l}} \varphi_r[\hat{n}_{k,k,\hat{n},n,l}^{\Delta}] \right|^2 \right]}{\mathbb{E} \left[\left| \sum_{k' \neq k} \sum_{l=1}^{L_k} C_{k,k',l,n} e^{j\Delta_{k,n,l}} \varphi_r[\hat{n}_{k,k',\hat{n},n,l}^{\Delta}] \right|^2 \right] + N_0} \right), \quad (15)$$

and the sum rate reads $R_{\text{sum}} = \sum_{k=1}^K R_k$. In all schemes treated here, coherent communication can be practically achieved by including per-user beamformed pilot symbols at the cost of a very small overhead, as it is quite state of art and usual in virtually any modern wireless communication standard. For simplicity, we shall not take into account this overhead or the degradation of quasi-coherent receivers, which is well known and not a specific feature of the systems under consideration.

⁶Actually, we have shown in our precious work [19] that, the phase perturbations over several strong paths are easy to compensate by standard carrier synchronization techniques given that a successful BA is achieved and the effective channel after BA has a very small time spreading. Due to the space limit, in (14) and also in our simulations, we will keep the phase perturbations such that the numerical results coincide with the conservative worst-case scenario.

1) Hybrid Precoding Formulation

Now the remaining problem is how to define the precoding/combining vectors. We assume that the BS communicates with the k -th UE along its top- p beams. We will show later that the parameter $p \geq 1$ is somehow a tradeoff between the transmitter power spreading, multiuser interference, and the system robustness to potential blockages. To simplify the practical implementation, we define the combining vector at the k -th UE as

$$\mathbf{v}_k = \frac{1}{\sqrt{p}} \mathbf{F}_N \cdot \sum_{p'=1}^p \check{\mathbf{v}}_{k,p'}, \quad (16)$$

where $\check{\mathbf{v}}_{k,p'} \in \mathbb{C}^N$ is an all-zero vector with a 1 at the component corresponding to the p' -th strong AoA, i.e., the AoA index of the p' -th strong component in Γ_k^* . Denoted by $\mathbf{V} \in \mathbb{C}^{NK \times K}$ as the aggregated receive beamforming matrix given by $\mathbf{V} = \text{diag}(\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_K)$. It follows that the receive data signal vector $\bar{\mathbf{y}}(t) = [y_1(t), y_2(t), \dots, y_K(t)]^T \in \mathbb{C}^K$ corresponding to the K UEs can be written as

$$\begin{aligned} \bar{\mathbf{y}}(t) &= \sqrt{E_0} \mathbf{V}^H \cdot \bar{\mathbf{H}}(t, \tau) \otimes (\mathbf{U}^{\text{RF}} \cdot \mathbf{W}^{\text{BB}} \cdot \mathbf{x}^d(t)) + \bar{\mathbf{z}}(t) \\ &\stackrel{\text{(a)}}{=} \sqrt{E_0} (\mathbf{V}^H \cdot \bar{\mathbf{H}}(t, \tau) \cdot \bar{\mathbf{U}} \cdot \mathbf{A}^{\text{RF}} \cdot \mathbf{W}^{\text{BB}}) \otimes \mathbf{x}^d(t) + \bar{\mathbf{z}}(t) \\ &\stackrel{\text{(b)}}{=} \sqrt{E_0} (\tilde{\mathbf{H}}(t, \tau) \cdot \mathbf{A}^{\text{RF}} \cdot \mathbf{W}^{\text{BB}}) \otimes \mathbf{x}^d(t) + \bar{\mathbf{z}}(t) \end{aligned} \quad (17)$$

where $\bar{\mathbf{z}}(t) \in \mathbb{C}^K$ indicates the noise vector, $\mathbf{U}^{\text{RF}} := \bar{\mathbf{U}} \cdot \mathbf{A}^{\text{RF}}$ is the analog beamforming matrix, $\tilde{\mathbf{H}}(t, \tau) := \mathbf{V}^H \cdot \bar{\mathbf{H}}(t, \tau) \cdot \bar{\mathbf{U}}$ denotes a constructed effective channel, and $\bar{\mathbf{H}}_s(t, \tau) \in \mathbb{C}^{NK \times M}$ represents the aggregated instantaneous channel of all the K UEs given by

$$\bar{\mathbf{H}}(t, \tau) = [\mathbf{H}_1(t, \tau)^T, \mathbf{H}_2(t, \tau)^T, \dots, \mathbf{H}_K(t, \tau)^T]^T, \quad (18)$$

where $\mathbf{H}_k(t, \tau)$, $k \in [K]$, is given in (1). In (17)(a), we define $\bar{\mathbf{U}} \in \mathbb{C}^{M \times pK}$ as the angular support, and $\mathbf{A}^{\text{RF}} = [\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_K] \in \mathbb{C}^{pK \times K}$ as the coefficient tuning for the analog part. More precisely, we assume $\bar{\mathbf{U}} = [\mathbf{U}_1, \dots, \mathbf{U}_K]$, where $\mathbf{U}_k \in \mathbb{C}^{M \times p}$, $k \in [K]$, takes on the form

$$\mathbf{U}_k = \left(\mathbf{F}_M \odot \mathbf{1}_{\{(k'-1)\hat{M}+1:k'\hat{M}, 1:M\}} \right) \cdot [\check{\mathbf{u}}_{k,1}, \check{\mathbf{u}}_{k,2}, \dots, \check{\mathbf{u}}_{k,p}], \quad (19)$$

where ($i' \equiv 1, \hat{M} = M$) for the FC architecture, and ($i' = k, \hat{M} = \frac{M}{M_{\text{RF}}}$) for the OSPS architecture. Also, we define $\check{\mathbf{u}}_{k,p'} \in \mathbb{C}^M$, $p' \in [p]$, as an all-zero vector with a 1 at the component corresponding to the p' -th strongest AoD of Γ_k^* .

Notice that in order to construct the beamforming vector at each k -th UE and the precoding vectors at the BS, only the AoA-AoD indices of the p strongest components in the estimated

channel gain matrix Γ_k^* are needed. Then, once these vectors are fixed, the resulting effective channel has much lower dimensions than the original physical $N \times M$ channel (from array to array). Therefore, it can be estimated using orthogonal uplink pilots and channel reciprocity as in regular TDD MU-MIMO (e.g., see [38, 39]). Namely, the constructed effective channel matrix $\tilde{\mathbf{H}}(t, \tau)$ in (17) (b) has dimension $K \times (pK)$, and can be estimated using pK uplink pilot sub-slots using TDD reciprocity.

2) Beam Steering (BST) Scheme

The BST scheme consists of simply steering the K data streams towards the K UEs along their strongest AoD. Hence, we have $p = 1$ in (16) and in (19), respectively. In such case, the analog tuning matrix and the baseband precoding matrices under the BST precoding scheme turn to be identity, i.e., $\mathbf{A}^{\text{RF}} = \mathbf{W}^{\text{BB}} = \mathbf{I}_K$. Note that in the BST scheme, we do not need any additional uplink channel estimation of $\tilde{\mathbf{H}}(t, \tau)$. Namely, once the UEs has fed back its strongest AoD control packet, the BS can immediately provide the BST precoder.

3) Analog Maximum Ratio Transmission (MRT) Scheme

In this scheme, we aims to maximize the desired signal power as well as to increase the scheme blockage robustness. To this end, the baseband precoding matrix remains identity, i.e., $\mathbf{W}^{\text{BB}} = \mathbf{I}_K$, while the k -th analog MRT tuning vector (i.e., the k -th column of \mathbf{A}^{RF}) is given by

$$\mathbf{a}_k = \left(\tilde{\mathbf{H}}(t, \tau)_{\{k,:\}} \right)^{\text{H}} \odot \mathbf{1}_{\{(k'-1)\hat{p}+1:k'\hat{p}\}} \cdot \Delta^{\text{RF}}, \quad (20)$$

where $\tilde{\mathbf{H}}(t, \tau)_{\{k,:\}}$ indicates the k -th row of $\tilde{\mathbf{H}}(t, \tau)$, and $\Delta^{\text{RF}} \in \mathbb{R}_+$ denotes the normalizing factor such that $\sum_{i=1}^{M_{\text{RF}}} \|\mathbf{u}_i\|^2 = M_{\text{RF}}$. The indicator vector $\mathbf{1}_{\{(k'-1)\hat{p}+1:k'\hat{p}\}} \in \mathbb{C}^{pK}$ has components 1 over the index $\{(k' - 1)\hat{p} + 1 : k'\hat{p}\}$ otherwise 0, where $(k' \equiv 1, \hat{p} = pK)$ for the FC architecture and $(k' = k, \hat{p} = p)$ for the OSPA architecture. Here the indicator vector ensures that, in the OSPA architecture, the analog beamforming matrix $\mathbf{U}^{\text{RF}} = \bar{\mathbf{U}} \cdot \mathbf{A}^{\text{RF}}$ satisfies the block diagonal structure as illustrated in (9).

4) Joint Analog Maximum Ratio and Baseband Zeroforcing (MR-ZF) Scheme

On top of the previous MRT scheme, in this joint MR-ZF scheme, we propose to make use of the baseband precoding to further reduce the multiuser interference. Accordingly, the analog MRT vectors in \mathbf{A}^{RF} are given by (20), while the baseband ZF matrix \mathbf{W}^{BB} takes on the form

$$\mathbf{W}^{\text{BB}} = \left(\tilde{\mathbf{H}}(t, \tau) \mathbf{A}^{\text{RF}} \right)^{\text{H}} \cdot \left(\tilde{\mathbf{H}}(t, \tau) \mathbf{A}^{\text{RF}} \left(\tilde{\mathbf{H}}(t, \tau) \mathbf{A}^{\text{RF}} \right)^{\text{H}} \right)^{-1} \cdot \Delta^{\text{ZF}}, \quad (21)$$

where $\Delta^{\text{ZF}} \in \mathbb{R}_+$ is the normalizing factor ensuring the total radiated power constraint, i.e., $\sum_{k=1}^K \|\mathbf{w}_k\|^2 = K$.

IV. HARDWARE IMPAIRMENTS

In all the above derivations, we have implicitly assumed that all the hardware components work in their ideal range without any distortion or power dissipation. However, in practical hardware systems, such assumption is not trivial to meet. For example, the implementation of HDA transceivers consists of a large number of power dividers and combiners in the analog part, particularly for the FC architecture. The power dissipation caused by these components has a severe impact on the transmit power and the power efficiency. Moreover, due to the superposition of multiple beamformed pilots/data, the input signal at the PAs may encounter a large PAPR. Also, different beamforming vectors will create different power levels for different PAs. As a result, the input power for some individual PAs may exceed their saturation limit (relevant to per-antenna power constraint) and even cause a disruption of the whole transmission. All these hardware impairment have a severe impact on the transmitter performance and should not be neglected. In this section, we will provide the mathematical model to evaluate the hardware efficiency of different transmitter architectures given in Fig. 1.

We assume that each analog path has simultaneous amplitude and phase control as shown in Fig. 1. Refer to (8), let $\tilde{\mathbf{x}} \in \mathbb{C}^M$ denote the pre-amplified beamformed signal⁷, given by

$$\tilde{\mathbf{x}} = \sqrt{\alpha_{\text{com}}} \cdot \tilde{\mathbf{U}}^{\text{RF}} \cdot \sqrt{\alpha_{\text{div}}} \cdot \mathbf{W}^{\text{BB}} \cdot \mathbf{x}, \quad (22)$$

where $\mathbf{x} = [x_1, \dots, x_K] \in \mathbb{C}^K$ denotes the transmit symbol, with $\mathbb{E}[|x_i|^2] = \epsilon$, $i \in [K]$. The factor α_{div} indicates the power splitting at the divider, with $\alpha_{\text{div}} = \frac{1}{M}$ for the FC architecture as shown in Fig. 1 (a) and $\alpha_{\text{div}} = \frac{M_{\text{RF}}}{M}$ for the OSPS architecture as shown in Fig. 1 (b). Moreover, the factor α_{com} models the power dissipation factor of the combiners, i.e., $\alpha_{\text{com}} = \frac{1}{M_{\text{RF}}}$ for the FC architecture, and $\alpha_{\text{com}} = 1$ for the OSPS architecture. Both α_{div} and α_{com} result from the hardware implementation and are based on the corresponding S-parameters of the dividers and combiners as in [5]. We assume that the baseband beamforming matrix \mathbf{W}^{BB} is of dimension $K \times K$ with $K = M_{\text{RF}}$, and the analog beamforming matrix $\tilde{\mathbf{U}}^{\text{RF}} = [\mathbf{u}_1, \dots, \mathbf{u}_{M_{\text{RF}}}] \in \mathbb{C}^{M \times M_{\text{RF}}}$ satisfies the specific FC/OSPS architecture as illustrated in (9).

⁷For notation simplicity, here we ignored the slot index s and the time index t .

We consider the rather simple BST precoding with $\mathbf{W}^{\text{BB}} = \mathbf{I}_K$. To first meet the total power constraint, for any $i \in [M_{\text{RF}}]$, we have $\|\mathbf{u}_i\|^2 = M$ for the FC architecture and $\|\mathbf{u}_i\|^2 = \frac{M}{M_{\text{RF}}}$ for the OSPS architecture, respectively. It follows that the effective pre-amplified radiated power of the beamformed signal $\tilde{\mathbf{x}}$ in (22) can be written as

$$\tilde{P} = \mathbb{E}[\tilde{\mathbf{x}}^H \tilde{\mathbf{x}}] = \alpha_{\text{com}} \alpha_{\text{div}} \cdot \mathbb{E}[\mathbf{x}^H (\tilde{\mathbf{U}}^{\text{RF}})^H \tilde{\mathbf{U}}^{\text{RF}} \mathbf{x}] = \alpha_{\text{com}} \alpha_{\text{div}} \cdot \text{tr} \left(\mathbb{E}[\mathbf{x} \mathbf{x}^H] \cdot (\tilde{\mathbf{U}}^{\text{RF}})^H \tilde{\mathbf{U}}^{\text{RF}} \right). \quad (23)$$

Accordingly, the pre-amplified radiated power for the FC and the OSPS architectures reads $\tilde{P}_{\text{FC}} = \epsilon M_{\text{RF}} \frac{1}{M_{\text{RF}}}$ and $\tilde{P}_{\text{OSPS}} = \epsilon M_{\text{RF}}$, respectively. As we can see, in order to achieve the same output power, the FC transmitter should compensate for an additional combiner power dissipation. More precisely, the transmitter should either boost the input signal as $M_{\text{RF}} \mathbf{x}$ or choose PAs with larger gain for the amplification stage. We consider the former approach and mathematically include the potential boosting factor M_{RF} as well as the factors $(\alpha_{\text{com}}, \alpha_{\text{div}})$ into the beamforming matrix $\tilde{\mathbf{U}}^{\text{RF}}$. Denoted by \mathbf{U}^{RF} as the integrated analog beamforming matrix, such that the pre-amplified beamformed signal in (22) can be written as $\tilde{\mathbf{x}} = \mathbf{U}^{\text{RF}} \cdot \mathbf{W}^{\text{BB}} \cdot \mathbf{x}$, which is consistent with our assumptions and formulations in Section II.

The beamformed signal (22) then goes through the amplification stage, where at each antenna branch a PA amplifies the signal before transmission. We assume that the PAs in different antenna branches have the same input-output relation. For any given antenna in the transmitter array, let P_{rad} denote the radiated power of the antenna, and P_{cons} denote the consumed power of the corresponding PA, which includes both the radiated power and the dissipated power. Following the approach in [28], the power consumed by the PA takes on the form

$$P_{\text{cons}} = \frac{\sqrt{P_{\text{max}}}}{\eta_{\text{max}}} \sqrt{P_{\text{rad}}}, \quad (24)$$

where P_{max} is the maximum output power of the PA with $P_{\text{rad}} \leq P_{\text{max}}$ and η_{max} is the maximum efficiency of the PA. Note that this relation holds for the most common PA implementations and is therefore a good choice for the following calculation. Considering that the PAs are often the predominant power consumption part, we define η_{eff} given by

$$\eta_{\text{eff}} = \frac{P_{\text{rad}}}{P_{\text{cons}}} \quad (25)$$

as the metric to effectively compare the power efficiency of the two transmitter architectures shown in Fig. 1. Note that due to the superposition of multiple beamforming vectors (particularly for the FC architecture) and the potentially high PAPR of the time-domain transmit waveform $\tilde{\mathbf{x}}$ in (22) (particularly with OFDM signaling), the input power for some individual PA may

exceed its saturation limit. This would result in non-linear distortion and even the disruption of the whole transmission. To compare the two transmitter architectures and ensure that all the underlying M PAs simultaneously work in their linear range, we generally have two options:

Option I: Both the FC architecture and the OSPS architecture utilize the same PA but apply a different input back-off $\alpha_{\text{off}} \in (0, 1]$, such that the peak power of the radiated signal is smaller than P_{max} . As a reference, we denote by $(P_{\text{rad},0}, \eta_{\text{max},0})$ as the parameters of a reference PA under the reference precoding/beamforming strategy with a power backoff factor $\alpha_{\text{off},0}$ (as illustrated later in Section V). For different scenarios (with certain α_{off}) the average radiated power and the consumed power take the form $P_{\text{rad}} = \frac{\alpha_{\text{off}}}{\alpha_{\text{off},0}} P_{\text{rad},0}$, $P_{\text{cons}} = \frac{\sqrt{P_{\text{max},0}}}{\eta_{\text{max},0}} \sqrt{P_{\text{rad}}}$. The transmitter power efficiency is given by

$$\eta_{\text{eff}} = \frac{P_{\text{rad}}}{P_{\text{cons}}} = \frac{\sqrt{P_{\text{rad}}} \cdot \eta_{\text{max},0}}{\sqrt{P_{\text{max},0}}}. \quad (26)$$

Option II: We choose to deploy different PAs for the FC architecture and the OSPS architecture. More precisely, we assume that the underlying PA has a maximum output power of $P_{\text{max}} = \frac{\alpha_{\text{off},0}}{\alpha_{\text{off}}} P_{\text{max},0}$, where α_{off} has the same value as in *Option I*. Consequently, the average radiated power and the consumed power of the underlying PA can be written as $P_{\text{rad}} = P_{\text{rad},0}$, $P_{\text{cons}} = \frac{\sqrt{P_{\text{max},0} \cdot \alpha_{\text{off},0} / \alpha_{\text{off}}}}{\eta_{\text{max}}} \sqrt{P_{\text{rad}}}$. The transmitter power efficiency is given by

$$\eta_{\text{eff}} = \frac{P_{\text{rad}}}{P_{\text{cons}}} = \frac{\sqrt{P_{\text{rad}}} \cdot \eta_{\text{max}}}{\sqrt{P_{\text{max},0} \cdot \alpha_{\text{off},0}}} \cdot \sqrt{\alpha_{\text{off}}}. \quad (27)$$

Note that the characteristics (P_{max} and η_{max}) of different PAs highly depend on the operation frequency, implementation, and technology. Aiming at illustrating how to apply the proposed analysis framework in practical system design, we will exemplify a set of PA parameters in Section V to evaluate the efficiency η_{eff} of the two architectures given in Fig. 1. For the comparison of BA and data communication algorithms, we are interested in the performance of the corresponding algorithms using the different transmitter architectures but with the same channel condition as in (11). Therefore, we assume the same total radiated power P_{tot} constraint for both architectures in Fig. 1. In practical systems, this assumption can be satisfied by applying a certain power backoff as in *Option I* or choosing different PAs as in *Option II*. This in addition fulfills the per-antenna power constraint, such that all the underlying PAs work in their linear range with an identical scalar gain. However, we will show in Section V that, under the same radiated power constraint, different architectures may have a different power efficiency.

V. NUMERICAL EVALUATION

We now present the numerical results to evaluate the proposed precoding schemes and to illustrate the performance of different transmitter architectures as shown in Fig. 1. The BA scheme was already extensively studied in [16, 18, 19] in terms of complexity, system-level scalability, and robustness to fast channel time-variations/large Doppler spread. Hence, here we focus only on the difference in time-to-successful BA required by the two BS architectures under comparison. We consider a system with a BS using $M = 128$ antennas and $M_{\text{RF}} = 4$ RF chains. The BS simultaneously schedules $K = M_{\text{RF}} = 4$ UEs, each of which uses $N = 16$ antennas and $N_{\text{RF}} = 1$ RF chain. We assume a short preamble structure used in IEEE 802.11ad [1, 40], where the beacon slot is of duration $t_0 S = 1.891 \mu\text{s}$. The system is assumed to work at $f_0 = 40$ GHz with a bandwidth of $B = 0.8$ GHz, namely, each beacon slot amounts to more than 1500 chips.

In the following simulations, otherwise stated, we will assume a fixed total radiated power constraint P_{tot} , where all the underlying PAs working in their linear range (w.r.t., per-antenna power constraint P_{max}) with an identical scalar gain. The MU-MIMO channel is generated in two ways:

1) In Section V-A and Section V-B, we use the channel model in (1) to generate the channel matrix between each UE k and the BS. Based on the practical mmWave MIMO channel measurements in [29], we assume $L_k = 3$, $k \in [K]$, multipath components for each UE, given by $(\gamma_{k,1} = 1, \eta_{k,1} = 100)$, $(\gamma_{k,2} = 0.6, \eta_{k,2} = 10)$ and $(\gamma_{k,3} = 0.4, \eta_{k,3} = 0)$ with respect to (4). Thus, the first link can be roughly regarded as the LOS path, while the remaining links represent the non-LOS (NLOS) paths. We also assume that the LOS paths for the simultaneously scheduled UEs are well separated in the beam domain, while all the NLOS paths are generated in a random way.

2) In Section V-C, we use the quasi-deterministic radio channel generator (QuaDRiGa) to generate the propagation channel matrix. The channel model is based on the 3GPP-3D urban micro-cell configuration [30]. In this case, the height of the BS antenna array is set to 10 m. The beam center of the BS orientates to the ground with an elevation angle⁸ of $\alpha_e = -20^\circ$ as shown in Fig. 2 (a). The simultaneous scheduled UEs are set to 1.5 m in height and $18 \sim 25$ m horizontally away from the BS with a downlink AoD difference of $\Delta\theta_{\text{min}} \approx 8^\circ$ [41]. Each UE

⁸In QuaDRiGa, the elevation angle 90° points to the zenith and 0° points to the horizon.

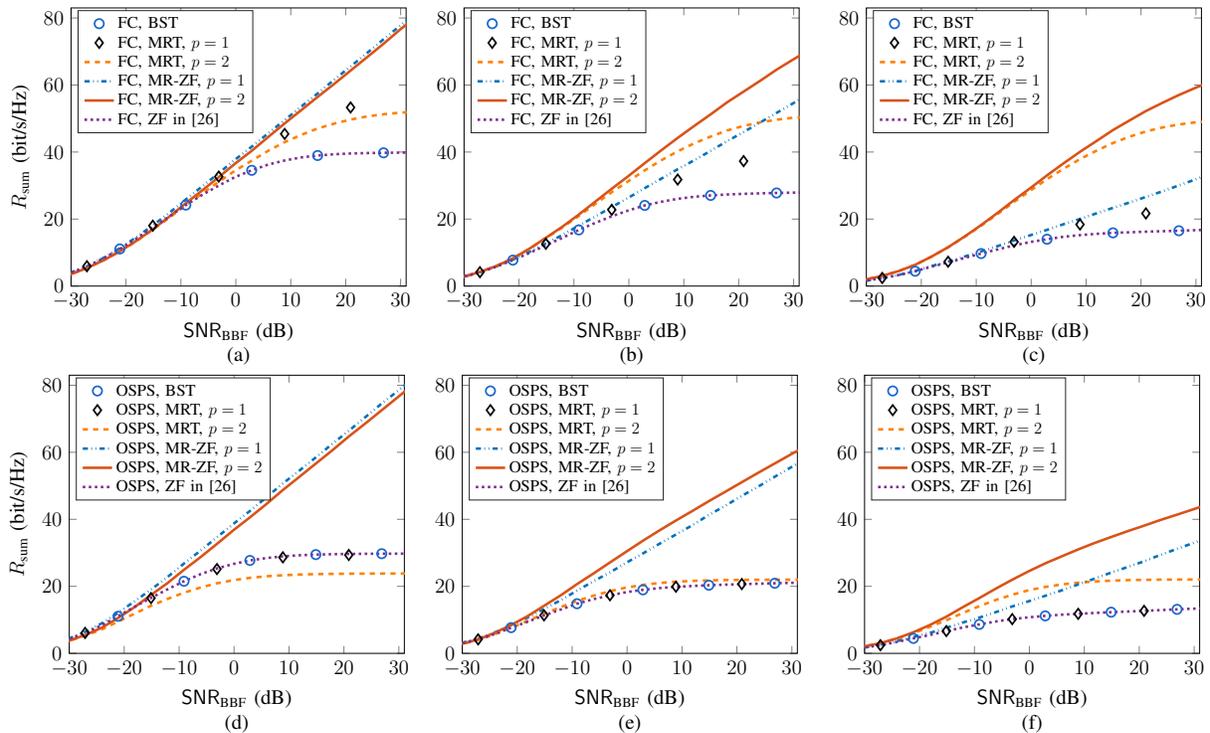


Fig. 6: The sum spectral efficiency vs. increasing SNR_{BBF} . The blockage probability of the strongest path is given by (a) 0.0, (b) 0.3, (c) 0.6 for the FC architecture, and (d) 0.0, (e) 0.3, (f) 0.6 for the OSPS architecture.

k moves towards the BS at a speed of $\Delta v_k = 1$ m/s. We will show that the numerical results based on our proposed channel model (1) are consistent with the results based on the QuaDRiGa generator, implying that the proposed work not only theoretically but also practically provides valuable references for mmWave system design.

A. Evaluation of the Proposed Precoding Schemes

The efficiency of the proposed precoding schemes are illustrated in Fig. 6. As a comparison, we also simulate the ZF precoder proposed in [26], where the effective channel is approximated by the initial BA vectors, and only a single path is selected between each UE and the BS. As we can see from Fig. 6(a), for the FC architecture with no blockage, all the schemes coincide with each other in the range of $\text{SNR}_{\text{BBF}} \leq 0$ dB. Whereas when $\text{SNR}_{\text{BBF}} > 0$ dB, the performance ranking of the underlying precoding schemes is as follow $(\text{MR-ZF}, p = 2) \approx (\text{MR-ZF}, p = 1) > (\text{MRT}, p = 1) > (\text{MRT}, p = 2) > (\text{BST}) \approx (\text{ZF in [26]})$. Here the MRT scheme with $p = 2$ performs worse than with $p = 1$ due to the fact of power spreading and the fact that with multiple receiving directions, the UE tends to have more interference. However, this effect is not observable in the MR-ZF scheme because of the further power coefficient tuning and interference cancellation,

which results from the baseband zeroforcing. Next, by increasing the blockage probability of the strongest path while remaining unblocked for all the less strong paths between each UE and the BS, as shown in Fig. 6 (b) and Fig. 6 (c), the curves with $p = 2$ drops much less than the others (equivalent to $p = 1$), and the scheme of MR-ZF with $p = 2$ achieves the best performance. For the OSPS architecture, when there is no blockage as shown in Fig. 6 (d), in the low SNR range ($\text{SNR}_{\text{BBF}} \leq -10$ dB), all the curves (roughly) coincide with each other. Whereas, by increasing $\text{SNR}_{\text{BBF}} > -10$ dB, the precoding schemes rank $(\text{MR-ZF}, p = 2) \approx (\text{MR-ZF}, p = 1) > (\text{MRT}, p = 1) \approx (\text{BST}) \approx (\text{ZF in [26]}) > (\text{MRT}, p = 2)$. Similar with the FC case, the MR-ZF scheme for the OSPS architecture achieves the best performance when increasing the blockage probability as shown in Fig. 6 (e) and Fig. 6 (f). As a brief summary w.r.t. the given scenario, for both architectures, when the channel SNR is weak and there is no blockage, we claim that the BST scheme is preferred since it is rather simple but adequate to achieve good performance. However, when the channel SNR is not too weak or there are potential blockages, the MR-ZF scheme with $p > 1$ outperforms the other schemes. As a side note, in practical implementation, the choice of p should not be too large since it plays a trade-off between blockage robustness, power spreading and the overhead for additional channel estimation.

B. Fully-Connected (FC) or One-Stream-Per-Subarray (OSPS)?

Note that the performance of different architectures highly depends on the channel condition and the underlying precoders. On top of the given scenario in this paper, we jointly evaluate the architecture performance in three aspects:

Training efficiency for the initial BA phase. Let P_D denote the detection probability, i.e., the probability of finding the strongest AoA-AoD pair between the BS and a generic UE. The BA results are illustrated in Fig. 7 (a). As a comparison, we also simulate a recent time-domain BA algorithm proposed in [42], which focuses on estimating the instantaneous channel coefficients with an orthogonal matching pursuit (OMP) technique. As we can see, the proposed BA scheme requires much less training overhead than that in [42]. In addition, due to the fact that the OSPS architecture has lower angular resolution and encounters larger sidelobe power leakage than the FC case, the former requires moderately ~ 10 more beacon slots than the latter for $P_D \geq 0.95$.

Spectral efficiency for the data communication phase. To compare the spectral efficiency of the two transmitter architectures as shown in Fig. 1, we consider a no-blockage scenario and focus on two precoding schemes, i.e., the simple BST scheme and the high-performance MR-ZF

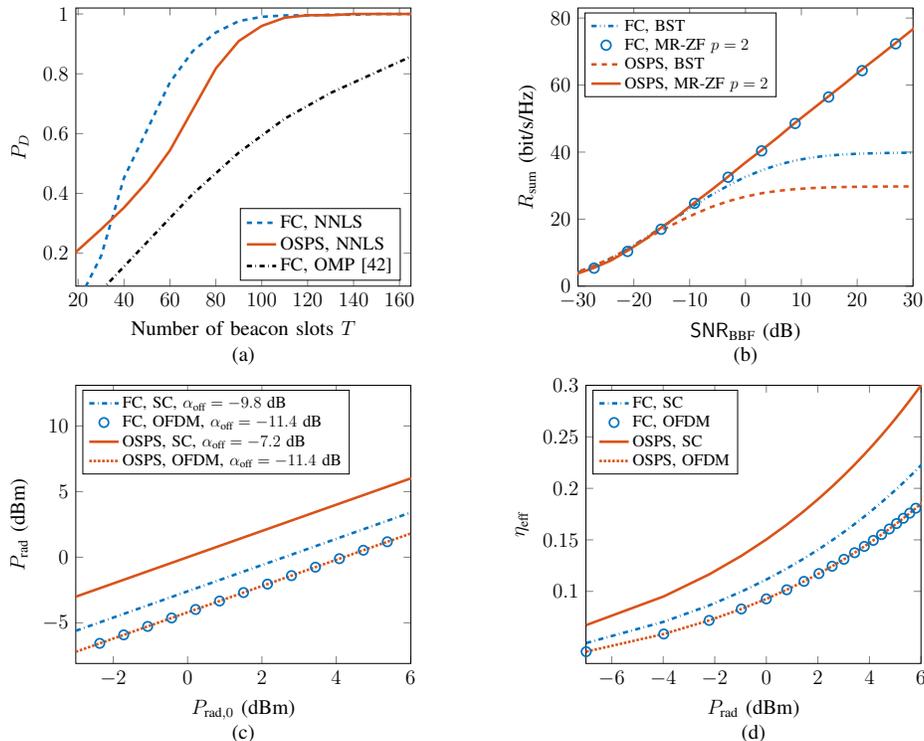


Fig. 7: The performance comparison of different transmitter architectures. (a) The initial BA detection probability vs. the training overhead with $\text{SNR}_{\text{BBF}} = -19$ dB. (b) The sum spectral efficiency vs. increasing SNR_{BBF} , without blockage. (c) The actual radiated power under *Option I* vs. the radiated power of the reference scenario. (d) The power efficiency under *Option II* vs. the actual radiated power.

scheme with $p = 2$. As we can see in Fig. 7 (b), in the range of $\text{SNR}_{\text{BBF}} \leq -10$ dB, which is more relevant in mmWave channels, all the 4 curves coincide with each other. Namely, for either the MR-ZF scheme or the BST scheme, the two architectures achieve a rather similar spectral efficiency. In contrast, when $\text{SNR}_{\text{BBF}} > -10$ dB, the MR-ZF scheme performs better. The two architectures with the MR-ZF precoding again achieve a rather similar performance.

Hardware power efficiency. To evaluate the architecture power efficiency, otherwise stated, we will consider the simple BST precoder. Also, since the modulation highly affects the power efficiency, we will take into account both the SC and the OFDM signaling in this section. We first assume a reference scenario as the baseline, i.e, the OSPS architecture using the BST precoder and a SC modulation. We use reference PAs with $P_{\text{max},0} = 6$ dBm and $\eta_{\text{max},0} = 0.3$. The backoff factor with respect to different waveforms and transmitter architectures can be written as $\alpha_{\text{off}} = 1/(P_{\text{PAPR}})$, where P_{PAPR} represents the PAPR of the input signal at a PA. The investigation for 3GPP LTE in [37] showed that with a probability of 0.999, the PAPR of the LTE SC waveform

(known as SC-FDMA) is smaller than ~ 7.2 dB and the PAPR of the LTE OFDM waveform (with 512 subcarriers employing QPSK) is smaller than ~ 11.4 dB. We set P_{PAPR} to these values for the OSPA architecture. For the FC architecture, however, the input signals of the PAs are the sum of the signals from different RF chains. Since each OFDM signal can be modeled as a Gaussian random process [37] and the signals from different RF chains are independent, the PAPR of the sum is the same as of one RF chain. For the case of SC signaling, there is no clear work in the literature that shows how the sum of SC signals behaves. We simulated the sum of $M_{\text{RF}} = 4$ SC signals using the same parameters as in [37]. The result shows that with probability of 0.999 the PAPR of the sum is smaller than ~ 9.8 dB. We apply these values and without loss of generality, we choose $\alpha_{\text{off},0} = -7.2$ dB as the reference scenario. As shown in (26), by deploying the same PAs (*Option I*), the two architectures achieve the same efficiency for a given P_{rad} . However, as illustrated in Fig. 7(c), the OSPA architecture with SC signaling (OSPA, SC) achieves the highest P_{rad} , followed by (FC, SC), (OSPA, OFDM), and (FC, OFDM). In contrast, by deploying different PAs (*Option II*)⁹, Fig. 7(d) shows that (OSPA, SC) achieves the highest power efficiency, followed by (FC, SC), (OSPA, OFDM) and (FC, OFDM).

To sum up, given the parameters in this paper, the two architectures achieve a similar sum spectral efficiency with certain precoders, but the OSPA architecture outperforms the FC case in terms of hardware complexity and power efficiency, only at the cost of a slightly longer latency for the initial BA.

C. Simulations Based on QuaDRiGa

In this section, we resort to the 3D geometry based channel generator QuaDRiGa [30] to show that our numerical results are quite consistent with practical mmWave communication channels.¹⁰ More precisely, we apply our BA and precoding schemes over $\sim 3 \times 10^5$ channel snapshots generated by QuaDRiGa. These channel snapshots correspond to a short segment of time evolution, where the BS is stationary and the speed of each UE along its moving direction is 1 m/s. The simulation results with respect to different transmitter architectures are shown in Fig. 8. As we can see from Fig. 8(a), for the initial BA with $P_D \geq 0.95$, the FC architecture requires ~ 10 less beacon slots than the OSPA case. Whereas, for the data communication phase

⁹Since the η_{max} of different PAs highly depends on the technology, for simplicity, we assume that different PAs working in their linear range have roughly the same maximum efficiency $\eta_{\text{max},0}$.

¹⁰Due to the QuaDRiGa generator limits, only the no-blockage scenario is considered in this section.

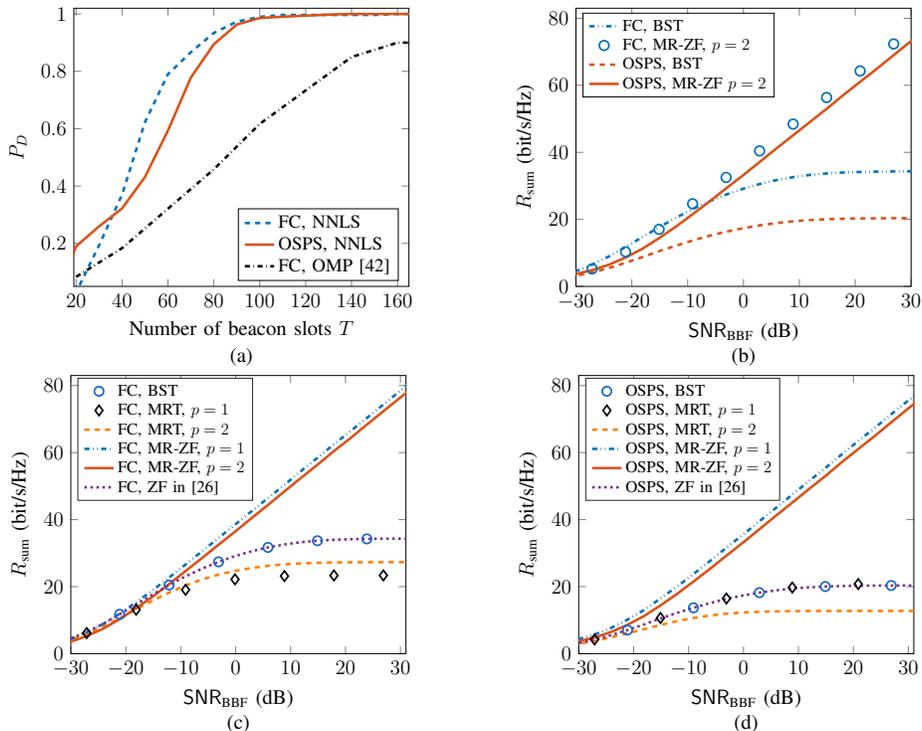


Fig. 8: The simulations based on QuaDRiGa: (a) The initial BA detection probability vs. the training overhead, with $\text{SNR}_{\text{BBF}} = -19$ dB. (b) The sum spectral efficiency of different transmitter architectures vs. increasing SNR_{BBF} . (c) The sum spectral efficiency of the FC architecture vs. increasing SNR_{BBF} . (d) The sum spectral efficiency of the OSPS architecture vs. increasing SNR_{BBF} .

as shown in Fig. 8 (b), by using either the BST or the MR-ZF precoder in the low SNR range ($\text{SNR}_{\text{BBF}} \leq -15$ dB), and using the MR-ZF precoder in the high SNR range ($\text{SNR}_{\text{BBF}} > -15$ dB), the two architectures achieve a quite similar performance. In addition, for both architectures as shown in Fig. 8 (c) and Fig. 8 (d), respectively, all the curves coincides with each other in the low SNR range, whereas the MR-ZF precoder outperforms the rest in the high SNR range. As we can see, all the results based on the QuaDRiGa generator are quite consistent with the results based on our proposed channel model. This consistency implies that our models, schemes, results and statements are not only theoretically reliable but also practically applicable.

VI. CONCLUSION

In this paper, we proposed an analysis framework to evaluate the performance of typical hybrid transmitters at mmWave frequencies. In particular, we focused on the

comparison of a fully-connected (FC) architecture and a partially-connected architecture with one-stream-per-subarray (OSPS) for a MU-MIMO base station using HDA beamforming. We jointly evaluated the performance of the two architectures in terms of the initial beam alignment (BA), the data communication, and the transmitter power efficiency. We used our recently proposed BA scheme and further proposed three simple precoding schemes on top of the effective channel after the BA. The precoding schemes are based on beam steering (BST), analog maximum ratio transmitting (MRT), and joint analog maximum ratio and baseband zero-forcing (MR-ZF), respectively. Particularly, both the BA scheme and the MR-ZF precoding scheme outperform the state-of-the-art counterparts in the literature. Given the parameters in this paper, our simulation results show that the two architectures achieve a similar sum spectral efficiency, but the OSPS architecture outperforms the FC case in terms of hardware complexity and power efficiency, only at the cost of a slightly longer latency for the initial BA. Therefore, the OSPS architecture emerges as a good choice for a simple and efficient design of MU-MIMO base stations operating at mmWave.

REFERENCES

- [1] K. Venugopal, A. Alkhateeb, N. G. Prelcic, and R. W. Heath, "Channel estimation for hybrid architecture-based wideband millimeter wave systems," *IEEE Journal on Selected Areas in Communications*, vol. 35, no. 9, pp. 1996–2009, 2017.
- [2] F. Sahrabi and W. Yu, "Hybrid digital and analog beamforming design for large-scale antenna arrays," *IEEE Journal of Selected Topics in Signal Processing*, vol. 10, no. 3, pp. 501–513, April 2016.
- [3] A. Li and C. Masouros, "Hybrid analog-digital millimeter-wave MU-MIMO transmission with virtual path selection," *IEEE Communications Letters*, vol. 21, no. 2, pp. 438–441, 2017.
- [4] S. S. Ioushua and Y. C. Eldar, "Hybrid analog-digital beamforming for massive MIMO systems," *arXiv preprint arXiv:1712.03485*, 2017.
- [5] J. Du, W. Xu, H. Shen, X. Dong, and C. Zhao, "Hybrid precoding architecture for massive multiuser MIMO with dissipation: sub-connected or fully connected structures?" *IEEE Transactions on Wireless Communications*, vol. 17, no. 8, pp. 5465–5479, 2018.
- [6] P. L. Cao, T. J. Oechtering, and M. Skoglund, "Precoding design for massive MIMO systems with sub-connected architecture and per-antenna power constraints," in *WSA 2018; 22nd International ITG Workshop on Smart Antennas*, March 2018, pp. 1–6.
- [7] R. W. Heath, N. Gonzalez-Prelcic, S. Rangan, W. Roh, and A. M. Sayeed, "An overview of signal processing techniques for millimeter wave MIMO systems," *IEEE journal of selected topics in signal processing*, vol. 10, no. 3, pp. 436–453, 2016.
- [8] X. Gao, L. Dai, and A. M. Sayeed, "Low RF-complexity technologies to enable millimeter-wave MIMO with large antenna array for 5G wireless communications," *IEEE Communications Magazine*, vol. 56, no. 4, pp. 211–217, APRIL 2018.
- [9] A. F. Molisch, V. V. Ratnam, S. Han, Z. Li, S. L. H. Nguyen, L. Li, and K. Haneda, "Hybrid beamforming for massive

- MIMO: A survey,” *IEEE Communications Magazine*, vol. 55, no. 9, pp. 134–141, 2017.
- [10] M. Majidzadeh, A. Moilanen, N. Tervo, H. Pennanen, A. Tölli, and M. Latva-aho, “Hybrid beamforming for single-user MIMO with partially connected RF architecture,” in *2017 European Conference on Networks and Communications (EuCNC)*, June 2017, pp. 1–6.
- [11] M. R. Castellanos, V. Raghavan, J. H. Ryu, O. H. Koymen, J. Li, D. J. Love, and B. Peleato, “Hybrid multi-user precoding with amplitude and phase control,” in *2018 IEEE International Conference on Communications (ICC)*, May 2018, pp. 1–6.
- [12] V. Raghavan, A. Partyka, A. Sampath, S. Subramanian, O. H. Koymen, K. Ravid, J. Cezanne, K. Mukkavilli, and J. Li, “Millimeter-wave MIMO prototype: Measurements and experimental results,” *IEEE Communications Magazine*, vol. 56, no. 1, pp. 202–209, 2018.
- [13] Z. Gao, L. Dai, and Z. Wang, “Channel estimation for mmwave massive MIMO based access and backhaul in ultra-dense network,” in *Communications (ICC), 2016 IEEE International Conference on*. IEEE, Conference Proceedings, pp. 1–6.
- [14] J. Rodríguez-Fernández, N. González-Prelcic, K. Venugopal, and R. W. Heath Jr, “Frequency-domain compressive channel estimation for frequency-selective hybrid mmWave MIMO systems,” *arXiv preprint arXiv:1704.08572*, 2017.
- [15] S. Haghighatshoar and G. Caire, “The beam alignment problem in mmWave wireless networks,” in *2016 50th Asilomar Conference on Signals, Systems and Computers*, Nov 2016, pp. 741–745.
- [16] X. Song, S. Haghighatshoar, and G. Caire, “A scalable and statistically robust beam alignment technique for mm-Wave systems,” *IEEE Trans. on Wireless Comm.*, vol. PP, pp. 1–1, 2018.
- [17] V. Va, J. Choi, and R. W. Heath, “The impact of beamwidth on temporal channel variation in vehicular channels and its implications,” *IEEE Transactions on Vehicular Technology*, vol. 66, no. 6, pp. 5014–5029, 2017.
- [18] X. Song, S. Haghighatshoar, and G. Caire, “A robust time-domain beam alignment scheme for multi-user wideband mmWave systems,” in *WSA 2018; 22th International ITG Workshop on Smart Antennas (to be published)*, March 2018, pp. 1–7.
- [19] —, “Efficient beam alignment for mmWave single-carrier systems with hybrid MIMO transceivers,” *IEEE Transactions on Wireless Communications*, 2019.
- [20] R. J. Weiler, M. Peter, W. Keusgen, and M. Wisotzki, “Measuring the busy urban 60 GHz outdoor access radio channel,” in *2014 IEEE International Conference on Ultra-WideBand (ICUWB)*, Sept 2014, pp. 166–170.
- [21] P. A. Eliasi, S. Rangan, and T. S. Rappaport, “Low-rank spatial channel estimation for millimeter wave cellular systems,” *IEEE Transactions on Wireless Communications*, vol. 16, no. 5, pp. 2748–2759, 2017.
- [22] O. El Ayach, R. W. Heath, S. Rajagopal, and Z. Pi, “Multimode precoding in millimeter wave MIMO transmitters with multiple antenna sub-arrays,” in *Global Communications Conference (GLOBECOM), 2013 IEEE*. IEEE, Conference Proceedings, pp. 3476–3480.
- [23] D. Zhang, Y. Wang, X. Li, and W. Xiang, “Hybridly connected structure for hybrid beamforming in mmWave massive MIMO systems,” *IEEE Transactions on Communications*, vol. 66, no. 2, pp. 662–674, 2018.
- [24] H.-L. Chiang, W. Rave, T. Kadur, and G. Fettweis, “Hybrid beamforming based on implicit channel state information for millimeter wave links,” *IEEE Journal of Selected Topics in Signal Processing*, vol. 12, no. 2, pp. 326–339, 2018.
- [25] V. Raghavan, S. Subramanian, J. Cezanne, A. Sampath, O. Koymen, and J. Li, “Directional hybrid precoding in millimeter-wave MIMO systems,” in *Global Communications Conference (GLOBECOM), 2016 IEEE*. IEEE, Conference Proceedings, pp. 1–7.
- [26] V. Raghavan, S. Subramanian, J. Cezanne, A. Sampath, O. H. Koymen, and J. Li, “Single-user versus multi-User precoding for millimeter wave MIMO systems,” *IEEE Journal on Selected Areas in Communications*, vol. 35, no. 6, pp. 1387–1401, June 2017.
- [27] X. Song, T. Kühne, and G. Caire, “Fully-connected vs. sub-connected hybrid precoding architectures for mmWave

- MU-MIMO,” in *2019 IEEE International Conference on Communications (ICC) (accepted)*.
- [28] N. N. Moghadam, G. Fodor, M. Bengtsson, and D. J. Love, “On the energy efficiency of MIMO hybrid beamforming for millimeter wave systems with nonlinear power amplifiers,” *arXiv preprint arXiv:1806.01602*, 2018.
- [29] T. Hälsig, D. Cvetkovski, E. Grass, and B. Lankl, “Statistical properties and variations of LOS MIMO channels at millimeter wave frequencies,” *arXiv preprint arXiv:1803.07768*, 2018.
- [30] S. Jaeckel, L. Raschkowski, K. Brner, and L. Thiele, “QuaDRiGa: A 3-D multi-cell channel model with time evolution for enabling virtual field trials,” *IEEE Transactions on Antennas and Propagation*, vol. 62, no. 6, pp. 3242–3256, 2014.
- [31] A. Gupta and R. K. Jha, “A Survey of 5G Network: Architecture and Emerging Technologies,” *IEEE Access*, vol. 3, pp. 1206–1232, 2015.
- [32] M. Agiwal, A. Roy, and N. Saxena, “Next Generation 5G Wireless Networks: A Comprehensive Survey,” *IEEE Communications Surveys Tutorials*, vol. 18, no. 3, pp. 1617–1655, thirdquarter 2016.
- [33] J. G. Proakis and M. Salehi, *Digital communications*. McGraw-Hill, 2008.
- [34] P. Bello, “Characterization of randomly time-variant linear channels,” *IEEE Transactions on Communications Systems*, vol. 11, no. 4, pp. 360–393, 1963.
- [35] A. Goldsmith, *Wireless communications*. Cambridge University Press, 2005.
- [36] A. M. Sayeed, “Deconstructing multiantenna fading channels,” *IEEE Transactions on Signal Processing*, vol. 50, no. 10, pp. 2563–2579, 2002.
- [37] H. G. Myung, J. Lim, and D. J. Goodman, “Peak-to-average power ratio of single carrier FDMA signals with pulse shaping,” in *Personal, Indoor and Mobile Radio Communications, 2006 IEEE 17th International Symposium on*. IEEE, Conference Proceedings, pp. 1–5.
- [38] C. Shepard, H. Yu, N. Anand, E. Li, T. Marzetta, R. Yang, and L. Zhong, “Argos: Practical many-antenna base stations,” in *Proceedings of the 18th annual international conference on Mobile computing and networking*. ACM, 2012, pp. 53–64.
- [39] T. L. Marzetta, E. G. Larsson, H. Yang, and H. Q. Ngo, *Fundamentals of massive MIMO*. Cambridge University Press, 2016.
- [40] E. Perahia, C. Cordeiro, M. Park, and L. L. Yang, “IEEE 802.11 ad: Defining the next generation multi-Gbps Wi-Fi,” in *Consumer Communications and Networking Conference (CCNC), 2010 7th IEEE*. IEEE, Conference Proceedings, pp. 1–5.
- [41] S. Jaeckel, L. Raschkowski, K. Börner, L. Thiele, and F. Burkhardt, “Quasi deterministic radio channel generator user manual and documentation,” *Fraunhofer Heinrich Hertz Institute Wireless Communications and Networks*, 2016.
- [42] K. Venugopal, A. Alkhateeb, R. W. Heath, and N. G. Prelcic, “Time-domain channel estimation for wideband millimeter wave systems with hybrid architecture,” in *Acoustics, Speech and Signal Processing (ICASSP), 2017 IEEE International Conference on*. IEEE, 2017, Conference Proceedings, pp. 6493–6497.