

Towards 6G MIMO: Massive Spatial Multiplexing, Dense Arrays, and Interplay Between Electromagnetics and Processing

Emil Björnson, *Fellow, IEEE*, Chan-Byoung Chae, *Fellow, IEEE*, Robert W. Heath Jr., *Fellow, IEEE*, Thomas L. Marzetta, *Life Fellow, IEEE*, Amine Mezghani, *Member, IEEE*, Luca Sanguinetti, *Senior Member, IEEE*, Fredrik Rusek, *Member, IEEE*, Miguel R. Castellanos, *Member, IEEE*, Dongsoo Jun, *Graduate Student Member, IEEE*, and Özlem Tuğfe Demir, *Member, IEEE*

Abstract—The increasing demand for wireless data transfer has been the driving force behind the widespread adoption of Massive MIMO (multiple-input multiple-output) technology in 5G. The next-generation MIMO technology is now being developed to cater to the new data traffic and performance expectations generated by new user devices and services in the next decade. The evolution towards “ultra-massive MIMO (UM-MIMO)” is not only about adding more antennas but will also uncover new propagation and hardware phenomena that can only be treated by jointly utilizing insights from the communication, electromagnetic (EM), and circuit theory areas.

This article offers a comprehensive overview of the key benefits of the UM-MIMO technology and the associated challenges. It explores massive multiplexing facilitated by radiative near-field effects, characterizes the spatial degrees-of-freedom, and practical channel estimation schemes tailored for massive arrays. Moreover, we provide a tutorial on EM theory and circuit theory, and how it is used to obtain physically consistent antenna and channel models. Subsequently, the article describes different ways to implement massive and dense antenna arrays, and how to co-design antennas with signal processing. The main open research challenges are identified at the end.

Index Terms—Ultra Massive MIMO, extremely large-scale aperture array, massive spatial multiplexing, electromagnetic theory for communication, antenna array design.

E. Björnson is with the School of Electrical Engineering and Computer Science, KTH Royal Institute of Technology, Stockholm, Sweden. E-mail: emilbjo@kth.se.

C.-B. Chae and D. Jun are with the School of Integrated Technology, Yonsei University, Seoul 03722, South Korea. E-mail: {cbchae, dongsoo.jun}@yonsei.ac.kr.

R. W. Heath, Jr. is with the Department of Electrical and Computer Engineering, University of California, San Diego, La Jolla, CA 92161, USA. E-mail: rwheathjr@ucsd.edu.

T. L. Marzetta is with the Department of Electrical and Computer Engineering, New York University, Brooklyn, NY 11201, USA. E-mail: tom.marzetta@nyu.edu.

A. Mezghani is with the Department of Electrical and Computer Engineering, University of Manitoba, Winnipeg, MB R3T 5V6, Canada. E-mail: amine.mezghani@umanitoba.ca.

L. Sanguinetti is with the Dipartimento di Ingegneria dell’Informazione, University of Pisa, Pisa, Italy. E-mail: luca.sanguinetti@unipi.it.

F. Rusek is with Sony Europe B.V., Lund, Sweden. E-mail: fredrik.x.rusek@sony.com.

M. R. Castellanos is with the Department of Electrical and Computer Engineering, North Carolina State University, Raleigh, NC 27695, USA. E-mail: mrcastel@ncsu.edu.

Ö. T. Demir is with the TOBB University of Economics and Technology, Ankara, Türkiye. E-mail: ozlemtugfedemir@etu.edu.tr.

The work of E. Björnson is supported by Swedish Research Council and Swedish Foundation for Strategic Research. The work of C.-B. Chae and D. S. Jun is in part supported by the Institute of Information and Communication Technology Promotion (IITP) grant funded by the Ministry of Science and ICT (MSIT), Korea (2021-0-00486, 2021-0-02208). The work of R. W. Heath, Jr. was supported in part by the National Science Foundation under grant nos. NSF-ECCS-2153698, NSF-CCF-2225555, NSF-CNS-2147955 and is supported in part by funds from federal agency and industry partners as specified in the Resilient & Intelligent NextG Systems (RINGS) program, as well as by support from Nokia, Samsung and Qualcomm. The work of T. L. Marzetta was supported by NYU WIRELESS. The work of A. Mezghani was in part supported by the Natural Sciences and Engineering Research Council of Canada as well as Research Manitoba. The work of L. Sanguinetti was partially supported by the Italian Ministry of Education and Research (MUR) in the framework of the FoReLab project (Departments of Excellence). The work of M. R. Castellanos was supported in part by Nokia, Motorola, and Samsung. The work of Ö. T. Demir was supported by 2232-B International Fellowship for Early Stage Researchers Programme funded by the Scientific and Technological Research Council of Türkiye.

I. INTRODUCTION

THE need for equipping transmitters and receivers with multiple antennas in wireless communication systems has been recognized for over a century. The first observed benefit was the adaptive directivity achievable by controlling the constructive and destructive superposition of electromagnetic (EM) signals using an antenna array [1], [2]. The transmitter can use this feature, traditionally referred to as *beamforming*, to focus a transmitted signal at the desired receiver while avoiding interference at specific locations. Similarly, the receiver can amplify signals impinging from a particular direction using multiple antennas while suppressing undesired interference. The second observed benefit was the higher robustness against channel fading achieved by using multiple antennas [3]–[6], as it becomes less likely that all transmit-receive antenna pairs experience deep fades simultaneously as we increase the number of antennas and the array size. This feature is called *spatial diversity* and *channel hardening* [7]. The third and most recently discovered benefit is multiple-input multiple-output (MIMO) communications [8]–[12], where antenna arrays are used to spatially multiplex many layers of data at the same time and frequency. This can be done in multi-user MIMO mode, where a multiple-antenna base station (BS) communicates with multiple user equipments (UEs) simultaneously. This is enabled using adaptive beamforming: the BS gives each transmitted signal a different spatial directivity, has the ability to amplify signals received from UEs in different directions, and can filter out interference in both transmission directions. There is also the single-user

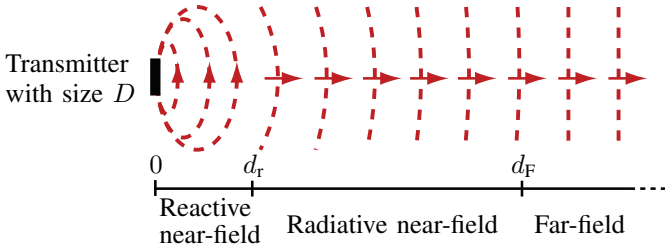


Fig. 1. The EM field looks different depending on the distance from the transmitting aperture antenna. The wavefront is almost planar in the far-field, while the spherical curvature is clearly noticeable in the radiative near-field but not reactive effects such as inductive coupling and evanescent waves.

MIMO mode, where a multi-antenna BS and multi-antenna UE exchange multiple data layers simultaneously by beamforming through different propagation paths.

The MIMO technology was first introduced in cellular and WiFi networks as a premium feature but is nowadays a mainstream technology. The 5G technology was built around the *Massive MIMO* concept [13] of having a surplus of antennas at the BS compared to the UE side, which makes it practically feasible to protect the data layers from mutual interference through spatial filtering, even under imperfect channel state information (CSI) and hardware impairments [14], [15]. A typical 5G BS in 2023 had 64 antenna ports and can support up to 16 data layers, such as 8 UEs assigned with two layers each. The driving force behind the MIMO adoption is the rapidly increasing demand for data traffic in cellular networks, currently growing by 40% per year [16]. The 5G MIMO technology, particularly in the 3.5 GHz band, can supply current BS sites with significantly higher capacity than in 4G, to support higher speeds per device and accommodate more simultaneously served devices. If the data traffic continues to grow at the current pace over the next ten years, 6G technology must deliver $1.4^{10} \approx 30$ times higher capacity than current networks. New emerging user devices (e.g., for augmented reality) and services (e.g., federated learning, hyper-reliable and low-latency communication) might create an even faster wireless traffic growth; thus, the next-generation MIMO technology should be developed to at least support 100 times higher capacity than in current networks. A portion of that can be achieved by expanding the bandwidth. However, since spectrum is scarce in bands suitable for wide-area coverage, the focus should be on increasing the sum spectral efficiency (SSE) [bit/s/Hz]. The SSE is the total data transmitted per second and per Hertz among all the spatially multiplexed layers.

The spectral efficiency (SE) per spatial layer is fundamentally limited by the signal-to-noise ratio (SNR), as expressed by the Shannon formula $\log_2(1 + \text{SNR})$ bit/s/Hz [17]. The logarithmic nature of this function places a fundamental constraint on massive improvements through the use of multiple antennas, except for UEs experiencing very low SNRs. However, the spectral efficiency (SSE) of ν data layers of this kind is upper-bounded by $\nu \cdot \log_2(1 + \text{SNR})$, a linearly increasing

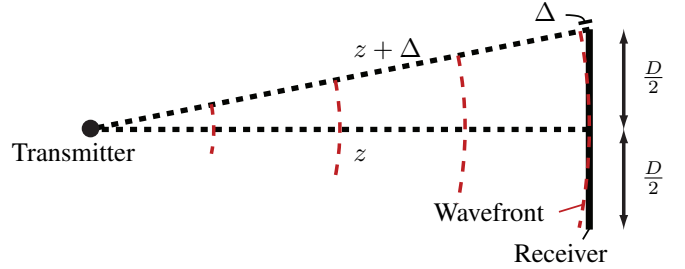


Fig. 2. The curvature of an impinging spherical wave creates a delay $\frac{\Delta}{c}$ between the center of the receiver and the edge. The delay turns into a phase-shift of $2\pi f_c \frac{\Delta}{c} = \frac{2\pi}{\lambda} \Delta$.

function of ν . To enhance SSE in 6G, the logical approach is to employ “more MIMO”—increasing spatial multiplexing with more layers to serve additional UEs. On the BS side, this involves using larger antenna arrays, either physically larger or relative to the wavelength (if the carrier frequency is increased compared to 5G). At first glance, this evolution may seem like an engineering challenge: assembling more antennas using current technology and applying the same algorithms found in textbooks such as [15], [18] with larger-dimensional matrices.

However, the reality is vastly different on multiple levels. Antenna array design in 6G requires not only a drastic increase in the number of antennas but also fundamental changes in EM properties. Despite the numerous antennas in 5G Massive MIMO systems, the aperture size is small enough to neglect near-field effects, focusing primarily on the far-field. In 6G, the significantly increased number of antennas and higher frequency bands will expand the near-field region, especially with dense BS deployments in urban and indoor environments. Consequently, MIMO systems with extremely larger antenna arrays, referred to as X-MIMO in industry terminology, are gaining attention. Various terms, such as extremely large aperture arrays (ELAA) [19], extremely large-scale MIMO (XL-MIMO) [20]–[23], and *Ultra-Massive MIMO (UM-MIMO)*, have been suggested in academic literature. In this paper, we use the latter terminology.

In addition to near-field propagation becoming more dominant, we approach fundamental limits on spatial degrees-of-freedom; more physically accurate models of channels, antennas, and hardware effects are necessary; and new implementation challenges emerge. These aspects will be further described in the remainder of this paper.

The remainder of the paper is organized as follows. Section II provides an overview of radiative near-field propagation effects and how they can be exploited for finite-depth beamforming and massive spatial multiplexing. Section III introduces the spatial degrees-of-freedom concept, which is the physical limit on the multiplexing capability of an array. We then describe efficient channel estimation techniques in Section IV, with a focus on exploiting array geometry and propagation characteristics. Having introduced the basic concepts, the paper then provides a tutorial on the underlying theory. Section V provides a linear system approach to EM theory. Section VI expands on this approach by using circuit theory to obtain a physically consistent end-to-end MIMO

channel representation. Section VII describes how to account for realistic antenna properties, such as mutual coupling, polarization, and near-field propagation for MIMO array modeling. Next, Section VIII describes four antenna array architectures that might be used in 6G UM-MIMO systems. Section IX takes a closer look at how the hybrid array architecture can be optimized jointly with the signal processing algorithms. Finally, the paper is concluded in Section X by describing some open research challenges.

II. MIMO COMMUNICATION IN THE RADIATIVE NEAR-FIELD

The behavior of a wireless channel is governed by Maxwell's equations, which can be solved to determine the electric and magnetic field distributions that a transmission creates at different locations in a specific environment. This fundamental approach often results in very complex expressions, but they can be fortunately simplified and tailored to the specific scenario, allowing for practical use in system design and optimization. As contemporary MIMO systems expand in size, the once-accurate simplified propagation models must be enriched to account for "new" phenomena that were always present but previously considered negligible. This section provides a tutorial on such properties, starting from the near- and far-fields of antennas and arrays, and subsequently describing the impact they have on future UM-MIMO communication systems. We will go deeper into many of these aspects in later sections.

We begin by considering the transmission from a point source. The electric field observed at a distance z from the source, in any direction perpendicular to the propagation direction, is proportional to [24]

$$\frac{-j\eta e^{-j\frac{2\pi}{\lambda}z}}{2\lambda z} \left(1 + \frac{j}{2\pi z/\lambda} - \frac{1}{(2\pi z/\lambda)^2} \right), \quad (1)$$

where η denotes the impedance of free space and λ is the signal's wavelength. The first term in (1) has a squared magnitude that decays proportionally to $1/z^2$, consistent with the classical pathloss behavior for free-space propagation [25]. The second factor in (1) is often overlooked in communications—for good reasons—because

$$\left| 1 + \frac{j}{2\pi z/\lambda} - \frac{1}{(2\pi z/\lambda)^2} \right|^2 = 1 - \frac{1}{(2\pi z/\lambda)^2} + \frac{1}{(2\pi z/\lambda)^4}, \quad (2)$$

which quickly approaches 1 when the propagation distance z increases. Already at a distance of 2λ from the transmitter, (2) equals 0.99. The region $z \geq 2\lambda$ where (1) can be approximated as $\frac{-j\eta e^{-j\frac{2\pi}{\lambda}z}}{2\lambda z}$ is known as the far-field, while $z < 2\lambda$ represents the near-field region. There is no strict boundary between these regions since (2) approaches 1 in a continuous manner.

If we replace the point source with a transmitting aperture antenna having a maximum length D that is larger than the wavelength, then the distances that characterize the near- and far-field change as well. In particular, the near-field can be divided into two parts: the *reactive* and *radiative* near-field. For most antenna types, terms similar to the second factor in (1) must be taken into account at distances $z \leq d_r = 0.62\sqrt{D^3/\lambda}$

[26], [27]. These terms represent EM components that remain around the transmitter instead of being radiated; for example, related to inductive coupling and evanescent fields. It is in the reactive near-field that the induction-based near-field communication (NFC) technology operates and it is commonly used for tags and keycards. This paper does not consider such technologies. A transmitting aperture antenna has similar far-field behavior as a point source when observed at a distance $z > d_F = \frac{2D^2}{\lambda}$, where d_F is greater than d_r for D larger than λ . This boundary is known as the *Fraunhofer distance* or *Rayleigh distance*. The far-field region is also known as the *Fraunhofer region*. In between the mentioned distance boundaries is a range $d_r < z \leq d_F$ that is called the *radiative near-field* or *Fresnel region*. Fig. 1 illustrates these different regions and emphasizes the core difference between the radiative near-field and conventional far-field: the curvature of the wavefront. It is spherical in both cases, which means that the received signal power decays with the distance z proportionally to $1/z^2$, but the curvature is only noticeable in the radiative near-field.

The implications of the strongly curved wavefront in the radiative near-field are easier to comprehend when considering an aperture antenna with the length D that receives a signal from a transmitting point source. This setup is shown in Fig. 2 for a transmitter at the distance z in the broadside direction. When the wavefront reaches the receiver's center, a distance Δ remains until it reaches the edge. The extra distance can be calculated as

$$\Delta = \sqrt{z^2 + \frac{D^2}{4}} - z = z\sqrt{1 + \frac{D^2}{4z^2}} - z \approx \frac{D^2}{8z}, \quad (3)$$

where the last expression follows from the first-order Taylor approximation $\sqrt{1+x} \approx 1 + \frac{x}{2}$, that is accurate for small x (i.e., when $z \gg \frac{D}{2}$). The Fraunhofer distance $z = d_F$, where $d_F \gg \frac{D}{2}$ so that we can use the above approximation, gives rise to the phase-shift

$$2\pi f_c \frac{\Delta}{c} = \frac{2\pi}{\lambda} \Delta \approx \frac{2\pi}{\lambda} \frac{D^2}{8d_F} = \frac{\pi}{8}, \quad (4)$$

where f_c is the carrier frequency and c is the speed of the EM radiation. There is nothing deeper behind the Fraunhofer distance than the fact that $\cos(\pi/8) \approx 0.92 \approx 1$ so that the spherical curvature creates a tiny phase variation over the antenna [28].

The spherical curvature also gives rise to power variations over the receive antenna. Since the squared magnitude of the EM field is proportional to $1/z^2$, the power difference is

$$\frac{z^2}{(z + \Delta)^2} \approx \frac{z^2}{\left(z + \frac{D^2}{8z}\right)^2} \quad (5)$$

between the center and edge. The power variations are negligible if $z = d_B = 2D$ [29] so that the propagation distance is twice as large as the surface, in which case (5) becomes 0.94. We notice that $d_F = d_B \frac{D}{\lambda}$, which means that the phase-variations are more prevalent when considering large antennas.

The aforementioned approximations are conventionally made in wireless communications without further discussion, but there are ways to be more precise. The gain of a receiving

aperture antenna can be computed by taking an integral of the impinging field over the aperture. For example, if the antenna spans the interval $x \in [-a/2, a/2]$, $y \in [-b/2, b/2]$ in the xy -plane and the electric field is denoted as $E(x, y)$, then the gain (relative to an isotropic reference antenna in the far-field) can be computed as [30], [31]

$$G = \frac{\left| \int_{-a/2}^{a/2} \int_{-b/2}^{b/2} E(x, y) dx dy \right|^2}{A_{\text{iso}} \int_{-a/2}^{a/2} \int_{-b/2}^{b/2} |E(x, y)|^2 dx dy}, \quad (6)$$

where $A_{\text{iso}} = \frac{\lambda^2}{4\pi}$ is the effective area of an isotropic antenna. Looking at this expression, it might seem logical that a larger antenna has a higher gain because of the expanded integration intervals, but this is not a necessity because phase variations in $E(x, y)$ over the antenna impact the numerator in (6). This phenomenon occurs even in the far-field when the impinging signal arrives from a non-broadside direction, creating linear phase variation similar to that of a plane wave. However, the effect is particularly dominant in the radiative near-field due to the wave's spherical curvature. For example, for a transmitting point source in the broadside direction $(0, 0, z)$, with $z > \max(d_r, d_B)$, the electric field observed at the receiver can be expressed as

$$E(x, y) = \frac{E_0}{\sqrt{4\pi z}} e^{-j\frac{2\pi}{\lambda} \sqrt{x^2 + y^2 + z^2}}, \quad (7)$$

where E_0 is a scaling factor for the electric intensity and $\sqrt{x^2 + y^2 + z^2}$ is the distance from the transmitter to the location $(x, y, 0)$ at the receiver. This expression assumes that the power variations over the antenna are negligible, but there are still phase variations across the antenna aperture when $\max(d_r, d_B) < z \leq d_F$. By substituting (7) into (6), the gain simplifies to

$$G = \frac{\left| \int_{-a/2}^{a/2} \int_{-b/2}^{b/2} e^{-j\frac{2\pi}{\lambda} \sqrt{x^2 + y^2 + z^2}} dx dy \right|^2}{A_{\text{iso}} ab}, \quad (8)$$

which becomes ab/A_{iso} if the phase in the numerator is approximately constant over the antenna aperture. This is a classical approximation that is valid for $z > d_F$ because the electric field can then be approximately as $E(x, y) \approx \frac{E_0}{\sqrt{4\pi z}} e^{-j\frac{2\pi}{\lambda} z}$. The ratio ab/A_{iso} between antenna areas is how the antenna gain is normally computed in the far-field. By contrast, a strictly smaller gain is obtained for $z < d_F = \frac{2(a^2 + b^2)}{\lambda}$ (with $D = \sqrt{a^2 + b^2}$ being the antenna's diagonal length) because of the noticeable spherical phase variations that make some parts of the impinging field cancel other parts. The way to circumvent that effect is to use an array of many small antennas instead of one large antenna so that each antenna achieves its physically maximum gain and then the gains are superimposed through receiver processing. In other words, we want to use small transmit and receive antennas that are in each others' far-field, but build large antenna arrays that can observe radiative near-field effects across the array.

Fig. 3 illustrates this situation for an array of 10×10 $\lambda/2$ -spaced antennas in the xy -plane, and the coloring shows the (normalized) real part $\cos(2\pi\sqrt{x^2 + y^2 + z^2}/\lambda)$ of the impinging electric field when the transmitter is 8λ from the

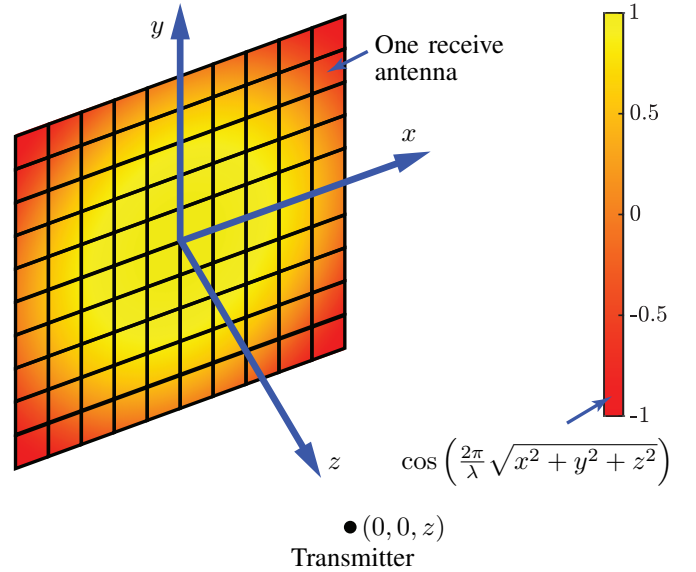


Fig. 3. The (normalized) real part of the electric field in (7) is shown for an antenna array deployed in the xy -plane. The transmitter is located in a broadside direction in the radiative near-field, leading to spherical phase variations.

receiver. The spherical-shaped phase variations result in values between $+1$ to -1 . If there were only one big antenna covering the entire surface, the red parts would cancel the yellow parts when computing (8), leading to a large loss in antenna gain. The gain of this large antenna would be only 35% of the maximum gain ab/A_{iso} when $a = b = 10\frac{\lambda}{2} = 5\lambda$, and the percentage shrinks as the antenna size increases. However, in Fig. 3, the surface is divided into 100 antennas so that the phase variation is small across each individual antenna. This division results in a negligible gain loss per antenna so that at least 95% of the maximum gain is achieved. This comes at the expense of requiring receiver hardware that can combine the antenna signals, which we anyway need for beamforming in mobile scenarios.

When the array is large, there can also be power variations between the antennas deployed at the surface, caused by having substantially different distances (and angles) to the transmitter. In contrast to the phase-shifts, this geometric pathloss effect cannot be mitigated through signal processing or making the antenna elements smaller. Fortunately, the effect is negligible in many practical situations. If we let D denote the largest dimension of the array (i.e., the diagonal in Fig. 3), then it follows that $d_F = d_B \frac{D}{\lambda} \gg d_B$ for large arrays with $D \gg \lambda$. This implies that there is a wide distance range $d_B < z \leq d_F$ where the power variations over the antenna array are negligible but not the phase variations. We will call this the *Fresnel region*. As an example, consider a half-wavelength-spaced BS array with the size 1×0.5 m, and operating at 30 GHz (i.e., $\lambda = 0.01$ m). This array contains 5000 antennas in the configuration 100×50 . It then follows that $d_B = 2D = 2\sqrt{1^2 + 0.5^2} \approx 2.24$ m so most UEs will be located beyond that distance, while they might be closer than the Fraunhofer distance that becomes $d_F = \frac{2D^2}{\lambda} =$

$\frac{2(1^2+0.5^2)}{0.01} = 250$ m. We will focus on the Fresnel region range $d_B < z \leq d_F$ in the remainder of this section, but note that there is one particular situation when the power variations over the antenna array cannot be neglected: when studying the asymptotic limits of large arrays because otherwise one obtains implausible results where more power is received than was transmitted [29]. One must also accommodate for such effects when modeling non-line-of-sight propagation environments [32], where different parts of a large array might see different scattering objects.

A. Beamfocusing in the radiative near-field of arrays

Suppose the antenna array is used to transmit a signal to a particular receiver. When M antennas transmit the signal with phase-shifts that create constructive interference at the receiver's location, the received power becomes M times larger than when transmitting with the same power from a single antenna. This is called the *array gain* or *beamforming gain* and is achievable in line-of-sight scenarios at any propagation distance $z > d_B = 2D$ (i.e., there are no pathloss variations over the array). The gain is the same in the far-field and radiative near-field, but there is an essential difference in how the radiated signal behaves at other locations, such as the shape of the focus area around the receiver. These geometric properties will be analyzed in this section.

Consider a uniform square array with M antennas arranged as $N \times N$ antennas, where $N = \sqrt{M}$ is an integer and Δ denotes the antenna spacing. We let $n \in \{1, \dots, N\}$ be the antenna index in the x dimension and $m \in \{1, \dots, N\}$ be the index in the y dimension. The antenna with index (n, m) is centered at the point $(\bar{x}_n, \bar{y}_m, 0)$, where

$$\bar{x}_n = \left(n - \frac{N+1}{2}\right) \Delta, \quad (9)$$

$$\bar{y}_m = \left(m - \frac{N+1}{2}\right) \Delta. \quad (10)$$

We consider an isotropic broadside receiver at the location $(0, 0, z)$. If each transmit antenna has the gain G towards the receiver, the channel coefficient becomes

$$h_{n,m} = \frac{\lambda\sqrt{G}}{4\pi z} e^{-j\frac{2\pi}{\lambda}\sqrt{\bar{x}_n^2 + \bar{y}_m^2 + z^2}} \quad (11)$$

$$\approx \frac{\lambda\sqrt{G}}{4\pi z} e^{-j\frac{2\pi}{\lambda}\left(z + \frac{\bar{x}_n^2}{2z} + \frac{\bar{y}_m^2}{2z}\right)}, \quad (12)$$

where (11) is the exact expression and the first-order Taylor approximation $\sqrt{1+x} \approx 1 + \frac{x}{2}$ can be utilized to simplify it to (12) since z is substantially larger than \bar{x}_n and \bar{y}_m when $z > d_B$. This is known as the Fresnel approximation [33].

If the data symbol $s \in \mathbb{C}$ is transmitted with power p , the received signal is

$$r = \sum_{n=1}^N \sum_{m=1}^N h_{n,m} \frac{e^{j\psi_{n,m}}}{\sqrt{M}} s + w_{n,m}, \quad (13)$$

where $w_{n,m} \sim \mathcal{N}_{\mathbb{C}}(0, \sigma^2)$ is independent complex Gaussian receiver noise, $\psi_{n,m}$ is the phase-shift assigned at antenna

(n, m) , and $1/\sqrt{M}$ divides the total power equally among the transmit antennas. The SNR becomes

$$\begin{aligned} \text{SNR} &= \frac{p}{\sigma^2} \left| \sum_{n=1}^N \sum_{m=1}^N h_{n,m} \frac{e^{j\psi_{n,m}}}{\sqrt{M}} \right|^2 \\ &= \frac{p}{\sigma^2} \frac{\lambda^2 G}{(4\pi z)^2} \frac{1}{M} \underbrace{\left| \sum_{n=1}^N \sum_{m=1}^N e^{-j\frac{2\pi}{\lambda}\sqrt{\bar{x}_n^2 + \bar{y}_m^2 + z^2}} e^{j\psi_{n,m}} \right|^2}_{=\text{AG}}, \end{aligned} \quad (14)$$

where AG denotes the array gain. It becomes $N^4/M = M$ when $\psi_{n,m} = \frac{2\pi}{\lambda}\sqrt{\bar{x}_n^2 + \bar{y}_m^2 + z^2}$ so that the transmitter cancels all the phase-shifts. This is the maximum array gain.

Suppose the transmitter instead focuses the signal on a point at the distance z but with some other small angle φ measured from the boresight in the horizontal plane. The focus point is at $(z \sin(\varphi), 0, z \cos(\varphi))$ and is obtained when the transmitter assigns the phase-shifts

$$\begin{aligned} \psi_{n,m} &= \frac{2\pi}{\lambda} \sqrt{(\bar{x}_n - z \sin(\varphi))^2 + \bar{y}_m^2 + z^2 \cos^2(\varphi)} \\ &= \frac{2\pi}{\lambda} \sqrt{\bar{x}_n^2 + \bar{y}_m^2 + z^2 - 2z\bar{x}_n \sin(\varphi)}. \end{aligned} \quad (15)$$

The array gain at the original receiver can then be computed, using the Fresnel approximation, as

$$\begin{aligned} &\frac{1}{M} \left| \sum_{n=1}^N \sum_{m=1}^N e^{-j\frac{2\pi}{\lambda}\sqrt{\bar{x}_n^2 + \bar{y}_m^2 + z^2}} e^{j\frac{2\pi}{\lambda}\sqrt{\bar{x}_n^2 + \bar{y}_m^2 + z^2 - 2z\bar{x}_n \sin(\varphi)}} \right|^2 \\ &\approx \frac{1}{M} \left| \sum_{n=1}^N \sum_{m=1}^N e^{-j\frac{2\pi}{\lambda}\left(z + \frac{\bar{x}_n^2}{2z} + \frac{\bar{y}_m^2}{2z}\right)} e^{j\frac{2\pi}{\lambda}\left(z + \frac{\bar{x}_n^2}{2z} + \frac{\bar{y}_m^2}{2z} - \bar{x}_n \sin(\varphi)\right)} \right|^2 \\ &= \frac{N^2}{M} \left| \sum_{n=1}^N e^{-j\frac{2\pi}{\lambda}\left(n - \frac{N+1}{2}\right)\Delta \sin(\varphi)} \right|^2 \\ &\approx \frac{N^2}{M} \left| \int_0^N e^{-j\frac{2\pi}{\lambda}n\Delta \sin(\varphi)} dn \right|^2 = M \text{sinc}^2\left(\frac{1}{\lambda}N\Delta \sin(\varphi)\right), \end{aligned} \quad (16)$$

where we also approximated the summation over many antennas by the corresponding integral. This approximation is tight when the antenna spacing is small, similar to how a Riemann sum approaches a Riemann integral. This array gain expression depends on the angle φ but is independent of the propagation distance, which implies that it is the same in the radiative near-field as in the far-field. The squared sinc-function determines how the array gain observed at the receiver tapers off when the transmitter aims the signal at a different location at the same distance. Since $\text{sinc}^2(0.443) \approx 0.5$, half the array gain is achieved at $\varphi = \pm \arcsin\left(\frac{0.443\lambda}{N\Delta}\right) \approx \pm \frac{0.443\lambda}{N\Delta}$. The half-power angular beamwidth then becomes

$$\text{BW}_{3\text{dB}} \approx \frac{0.886\lambda}{N\Delta} \text{ rad}, \quad (17)$$

which is inversely proportional to the width $N\Delta$ of the array and proportional to the wavelength. Although the angular

beamwidth is the same at all propagation distance for which $z > d_B$, the physical beamwidths (in meters) is approximately $BW_{3\text{dB}}z$; thus, the width of the beam around the receiver is proportional to the distance to the receiver and much smaller in the radiative near-field than in the far-field.

When the antenna spacing is $\Delta = \lambda/2$, the beamwidth expression simplifies to $BW_{3\text{dB}} = \frac{1.772}{N}$. It is clear that more antennas per dimension leads to a narrower beamwidth.

The radiated signal from an array is often illustrated as a cone with an angular width of $BW_{3\text{dB}}$. However, this description is incomplete because the array gain also tapers off in the depth domain. We can characterize this phenomenon similarly to the beamwidth analysis by supposing that the transmitting array focuses the emitted signal on another point $(0, 0, F)$ in the same direction but at a distance $F \neq z$. Let us define the focal point deviation

$$z_{\text{eff}} = \left| \frac{1}{F} - \frac{1}{z} \right|^{-1} = \frac{Fz}{|F - z|}. \quad (18)$$

The array gain at the original receiver can then be computed, using the Fresnel approximation, as [31]

$$\begin{aligned} & \frac{1}{M} \left| \sum_{n=1}^N \sum_{m=1}^N e^{-j\frac{2\pi}{\lambda} \sqrt{\bar{x}_n^2 + \bar{y}_m^2 + z^2}} e^{j\frac{2\pi}{\lambda} \sqrt{\bar{x}_n^2 + \bar{y}_m^2 + F^2}} \right|^2 \\ & \approx \frac{1}{M} \left| \sum_{n=1}^N \sum_{m=1}^N e^{-j\frac{2\pi}{\lambda} \left(z + \frac{\bar{x}_n^2}{2z} + \frac{\bar{y}_m^2}{2z} \right)} e^{j\frac{2\pi}{\lambda} \left(F + \frac{\bar{x}_n^2}{2F} + \frac{\bar{y}_m^2}{2F} \right)} \right|^2 \\ & \approx \frac{1}{M} \left| \int_{-N/2}^{N/2} e^{j\frac{\pi}{\lambda} \frac{n^2 \Delta^2}{z_{\text{eff}}}} dn \int_{-N/2}^{N/2} e^{j\frac{\pi}{\lambda} \frac{m^2 \Delta^2}{z_{\text{eff}}}} dm \right|^2 \\ & = M \left(\frac{8z_{\text{eff}}}{d_F} \right)^2 \left(C^2 \left(\sqrt{\frac{d_F}{8z_{\text{eff}}}} \right) + S^2 \left(\sqrt{\frac{d_F}{8z_{\text{eff}}}} \right) \right)^2, \quad (19) \end{aligned}$$

where $d_F = 4N^2\Delta^2/\lambda$ is the Fraunhofer distance of the considered array, while $C(x) = \int_0^x \cos(\pi t^2/2) dt$ and $S(x) = \int_0^x \sin(\pi t^2/2) dt$ denote the Fresnel integrals.

The array gain expression in (19) has the structure $A(x) = (C^2(\sqrt{x}) + S^2(\sqrt{x}))/x^2$, where $x = d_F/(8z_{\text{eff}})$. This is a decreasing function for $x \in [0, 2]$ with $A(0) = 1$ and $A(1.25) \approx 0.5$. Hence, half the array gain is achieved when

$$1.25 = \frac{d_F}{8z_{\text{eff}}} = \frac{d_F|F - z|}{8Fz} \rightarrow z = \frac{d_FF}{d_F \pm 10F}. \quad (20)$$

This implies that when the transmitter focuses a signal on the point $(0, 0, F)$, the array gain will only be large at some of the potential locations $(0, 0, z)$ in the same direction. The smallest solution is $z = \frac{d_FF}{d_F + 10F} < F$, which is the starting point of the beam in the depth domain. If $F < d_F/10$, we also have the solution $z = \frac{d_FF}{d_F - 10F}$, which marks the end of the beam [31]. Hence, for transmission to receivers at locations $d_F < F < d_F/10$ in the radiative near-field, the transmitted signal behaves as a beam with finite half-power beamdepth:

$$BD_{3\text{dB}} = \frac{d_FF}{d_F - 10F} - \frac{d_FF}{d_F + 10F} = \frac{20d_FF^2}{d_F^2 - 100F^2}. \quad (21)$$

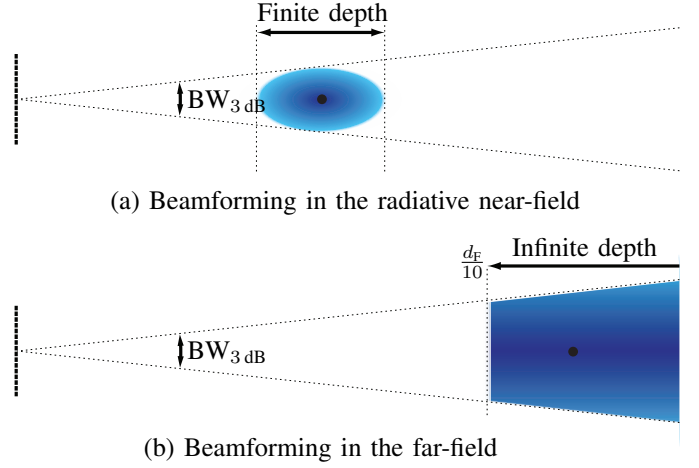


Fig. 4. Beamforming leads to a limited beamwidth regardless of whether the signal is focused on a receiver in the near-field or far-field. However, if the focus point is at a closer distance than $d_F/10$, the beamdepth will be finite. This is not the case when the focus point is beyond $d_F/10$, because then the beam continues until infinity.

This is a unique phenomenon for the beamdepth because for transmission to points $F > d_F/10$ (which includes the far-field), the beam continues all the way to infinity. Interestingly, as $F \rightarrow \infty$ so that the beam is focused at a faraway location, the lower limit $\frac{d_FF}{d_F + 10F}$ converges to $d_F/10$. This is an alternative boundary between the near-field and far-field when considering the beamdepth.

Fig. 4 illustrates the differences between beamforming in the radiative near-field and far-field, which in this case ends at $d_F/10$. The beamdepth is finite in the near-field but semi-infinite in the far-field, meaning that it starts at $d_F/10$ and then continues to infinity. These phenomena resemble the depth of focus of camera lenses, which is finite when focusing on a nearby object (leading to a blurry background) and semi-infinite when the object is far away [34]–[36].

The array processing that underpins near-field beamfocusing has been studied for decades, starting with the works [37]–[39] on microphone arrays. It is the applications within long-range wireless communication networks that are novel [40]–[44] and considered for the 6G era.

B. Near-field spatial multiplexing

The fact that signals transmitted toward receivers in the near-field have a smaller focus area (both in width and depth) means a drastically reduced risk of causing interference between concurrent signal transmissions. This is by itself an enabler of the emerging *massive spatial multiplexing* paradigm, where we are not serving tens of UEs per BS as in 5G but hundreds or a thousand. The classical SE formulas and transmit/receiver signal processing schemes (see [15], [18]) can be utilized, but they will result in substantially higher values thanks to the more favorable propagation conditions obtained in the radiative near-field compared to the far-field.

To exemplify these effects, we begin by considering an uplink single-cell multi-user MIMO setup with K single-

antenna UEs. The channel between the M -antenna BS and UE k is denoted by $\mathbf{h}_k \in \mathbb{C}^M$. The received signal $\mathbf{y} \in \mathbb{C}^M$ is modeled as

$$\mathbf{y} = \sum_{k=1}^K \mathbf{h}_k s_k + \mathbf{n}, \quad (22)$$

where s_k is the data signal transmitted by UE k with power p_k and $\mathbf{n} \sim \mathcal{N}_{\mathbb{C}}(\mathbf{0}, \sigma^2 \mathbf{I}_M)$ is circularly symmetric complex Gaussian noise with variance σ^2 and \mathbf{I}_M is the M -dimensional identity matrix. The BS can apply a receive combining vector $\mathbf{v}_k \in \mathbb{C}^M$ to the received signal in (22) as $\mathbf{v}_k^H \mathbf{y}$ to detect the signal s_k . By treating the co-user interference as noise, the SE for UE k becomes

$$\begin{aligned} & \log_2 \left(1 + \frac{p_k |\mathbf{v}_k^H \mathbf{h}_k|^2}{\sum_{i=1, i \neq k}^K p_i |\mathbf{v}_k^H \mathbf{h}_i|^2 + \sigma^2 \|\mathbf{v}_k\|^2} \right) \\ & \leq \log_2 \left(1 + p_k \mathbf{h}_k^H \left(\sum_{i=1, i \neq k}^K p_i \mathbf{h}_i \mathbf{h}_i^H + \sigma^2 \mathbf{I}_M \right)^{-1} \mathbf{h}_k \right), \end{aligned} \quad (23)$$

$$(24)$$

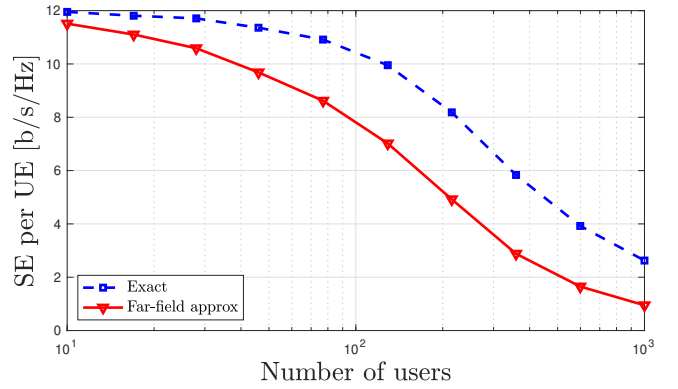
where the upper bound is achieved by using the linear minimum mean-squared error (LMMSE) receive combining vector:

$$\mathbf{v}_k^{\text{LMMSE}} = p_k \left(\sum_{i=1}^K p_i \mathbf{h}_i \mathbf{h}_i^H + \sigma^2 \mathbf{I}_M \right)^{-1} \mathbf{h}_k. \quad (25)$$

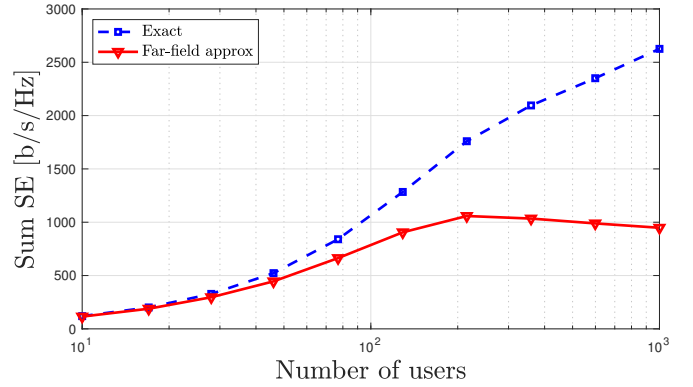
This combining scheme finds the SE-maximizing tradeoff between amplifying the signal power by aligning the receiver to \mathbf{h}_k and rejecting interference by spatial whitening of the received signal using the inverse of $\mathbb{E}\{\mathbf{y}\mathbf{y}^H\}$. The provided formulas are typical ones for uplink multi-user MIMO because the difference lies in the channels.

Fig. 5 shows the uplink SE achieved in a single-cell multi-user MIMO setup operating at 30 GHz. The BS has a half-wavelength-spaced array of 1×0.5 m, containing 5000 antennas in the configuration 100×50 . The single-antenna UEs are uniformly distributed in the horizontal plane in the angular sector $\varphi \in [-\pi/3, +\pi/3]$ and in the distance range 15–500 m. The Fraunhofer distance is $d_F = 250$ m in this setup, so some of the UEs are located in the radiative near-field and others in the far-field. The propagation parameters are otherwise the same as in [45].

The figure contains one curve obtained using an exact line-of-sight channel model and one curve obtained using the conventional far-field approximation, which is mismatched for the UEs actually located in the radiative near-field. Fig. 5(a) shows the average SE per UE, as a function of the number of UEs. A logarithmic scale is used on the horizontal axis since we consider the range 10–1000 UEs. The SE per UE reduces as more UEs are added to the setup, due to the increased interference. The optimal LMMSE receive combining from (25) is utilized. There is a substantial gap between the curves where the exact model consistently provides better results. The far-field approximation basically moves near-field UEs



(a) The average SE per UE.



(b) The average sum SE for all UEs.

Fig. 5. The average uplink SE per UE and sum SE in a setup with 5000 antennas. The exact near-field propagation model leads to much higher values than in an identical setup where a mismatched far-field approximation is utilized.

outwards to the far-field and thereby makes the resulting beam-focus areas wider and deeper, leading to increased interference. This showcases how the ability to utilize the depth domain to distinguish between UE channels makes it easier to deal with interference.

Fig. 5(b) shows the SSE, which is the SE values from Fig. 5(a) multiplied by the respective number of UEs. The massive difference between using the exact and mismatched far-field model becomes evident in this case: With the exact model, the SSE keeps growing with the number of UEs while it saturates at around 200 UEs with the mismatched model. From a network operational perspective, we want to transmit as much data as possible, and the new propagation phenomena observed in the radiative near-field enable higher SEs per UE and efficient spatial multiplexing of many more UEs than if the same MIMO system would operate in the far-field. It is through the massive spatial multiplexing of 1000 UEs that one can reach groundbreaking SSE levels in 6G.

While the depth perception facilitated the spatial multiplexing of many UEs, the tiny beamwidth can enable multiple data streams to be transmitted to a single UE—even in line-of-sight scenarios where we are used to only support a single spatial layer in the far-field. For the sake of argument, suppose the BS

and UE are both equipped with uniform linear arrays (ULAs) with M antennas. The UE has half-wavelength spacing since it is supposed to be a compact device, while the spacing Δ at the BS can be optimized if its ULA is deployed on the facade of a building. When the BS transmits, the half-power beamwidth is given in (17) and is inversely proportional to Δ . If we make it sufficiently small, we can beamform a different signal to each antenna in the receiver array.

We denote the distance between the transmitter and receiver as d , and let the ULAs be deployed perpendicularly to the propagation direction. The BS then sees two adjacent UE antennas with an angular difference of $\sin(\varphi) \approx (\lambda/2)/d$. We can select Δ so that the array gain is zero at the adjacent antenna. The expression in (16) is zero when φ is such that the argument of the sinc-function is 1. By solving for Δ , we obtain

$$\frac{1}{\lambda} N \Delta \sin(\varphi) = 1 \quad \Rightarrow \quad \Delta = \frac{\lambda}{N \sin(\varphi)} \approx \frac{2d}{N}. \quad (26)$$

Hence, the total length of the ULA should be $N\Delta = 2d$, which makes it longer than the propagation distance. More practical deployment scenarios can be achieved by tuning the antenna spacing in the ULAs at both the transmitter and receiver. If the spacings are Δ_t and Δ_r , respectively, the same result can be achieved if [46]–[48]

$$\Delta_t \Delta_r = \frac{\lambda d}{M}. \quad (27)$$

The main message is that one can benefit from increasing the antenna spacing in single-user MIMO systems if that pushes the propagation into the radiative near-field, where the beamwidth can be smaller than the array.

Fig. 6 shows how the SE over a single-user MIMO channel with 16 antennas at the transmitter and receiver. The propagation distance is $d = 50$ m, the frequency is 30 GHz, and the SNR for a single layer is 20 dB. The UE has a ULA with half-wavelength spacing while the antenna spacing at the BS is varied on the horizontal axis. The first point on the curve represents half-wavelength spacing and the figure shows that drastically higher SEs can be achieved by increasing the spacing. The reason is that the MIMO channel matrix transitions from having one non-zero singular value in the beginning to 16 equally large singular values at the peak value at around 8 m. At that point, the BS can transmit a different beam toward each receive antenna, and it will only be “heard” by the designated antennas thanks to the narrow beamwidth. Beyond that antenna spacing, the SE begins to decay again due to sidelobe effects.

The main point of this example is that large SEs can also be achieved for a single UE in the radiative near-field, since we can make the beamwidth so narrow that we can beamform different signals toward different parts of the receiver.

III. SPATIAL DEGREES-OF-FREEDOM

A key reason for increasing the number of antennas in future UM-MIMO systems is to enable more spatial layers, as demonstrated in the last section. We will now take a closer look at the maximum number of spatial layers that can be

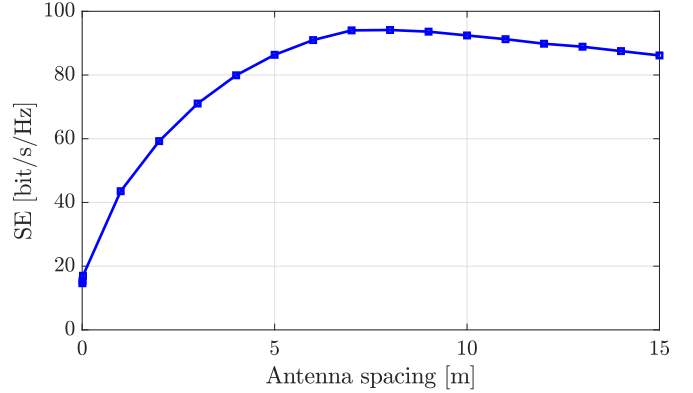


Fig. 6. The SE of single-user MIMO channel with 16 antennas at each side depends on the antenna spacing. The UE has half-wavelength spacing while the BS has a varying spacing, which can be fine-tuned to maximize the SE.

transmitted efficiently over a given channel, which is called the *spatial degrees-of-freedom (DoF)*. As a background to this concept, we begin by reviewing the concept of *spectral DoF*.

A. Spectral degrees-of-freedom

Wireless channels are also known as *waveform channels* because they take a transmitted analog waveform/signal as the input and produce a received analog waveform at the output. Waveforms can be represented either as time-domain signals or as spectra in the frequency domain, and the Fourier transform serves as the bridge between these identical representations. Let us consider a waveform channel that accepts complex-valued continuous-time signals $s(t)$ that is approximately limited to the T -length time interval $[-T/2, T/2]$ and strictly band-limited to the spectral interval $[-B/2, B/2]$. A band-limited signal cannot be entirely time-limited, but according to the Shannon-Nyquist sampling theorem [17], we can approximate $s(t)$ as

$$s(t) \approx \sum_{n=-TB/2}^{TB/2-1} s\left(\frac{n}{B}\right) \text{sinc}(Bt - n), \quad (28)$$

where the approximation error becomes negligible as $TB \rightarrow \infty$. Notably, (28) is a band-limited orthonormal series expansion characterized by the finite set of coefficients

$$\left\{ s\left(\frac{n}{B}\right) : n = -\frac{TB}{2}, \dots, \frac{TB}{2} - 1 \right\}. \quad (29)$$

The cardinality of this set is

$$\eta = TB \quad (30)$$

and is called the *dimension* or DoF of the waveform. Since the signal $s(t)$ is completely determined by B complex-valued equal-spaced samples per second, we can call these the spectral DoF.

In communication systems, we want to design the waveform $s(t)$ to carry data over the channel. A practical system might operate over a real-valued passband channel with bandwidth

B around some carrier frequency f_c ; that is, it accepts real-valued waveforms that are band-limited to the spectral interval $[f_c - B/2, f_c + B/2]$. This channel can be identically represented as a complex baseband waveform channel of the kind described above [49].¹ Hence, the data signal is completely determined by B complex-valued equal-spaced samples per second; these are the spectral DoF available for shaping the communication signal at the transmitter. For example, the transmitter can place information into these samples using a 16-QAM (Quadrature Amplitude Modulation) scheme, which is a complex constellation with 16 different states. Since each sample represents $\log_2(16) = 4$ bits, the transmitter can convey four bits per spectral DoF. If the channel has the bandwidth $B = 10$ MHz, which leads to $10 \cdot 10^6$ DoF per second, the data rate becomes $4 \cdot 10 \cdot 10^6 = 40$ Mbps. More information can be conveyed by increasing the constellation size, but it is essential that the data rate is below the channel capacity, so the receiver can decode the data without error. The SE is the theoretical limit on the number of bits to convey per sample. We previously said that its unit is bit/s/Hz, but it can be equivalently expressed as bit/sample or bit/DoF.

In 5G, the typical bandwidth in the 3.5 GHz band is 100 MHz, and the maximum constellation size is 1024-QAM. This corresponds to $100 \cdot 10^6$ DoF/s and 10 bit/DoF. Consequently, the maximum data rate is 1 Gbps and it is determined by two variables: the spectral DoF (i.e., bandwidth) and SE. It is not realistic to increase the SE beyond 10 bit/DoF in future networks, because one needs an SNR greater than 30 dB to reach that number. When moving beyond that SNR value, the system is typically limited by hardware fidelity rather than noise. We also cannot expect to increase the bandwidth by more than a factor of 10 compared to 5G. Hence, the important question is: *How can we increase the data rate by a factor of 100× or 1000×?* The answer is that we need to create new signal dimensions through MIMO instead.

B. (Massive) MIMO: An information-theoretic perspective

Two types of MIMO communication systems were introduced earlier in this paper: single-user MIMO involves a multi-antenna BS and a multi-antenna UE, while multi-user MIMO involves a multi-antenna BS and multiple UEs. In both cases, the expansion of the signal space into the spatial (antenna) domain creates multiple parallel spatial channels representing the *spatial DoF*. The simultaneous transmission of independent data signals over these spatial channels, as opposed to different time slots and/or frequency subbands, results in a traffic multiplier or *spatial multiplexing gain*, as long as effective measures are taken to mitigate interference between the transmitted signals. Fig. 7 illustrates this as a two-dimensional DoF plane, where the spatial and spectral dimensions are orthogonal. The basic building blocks of a signal is 1 DoF spanning over one Hertz and one antenna dimension. Although we might always have fewer spatial DoF (e.g., hundreds to thousands) than spectral DoF (e.g., millions to billions) in wireless systems, the total number of DoF is

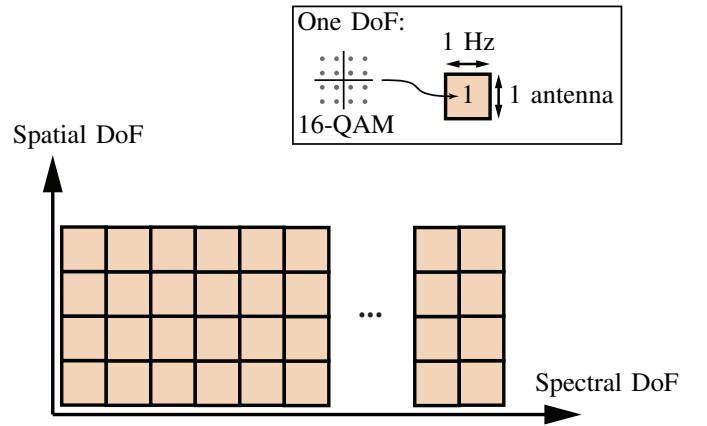


Fig. 7. The basic building blocks of wireless waveforms are called DoF, which is the number of complex-valued coefficients that describes the signal per second. There are both spectral and spatial DoF and their product is the total number of DoF. Each DoF can represent an amount of data that equals the SE of the channel, for example, represented using a 16-QAM constellation.

their product. We will now take a closer look at the spatial DoF delivered by the two MIMO categories.

The capacity of single-user MIMO was first established in [10], [11]. Such a channel with M_t transmit antennas and M_r receive antennas is represented by the MIMO channel matrix $\mathbf{H} \in \mathbb{C}^{M_r \times M_t}$, which can be modeled in various ways depending on propagation environment. The channel capacity is generally determined by the $M_{\min} = \min(M_r, M_t)$ singular values $\mu_1, \dots, \mu_{M_{\min}}$ of \mathbf{H} , such that [49, Sec. 7.1]

$$C = \sum_{i=1}^{M_{\min}} \log_2 \left(1 + \frac{p_i \mu_i^2}{\sigma^2} \right), \quad (31)$$

where the summation is over the different spatial layers and the corresponding transmit powers are $p_1, \dots, p_{M_{\min}}$. For an ideal MIMO channel where all the singular values are equal and each entry of \mathbf{H} has an equal squared magnitude β , the capacity becomes

$$C = M_{\min} \log_2 \left(1 + \text{SNR} \frac{M_r M_t}{M_{\min}^2} \right), \quad (32)$$

where $\text{SNR} = p\beta/\sigma^2$, p is the total transmit power, and σ^2 is the noise variance. This is the kind of MIMO channel that was obtained by optimizing the antenna spacing according to (27). The number of signals that are spatially multiplexed, or equivalently, the number of spatial DoFs, increases with the minimum of M_r and M_t . Hence, the total DoF of a single-user MIMO system with bandwidth B is $\min(M_r, M_t) \times B$ per second. With $M_t = M_r = 64$, there is a potential for a 64-fold increase in capacity compared to a single-antenna system.

The capacity of multi-user MIMO systems was characterized in the early 2000s [50]–[52]. The capacity requires non-linear signal processing at the transmitter or receiver, which is challenging to implement in practical systems. Hence, there

¹Due to Doppler spread and use of other pulses than the sinc-function, the bandwidth can be slightly larger than B .

is a variety of alternative SE expressions for specific linear processing schemes, such as the ones provided in (23). We will consider an uplink expression from [18, Table 3.2] for the case of i.i.d. Rayleigh fading channels, K single-antenna UEs and M BS antennas. If all UEs have the same average SNR, the SSE becomes

$$K \log_2 \left(1 + \frac{M \text{SNR}}{K \text{SNR} + 1} \right). \quad (33)$$

The multiplicative factor K preceding the logarithm indicates the availability of K spatial DoFs. What is more intriguing is that by increasing both M and K simultaneously (with some fixed ratio M/K that is preferably large), we can maintain a nearly constant SE per UE while providing a K -fold increase in SSE. If we could somehow allow the UEs to collaborate in a multi-user MIMO system, we obtain a single-user MIMO setup whose capacity must be equal or higher. This indicates that the spatial DoFs cannot surpass $\min(M, K)$ in a multi-user MIMO system.

The basic principles of MIMO outlined above were established for specific idealistic channel models and perfect CSI. However, the spatial multiplexing capability of single-user MIMO is traditionally hampered by having a low-rank channel matrix, while imperfect CSI acquisition is a main limiting factor for multi-user MIMO. The Massive MIMO concept was introduced in [13] to address these challenges. Firstly, it relies on deploying a significantly larger number of BS antennas than spatially multiplexed devices (e.g., $M/K \geq 8$). In addition to the reason described above, there are two further practical advantages: it reduces interference between UEs and provides the so-called *channel hardening*, which ensures minimal SNR fluctuations after the precoding/combining has been applied. Secondly, the protocols were designed for time-division duplexing (TDD) to enable CSI acquisition for arbitrarily many BS antennas through uplink pilot transmission [15], [18]. We will return to the channel estimation challenge in Section IV.

Current 5G BSs are built using the Massive MIMO principle with digital planar arrays of $M = 32$ or $M = 64$ antennas. When increasing these numbers in future networks to hundreds or thousands, it would be desirable to pack the antennas more closely. However, this will make the channel coefficients more similar (e.g., statistically correlated) and we cannot increase the capacity indefinitely by doing so. This is analogous to the waveform channel where, given the bandwidth constraint B and the transmission interval T , increasing the number of time samples will also not increase the spectral DoF indefinitely. This is called temporal oversampling. The available spectral DoF in the given time interval is always limited to TB . A fundamental question arises: *given an area limitation for the antenna array, what is the intrinsic number of DoF available in the channel?* To answer this question, it is necessary to extend the Shannon-Nyquist sampling theorem to spatially bandlimited fields.

C. Shannon-Nyquist sampling theorem for electric fields

While a single receive antenna takes samples of an impinging waveform at different times, an antenna array can take

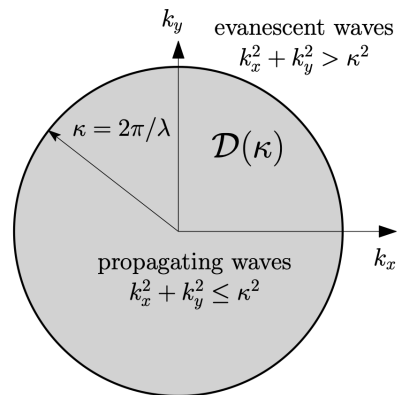


Fig. 8. Wavenumber support of the electric field $s(x, y, z)$.

many spatially separated samples of the waveform at the same time. We can apply the classical Shannon-Nyquist sampling theorem also in this situation, with the only difference being that the samples are collected over a limited spatial interval rather than a limited time interval. When transforming a spatial signal into the frequency domain, we obtain the spatial frequencies, also known as *wavenumbers*. Before we can apply the sampling theorem to this situation, we need an accurate description of electric fields in the wavenumber domain, which is tied to its physical nature and geometry of the antenna array, as it imposes different boundary conditions on the set of Maxwell's equations. In both the EM and communications literature, various methodologies for the treatment of deterministic and stochastic spatial fields have been explored. A non-exhaustive list of relevant publications in this area is [53]–[63].

To understand the basic principles, we consider a scenario where EM waves propagate through a homogeneous, isotropic, source-free, and scattered infinite medium. Monochromatic waves with no polarization then behave as acoustic waves and the electric field $\{s(x, y, z) : (x, y, z) \in \mathbb{R}^3\}$ satisfies the scalar Helmholtz equation in the frequency domain [64, Eq. (1.2.17)]

$$(\nabla^2 + \kappa^2) s(x, y, z) = 0, \quad (34)$$

where $\kappa = \frac{2\pi}{\lambda}$ is the angular wavenumber of the considered signal (i.e., the angular variation in radians per unit of length). The solution to (34) takes the form

$$s(x, y, z) = \alpha e^{j(k_x x + k_y y + k_z z)} \quad (35)$$

where $\alpha \in \mathbb{C}$ is an unknown complex scaling factor and $(k_x, k_y, k_z) \in \mathbb{R}^3$ are three wavenumber coefficients that characterize the solution. These represent the wavenumbers observed in the x , y , and z dimensions, respectively. Plugging (35) into (34) yields the condition

$$k_x^2 + k_y^2 + k_z^2 = \kappa^2. \quad (36)$$

This implies that the wavenumber κ of the original wave is divided between the three dimensions. We can observe the value ranges $k_x, k_y, k_z \in [-\kappa, \kappa]$, but if we increase the magnitude of one wavenumber, the others must reduce accordingly. In fact, the constraint in (36) allows us to eliminate one of

the three coefficients. Specifically, we consider the half-space where k_z is positive such that

$$k_z = \sqrt{\kappa^2 - k_x^2 - k_y^2}. \quad (37)$$

It follows that (k_x, k_y) must have compact support given by

$$\mathcal{D}(\kappa) = \left\{ (k_x, k_y) \in \mathbb{R}^2 : k_x^2 + k_y^2 \leq \kappa^2 \right\}, \quad (38)$$

which is a disk of radius κ centered on the origin, as illustrated in Fig. 8. From (35) and (37), we thus have that

$$s(x, y, z) = \alpha e^{j(k_x x + k_y y + \sqrt{\kappa^2 - k_x^2 - k_y^2} z)}, \quad (39)$$

which is the equation of an incident *plane-wave* impinging on the spatial point (x, y, z) . Note that, by imposing the condition $(k_x, k_y, k_z) \in \mathbb{R}^3$, we exclude the *evanescent waves* from the analysis and consider only *propagating waves*. This is consistent with our previous assumption of studying the radiative near-field and far-field. If no directionality is enforced on the EM waves (corresponding to an isotropic scattering environment), then the support of $s(x, y, z)$ is limited to $|\mathcal{D}(\kappa)| = \pi\kappa^2$, revealing the spatially band-limited nature of EM fields. Scattering mechanisms act as a filter limiting further the field to a smaller support [62]. Compared to the classical definition of a band-limited signal in the frequency domain, here the notion of band-limited support applies in the wavenumber domain.

The above findings are crucial for calculating the spatial DoF of electric fields through the generalization of the sampling theorem. Assume for example that $s(x, y, z)$ is observed over a one-dimensional (1D) line segment of length L_x along the x -axis. Then, the spatial DoF that characterize all possible fields that can be observed are given by [59]–[62]

$$\eta_{1D} = \frac{2}{\lambda} L_x = \frac{\kappa}{\pi} L_x. \quad (40)$$

The first expression is the product between the length L_x of the spatial interval where the field is observed, while $\frac{2}{\lambda}$ is the length of the interval $[-1/\lambda, 1/\lambda]$ of spatial frequencies (i.e., non-angular wavenumbers) that the field might contain. This is the spatial-wavenumber counterpart to the product TB between time duration T and bandwidth B in (30). The second expression in (40) expresses the same relation using the angular wavenumber κ .

Suppose that the electric field $s(x, y, z)$ is instead observed over a two-dimensional (2D) rectangle with side lengths L_x and L_y . In this case, the spatial DoF are given by [59]–[62]

$$\eta_{2D} = \frac{\pi}{\lambda^2} L_x L_y, \quad (41)$$

which are proportional to the surface area $L_x L_y$ normalized by the squared wavelength. Hence, each portion of the array aperture with area λ^2 can observe π DoF from the impinging electric field.

One may expect that the spatial DoFs of a two-dimensional array would be the product of the DoFs that can be observed horizontally and vertically, using the one-dimensional formula in (40). That computation leads to

$$\tilde{\eta} = \left(\frac{2}{\lambda} L_x \right) \left(\frac{2}{\lambda} L_y \right) = \frac{4}{\lambda^2} L_x L_y, \quad (42)$$

which is different from (41). Specifically $\eta_{2D}/\tilde{\eta} = \pi/4 \approx 0.79$, so the correct spatial DoFs is smaller. The difference is exactly the ratio between the areas of the disk $\mathcal{D}(\kappa)$ and the square circumscribing it, regardless of the dimensions of the rectangular aperture. The intuition is that the wavenumbers observed horizontally and vertically in a planar array are correlated. For example, a wave that arrives from a large elevation angle (i.e., near the y -axis) can give rise to rather small horizontal variations, similar to how there are smaller distances when moving around the Earth near the North Pole compared to the equator.

D. Implications for MIMO Systems

The signal transmission in wireless communications generates electric fields at the transmitter and samples them at the receiver. The spatial DoF determine how many coefficients are required to characterize these electric fields, not in general but from the perspective of a particular antenna array; if two different electric fields look the same to the array, then we cannot use their difference to carry any additional data. If we are given a particular deployment area of $L_x \times L_y$ meters to deploy an array, *how should we deploy the antennas to obtain all the available spatial DoF?*

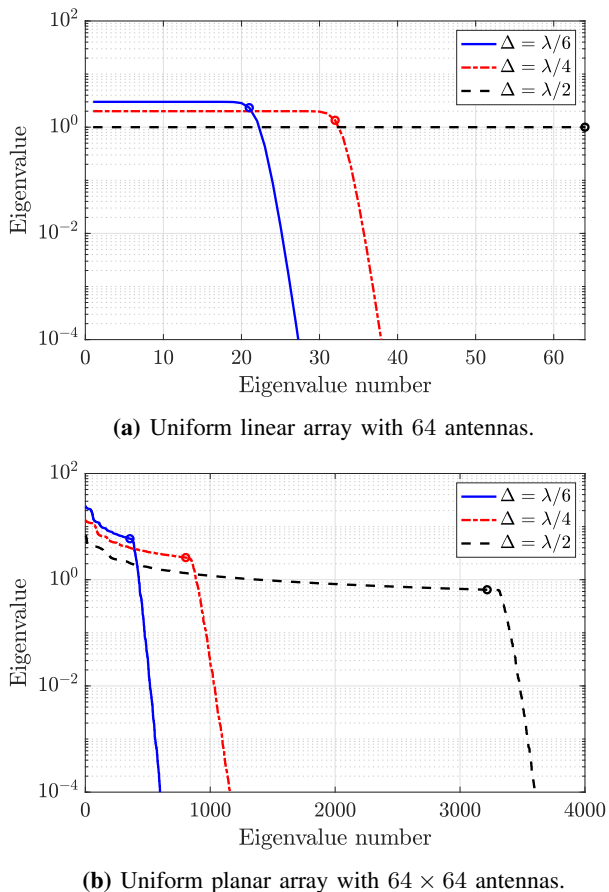
The common practice in array design has been to use uniform planar arrays (UPAs) with the spacing $\Delta = \lambda/2$ both horizontally and vertically. We can fit

$$M = \frac{L_x}{\Delta} \frac{L_y}{\Delta} = \frac{4}{\lambda^2} L_x L_y \quad (43)$$

antennas into this area. This value matches with (42) and is, thus, larger than the available spatial DoF that are given by (41). Consequently, the conventional approach to array design is sufficient to capture all the available spatial DoF.

Fig. 9 illustrates what would happen if we reduce the antenna spacing to $\Delta < \lambda/2$. We consider an isotropic scattering environment where plane waves can arrive at the M receive antennas from any direction with equal probability. In this situation, one can compute the spatial correlation matrix $\mathbf{R} = \mathbb{E}\{\mathbf{h}\mathbf{h}^H\}$ of the channel vector $\mathbf{h} \in \mathbb{C}^M$ and study its eigenvalues. The number of large eigenvalues represents the number of spatial DoF that the array observes, and this is the maximum value since an isotropic environment excites all possible channel dimensions. Fig. 9(a) considers a uniform linear array (ULA) with $M = 64$ antennas and varying antenna spacings. In the case of $\Delta = \lambda/2$, the array covers a spatial interval of $M\lambda/2$ meters, and the spatial DoF in (40) then becomes M , which results in all eigenvalues being equally large. However, if we decrease the antenna spacing to $\lambda/4$ or $\lambda/6$, the number of spatial DoF respectively reduces to $M/2$ and $M/3$. These numbers are illustrated with circles in the figure and clearly predict the number of large eigenvalues; that is, the number of channel dimensions. Fig. 9(b) considers the case of a UPA with 64×64 antennas. We see the same general trends as in the case of a ULA, with the main difference that even the case with $\lambda/2$ spacing leads to only 79% large eigenvalues.

The main conclusion from this section is that we should continue using arrays with $\lambda/2$ -spacing in future UM-MIMO



(a) Uniform linear array with 64 antennas.

(b) Uniform planar array with 64×64 antennas.

Fig. 9. The normalized eigenvalues of the spatial correlation matrix \mathbf{R} in decreasing order. We consider an isotropic scattering environment with different array geometries and antenna spacings. The numbers of large eigenvalues are determined by the spatial DoF captured by the particular array type. The theoretical values from (40) and (41) are shown using circles.

systems. We have seen earlier in the paper (e.g., in Fig. 6) that one can possibly benefit from increasing the antenna spacing beyond that limit, to observe near-field propagation effects that improve the channel conditions. However, one should not reduce the antenna spacing in the hope of increasing the spatial DoF because that is not possible—it corresponds to spatial oversampling. If the array dimensions are limited, there are still some good reasons for filling it with more antennas than what is needed to capture all the available spatial DoF. Small antennas have a more isotropic-like radiation pattern, which improves the array’s ability to transmit and receive signals in any direction. But if we shrink the antenna sizes, we should compensate by adding more antennas into the aperture to keep the total antenna area fixed.

IV. CHANNEL ESTIMATION FOR LARGE MIMO ARRAYS

CSI is necessary to make efficient use of a large number of antennas, so that the transmitted signals can be beamformed to the intended location in the downlink, and the received signal can be combined coherently in the uplink. The basic way of acquiring CSI is to transmit a predefined pilot sequence

and estimate the channel coefficients at the receiver side. The length of the pilot sequence must equal the number of transmit antennas, while arbitrarily many receive antennas can collect channel estimates simultaneously. New communication systems are nowadays built using TDD spectrum, where the same band is used in both uplink and downlink. Hence, we have the liberty to choose in which direction to transmit pilots. Since the 5G Massive MIMO technology builds on serving tens of UEs with a large number of BS antennas, the pilot sequences (called sounding reference symbols in 5G) are transmitted in the uplink where the pilot sequence length equals the number of UE antennas [13].

The same procedure can be used in the next generation of MIMO technology, but if we increase the number of antennas per UE and the number of spatially multiplexed UEs, the pilot sequence length will increase accordingly. In this section, we will outline how the pilot sequence can be reduced depending on what prior information exists regarding the channel properties.

For brevity in presentation, we focus on the channel estimation from an M -antenna UE to one of the many antennas at the BS. The approach described in this section can be applied separately for each BS antenna. The considered channel vector is denoted by $\mathbf{h} \in \mathbb{C}^M$. The channel is made of a superposition of multipath components. In the far-field, it can be described as a discrete summation of plane waves arriving from different directions. In the radiative near-field, the channel can be expressed as a continuous summation of plane waves [55], which is also an accurate representation of spherical waves. Hence, we may write the channel vector as

$$\mathbf{h} = \iint_{-\pi/2}^{\pi/2} g(\varphi, \theta) \mathbf{s}(\varphi, \theta) d\theta d\varphi \quad (44)$$

where $\mathbf{s}(\varphi, \theta)$ is the far-field array response vector for the azimuth angle φ and elevation angle θ , while the *angular spreading function* $g(\varphi, \theta)$ specifies the gain and phase-shift from each direction. The integration limits for the azimuth angle are limited to $\varphi \in [-\frac{\pi}{2}, \frac{\pi}{2}]$ because waves can only arrive from the front of the array.

The channel realization is determined by the angular spreading function. Since small-scale UE mobility is hard to model accurately and the interaction with the surrounding multipath environment is complex, it is customary to model small-scale variations stochastically. In this section, we consider the block fading model, where the channel vector \mathbf{h} is constant within one block of time-frequency resources and takes independent realization across blocks from a stationary stochastic distribution. We will model $g(\varphi, \theta)$ as a spatially uncorrelated circularly symmetric Gaussian stochastic process with cross-correlation

$$\mathbb{E}\{g(\varphi, \theta)g^*(\varphi', \theta')\} = \beta f(\varphi, \theta)\delta(\varphi - \varphi')\delta(\theta - \theta'), \quad (45)$$

where β denotes the average channel gain and $f(\varphi, \theta)$ is the normalized *spatial scattering function* [55]. This is a probability density function (PDF) that provides a statistical representation of the multipath environment in terms of how likely it is for signals to arrive from different directions. As with any PDF, it holds that $\iint f(\varphi, \theta)d\theta d\varphi = 1$. Based on

these assumptions, we obtain the classical correlated Rayleigh fading channel distribution

$$\mathbf{h} \sim \mathcal{N}_{\mathbb{C}}(\mathbf{0}, \mathbf{R}), \quad (46)$$

but the key reason for the aforementioned assumptions is that we can give the spatial correlation matrix a geometric model:

$$\mathbf{R} = \mathbb{E}\{\mathbf{h}\mathbf{h}^H\} = \beta \iint_{-\pi/2}^{\pi/2} f(\varphi, \theta) \mathbf{s}(\varphi, \theta) \mathbf{s}^H(\varphi, \theta) d\theta d\varphi \quad (47)$$

which follows from the property in (45). The average gain of the channel is $\mathbb{E}\{\|\mathbf{h}\|^2\} = \text{tr}(\mathbf{R}) = M\beta$.

We will consider different ways of estimating the realizations of \mathbf{h} depending on which parts of the statistical characterization are known. In all these cases, the estimation is based on transmitting a predefined pilot sequence of some length $\tau_p \leq M$. It is desirable to make this sequence as short as possible to not spend unnecessarily many signal resources (i.e., DoF) on channel estimation, and we will later show that different estimators allow for different pilot lengths. We let $\Phi \in \mathbb{C}^{\tau_p \times M}$ denote the pilot sequence matrix, where the (n, m) th entry represents the pilot symbol transmitted from antenna m of the UE at time instance n in the τ_p -length sequence. We assume that the average pilot power is normalized to one, in the sense that

$$\text{tr}(\Phi^H \Phi) = \tau_p. \quad (48)$$

If we collect all the τ_p received symbols at the BS antenna in a vector $\mathbf{y} \in \mathbb{C}^{\tau_p}$, we can express it as

$$\mathbf{y} = \sqrt{p}\Phi\mathbf{h} + \mathbf{n}, \quad (49)$$

where $p > 0$ is the pilot power and $\mathbf{n} \sim \mathcal{N}_{\mathbb{C}}(\mathbf{0}, \sigma^2 \mathbf{I}_{\tau_p})$ is the independent noise.

A. Least-squares estimation

The simplest channel estimation method is least-squares (LS) estimation which does not require any statistical information regarding the channel vector. This estimator finds the solution to the *least squares* problem

$$\underset{\hat{\mathbf{h}} \in \mathbb{C}^M}{\text{minimize}} \quad \left\| \mathbf{y} - \sqrt{p}\Phi\hat{\mathbf{h}} \right\|^2 \quad (50)$$

between the actual received signal \mathbf{y} and the potential received signal $\sqrt{p}\Phi\hat{\mathbf{h}}$ in the absence of noise. If the pilot length is $\tau_p = M$ so that the pilot matrix Φ can be selected to be invertible, the solution to (50) is obtained as

$$\hat{\mathbf{h}}_{\text{LS}} = \frac{\Phi^{-1}\mathbf{y}}{\sqrt{p}} = \mathbf{h} + \underbrace{\frac{\Phi^{-1}\mathbf{n}}{\sqrt{p}}}_{=\tilde{\mathbf{h}}_{\text{LS}}}, \quad (51)$$

where $\tilde{\mathbf{h}}_{\text{LS}} \in \mathbb{C}^M$ is the channel estimation error. This choice makes the term inside the square in (50) zero. The channel estimation quality can be quantified by the average power of

the channel estimation error, which is called the mean-squared error (MSE). The MSE of the LS estimator is computed as

$$\begin{aligned} \text{MSE}_{\text{LS}} &= \mathbb{E} \left\{ \left\| \mathbf{h} - \hat{\mathbf{h}}_{\text{LS}} \right\|^2 \right\} = \mathbb{E} \left\{ \left\| \tilde{\mathbf{h}}_{\text{LS}} \right\|^2 \right\} \\ &= \frac{\sigma^2}{p} \text{tr} \left((\Phi^H \Phi)^{-1} \right) \end{aligned} \quad (52)$$

by utilizing the fact that $\mathbb{E}\{\mathbf{n}\mathbf{n}^H\} = \sigma^2 \mathbf{I}_M$. This expression reveals that the channel estimation quality depends on the selection of the pilot matrix. It can be shown that selecting Φ as any unitary matrix minimizes the MSE, under the average power constraint in (48). The resulting minimum MSE is

$$\text{MSE}_{\text{LS}}^* = \frac{M\sigma^2}{p} \quad (53)$$

which is proportional to the number of UE antennas (i.e., the number of unknowns) and a linearly decreasing function of the pilot SNR, $\frac{p}{\sigma^2}$. The key observation is that the optimal pilot matrix is obtained by treating each dimension of the M -dimensional vector space equally by allocating the same amount of power to all of them.

The main drawback of the LS estimator is that the pilot length must equal the number of UE antennas, otherwise, we cannot invert the pilot matrix. To reduce the pilot length and obtain better-quality channel estimates, we need to exploit more about the structure of the channel, as elaborated in the following part.

B. Minimum mean-squared error estimation

The considered channel \mathbf{h} follows the correlated Rayleigh fading model in (46) with the spatial correlation matrix \mathbf{R} defined in (47). If this correlation matrix is completely known at the BS, the minimum MSE (MMSE) estimate can be computed as [65]

$$\hat{\mathbf{h}}_{\text{MMSE}} = \sqrt{p}\mathbf{R}\Phi^H \left(p\Phi\mathbf{R}\Phi^H + \sigma^2\mathbf{I}_M \right)^{-1} \mathbf{y}. \quad (54)$$

As the name suggests, the MMSE estimator minimizes the MSE among all conceivable estimators that have access to the statistical characterization. The true channel \mathbf{h} can be decomposed as $\mathbf{h} = \hat{\mathbf{h}}_{\text{MMSE}} + \tilde{\mathbf{h}}_{\text{MMSE}}$, where the estimation error $\tilde{\mathbf{h}}_{\text{MMSE}}$ is independent of the estimate $\hat{\mathbf{h}}_{\text{MMSE}}$. Consequently, the MSE can be computed as

$$\begin{aligned} \text{MSE}_{\text{MMSE}} &= \text{tr} \left(\mathbb{E} \left\{ \tilde{\mathbf{h}}_{\text{MMSE}} \tilde{\mathbf{h}}_{\text{MMSE}}^H \right\} \right) \\ &= \text{tr} \left(\mathbb{E} \left\{ \mathbf{h}\mathbf{h}^H \right\} \right) - \text{tr} \left(\mathbb{E} \left\{ \hat{\mathbf{h}}_{\text{MMSE}} \hat{\mathbf{h}}_{\text{MMSE}}^H \right\} \right) \\ &= M\beta - \text{tr} \left(p\mathbf{R}\Phi^H \left(p\Phi\mathbf{R}\Phi^H + \sigma^2\mathbf{I}_M \right)^{-1} \Phi\mathbf{R} \right). \end{aligned} \quad (55)$$

This MSE depends on the pilot matrix Φ and, thus, it can be minimized by properly designing the pilot matrix. We let $\mathbf{R} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^H$ denote the eigendecomposition of the spatial correlation matrix \mathbf{R} , where the unitary matrix $\mathbf{U} \in \mathbb{C}^{M \times M}$ contains the eigenvectors as its columns, and the corresponding eigenvalues are located in decreasing order along the diagonal of the matrix $\mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_M)$. It can be shown that

the MSE is minimized by selecting the $\tau_p \times M$ pilot matrix as $\Phi = \mathbf{D}\mathbf{U}^H$ [66], where $\mathbf{D} \in \mathbb{C}^{\tau_p \times M}$ is a rectangular diagonal matrix where the diagonal entries are

$$d_m = \sqrt{\max\left(0, \mu - \frac{\sigma^2}{p\lambda_m}\right)}, \quad m = 1, \dots, \tau_p, \quad (56)$$

where $\mu > 0$ is selected such that

$$\sum_{m=1}^{\tau_p} \max\left(0, \mu - \frac{\sigma^2}{p\lambda_m}\right) = \tau_p. \quad (57)$$

This pilot matrix is matched to the eigendecomposition of the spatial correlation matrix since \mathbf{U} is utilized, while d_m^2 determines how much power is allocated to estimating the channel components along the m th eigenvector. The power allocation in (56) has a water-filling structure, where more power is allocated to the channel directions with larger eigenvalues.

If we substitute $\Phi = \mathbf{D}\mathbf{U}^H$ into (54), we can simplify the estimator as

$$\hat{\mathbf{h}}_{\text{MMSE}} = \mathbf{U}\mathbf{A}\mathbf{y} \quad (58)$$

where $\mathbf{A} = \sqrt{p}\mathbf{A}\mathbf{D}^T (p\mathbf{D}\mathbf{A}\mathbf{D}^T + \sigma^2\mathbf{I}_{\tau_p})^{-1}$ is a diagonal matrix. This MMSE estimator carries out two operations. First, it computes the MMSE estimates of the channel components in the τ_p strongest eigendirections as $\mathbf{A}\mathbf{y}$. Second, it brings this estimate back to the original channel space by multiplying it with the eigenvector matrix \mathbf{U} .

The core difference between the MMSE and LS estimator is that the former knows the statistical strength of the channel in different in each eigendirection. It can therefore fine-tune the estimator and allocate more pilot power to stronger eigendirections, thereby reducing the MSE. While the LS estimator necessitates $\tau_p = M$, the MMSE estimator can be applied with any τ_p and will then only transmit pilots along the τ_p strongest eigendirections. To highlight the practical importance of this, we will consider a scenario where the rank of \mathbf{R} is strictly smaller than M , as previously illustrated in Fig. 9.

In Fig. 10, we plot the normalized mean square error (NMSE), obtained by dividing the MSE by $\mathbb{E}\{\mathbf{h}\mathbf{h}^H\} = \text{tr}(\mathbf{R})$ for the LS and MMSE estimators in different propagation environments. We consider an 8×8 UPA with the antenna spacing $\Delta = \lambda/4$. Two different propagation environments are considered: isotropic and clustered scattering. The isotropic environment assumes that the multipath components are equally strong in all directions (as in Fig. 9). The clustered environment is generated using a model from [67] where there are three scattering clusters located in the azimuth directions $0, \pi/4$, and $-\pi/4$ and each having a 10° angular standard deviation. The effective signal-to-noise ratio (SNR) is $p\text{tr}(\mathbf{R})/(M\sigma^2) = 10$ dB.

The figure shows the NMSE as a function of the pilot length τ_p . When $\tau_p < M$, the LS estimator is applied with the pseudoinverse of the pilot matrix instead of the true inverse. This leads to poor estimation performance, which is why we ruled out this situation previously. The MMSE estimator outperforms the LS estimator since it exploits the spatial correlation matrix, but it particularly enables good estimation quality at significantly smaller pilot lengths. The

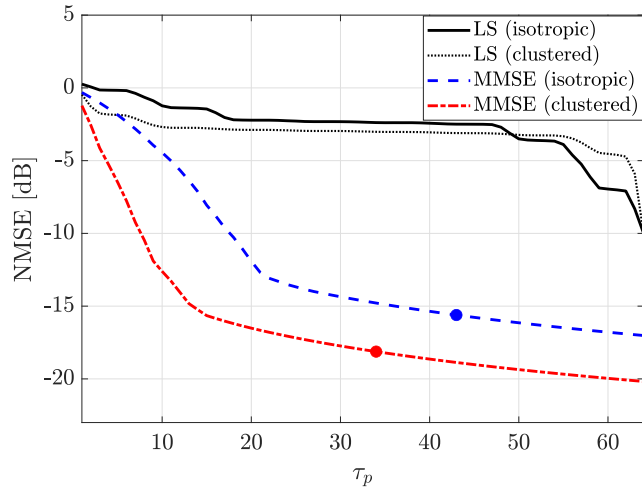


Fig. 10. The NMSE versus the pilot length τ_p for the LS and MMSE estimators in isotropic and clustered scattering environments. The dots show the NMSE performance when τ_p is equal to the rank of the corresponding spatial correlation matrix.

dots on the plot represent the MMSE estimator's performance when τ_p equals the rank of the spatial correlation matrices. In the case of isotropic scattering, the value matches the spatial DoF. For clustered scattering, the rank is smaller since the scattering environment only excites a subset of the possible channel dimensions. If we increase the pilot length beyond the dot-marked number, the pilot matrix will not explore any new channel dimensions but only benefit from increased pilot energy. The MMSE estimator provides lower MSEs when the propagation environment has a structure, such as clustered scattering because it can then focus the pilot power into a few important eigendirections.

C. Reduced-subspace least-square estimation

In the preceding section, we demonstrated how the MMSE estimator effectively utilizes the spatial correlation matrix to minimize the MSE. This $M \times M$ matrix can be estimated in practice by collecting many channel realizations and forming a sample covariance matrix [15], but much more than M observations are needed to obtain an accurate estimate. Hence, it is challenging to acquire the per-UE spatial correlation matrix information in practice, particularly in scenarios with a larger antenna number, rapid UE mobility that changes the statistics, or during short data packet transmissions [67].

In such situations, an alternative approach is to only exploit the spatial correlations induced by the array geometry and general characteristics of the propagation environment. For instance, it may be known that multipath components can only be observed within certain angular intervals. This implies that any plausible channel vector lies in a lower-dimensional subspace of \mathbb{C}^M . If we know the basis vectors of this subspace, the channel estimation can be performed exclusively within this subspace. This approach is called *reduced-subspace least squares (RS-LS)* estimation [67].

We collect the orthonormal basis vectors of the subspace as columns of the matrix $\bar{\mathbf{U}} \in \mathbb{C}^{M \times \bar{r}}$, where $\bar{r} \leq \tau_p < M$ is the dimension of the subspace. Hence, any UE channel can be expressed as $\mathbf{h} = \bar{\mathbf{U}}\mathbf{v}$ for some $\mathbf{v} \in \mathbb{C}^{\bar{r}}$. The RS-LS estimator builds on estimating \mathbf{v} and consists of two steps. First, we compute the LS estimate of \mathbf{v} in the subspace spanned by the columns of $\bar{\mathbf{U}}$ as

$$\hat{\mathbf{v}}_{\text{LS}} = \frac{1}{\sqrt{p}} \left(\bar{\mathbf{U}}^H \Phi^H \Phi \bar{\mathbf{U}} \right)^{-1} \bar{\mathbf{U}}^H \Phi^H \mathbf{y}, \quad (59)$$

where we assumed $\tau_p \geq \bar{r}$ so that the inverse is well-defined. The matrix $\left(\bar{\mathbf{U}}^H \Phi^H \Phi \bar{\mathbf{U}} \right)^{-1} \bar{\mathbf{U}}^H \Phi^H$ in (59) gives the orthogonal projection of the channel onto the considered subspace. Next, we return this estimate to the original M -dimensional space by multiplying $\hat{\mathbf{v}}_{\text{LS}}$ by $\bar{\mathbf{U}}$:

$$\begin{aligned} \hat{\mathbf{h}}_{\text{RS-LS}} &= \bar{\mathbf{U}} \hat{\mathbf{v}}_{\text{LS}} = \frac{1}{\sqrt{p}} \bar{\mathbf{U}} \left(\bar{\mathbf{U}}^H \Phi^H \Phi \bar{\mathbf{U}} \right)^{-1} \bar{\mathbf{U}}^H \Phi^H \mathbf{y} \\ &= \underbrace{\bar{\mathbf{U}} \mathbf{v}}_{=\mathbf{h}} + \frac{1}{\sqrt{p}} \bar{\mathbf{U}} \left(\bar{\mathbf{U}}^H \Phi^H \Phi \bar{\mathbf{U}} \right)^{-1} \bar{\mathbf{U}}^H \Phi^H \mathbf{n}. \end{aligned} \quad (60)$$

The RS-LS estimator effectively eliminates noise from all unused channel dimensions when $\bar{r} < M$, and the required pilot length is reduced to $\tau_p \geq \bar{r}$ (compared to $\tau_p \geq M$ for the original LS estimator). The remaining question is how to determine the basis vectors for the reduced subspace. One option is to collect many different spatial correlation matrices over time and then compute the union of their span [67, Lem. 3]. Alternatively, one can consider the worst-case scenario of isotropic scattering, where all conceivable channel dimensions might exist in the channel vector. Fig. 9 previously showed that such channels have a low-rank behavior when using UPAs or when the antenna spacing is less than $\lambda/2$. This property is not UE-channel-specific but rather array-dependent, thereby enabling the removal of noise from unused directions.

The MSE of the RS-LS estimator in (60) is

$$\begin{aligned} \text{MSE}_{\text{RS-LS}} &= \frac{\sigma^2}{p} \text{tr} \left(\bar{\mathbf{U}} \left(\bar{\mathbf{U}}^H \Phi^H \Phi \bar{\mathbf{U}} \right)^{-1} \bar{\mathbf{U}}^H \right) \\ &= \frac{\sigma^2}{p} \text{tr} \left(\left(\bar{\mathbf{U}}^H \Phi^H \Phi \bar{\mathbf{U}} \right)^{-1} \right) \end{aligned} \quad (61)$$

and depends on the pilot matrix Φ . The MSE is minimized by selecting the pilot matrix as

$$\Phi^* = \sqrt{\frac{\tau_p}{\bar{r}}} \mathbf{S} \bar{\mathbf{U}}^H \quad (62)$$

where $\mathbf{S} \in \mathbb{C}^{\tau_p \times \bar{r}}$ is an arbitrary matrix with orthonormal columns [68]. Substituting the optimal pilot matrix into the MSE expression in (61), we obtain

$$\text{MSE}_{\text{RS-LS}}^* = \frac{\bar{r}^2 \sigma^2}{\tau_p p}, \quad (63)$$

which is $(M/\bar{r})^2$ times smaller than the MSE achieved by the LS estimator in (53) when $\tau_p = M$.

Fig. 11 shows the NMSE as a function of the antenna spacing Δ for the 8×8 UPA considered previously for the same channel with clustered scattering. The LS estimator

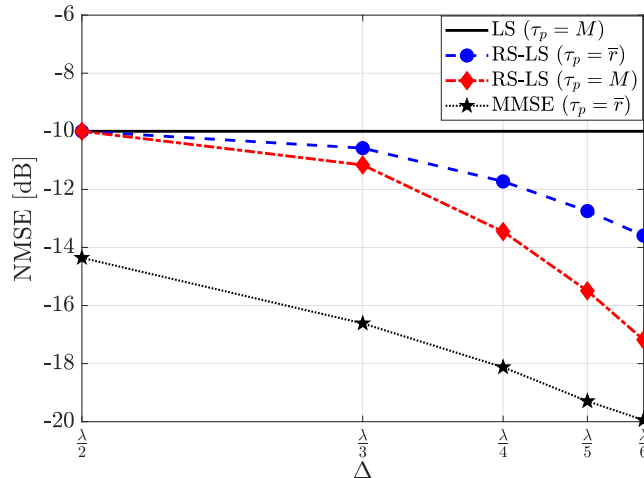


Fig. 11. The NMSE versus the antenna spacing Δ for the LS, RS-LS, and MMSE estimators in a clustered scattering environment.

serves as a reference with $\tau_p = M$ and is unaffected by the antenna spacing, as evident from the MSE expression in (53). By contrast, the other estimators exhibit smaller NMSE values as the antenna spacing decreases, because of the higher spatial correlation that these estimators exploit to varying degrees.

The isotropic spatial correlation matrix is utilized to construct the reduced subspace in the RS-LS estimator. Two RS-LS results are presented: one with $\tau_p = M$, and the other with a pilot length equal to the dimension of the reduced subspace. Since subspace size \bar{r} decreases when Δ becomes smaller, the pilot length used by the RS-LS estimator and MMSE estimator with $\tau_p = \bar{r}$ also decreases. Despite this, the NMSE decreases as Δ reduces thanks to the increased spatial correlation. We note that the NMSE with the RS-LS is better with $\tau_p = M$ than with $\tau_p = \bar{r}$, but the latter might still be preferable in practice since fewer pilot resources are required.

In summary, the RS-LS estimator is a meaningful alternative to the conventional LS estimator, as both methods do not require UE-specific statistical information.

D. Compressed-sensing-based channel estimation

There is a middle ground between the MMSE estimator, which requires the complete UE-specific channel statistics, and the RS-LS estimator which only utilizes channel statistics at the UE population level. If the BS knows that the UE channel features multipath propagation caused by only a small number of scattering clusters, this information can be utilized by the estimator. The goal is then to jointly sense the locations of the clusters and estimate their related parameters. In such cases, channel estimation based on compressed sensing methods can be effective in achieving a good estimation quality with a relatively small pilot overhead [69].

The first step in compressed-sensing-based channel estimation is to create a dictionary of vectors representing the channel from scattering clusters at different plausible locations. The goal is then to identify a linear combination of a small number of these dictionary vectors that results in a channel vector that

resembles the one observed during the pilot transmission. An example of this is shown in Fig. 12, where we consider a channel comprising $L = 3$ paths from clusters located in the far-field and a UPA at the UE. The UE antennas are deployed as 8×8 UPA with $\lambda/4$ antenna spacing. Each path is represented by an array response vector, scaled by a complex channel gain. Hence, the dictionary contains many such array response vectors with uniform sampling of the plausible azimuth and elevation angles. The azimuth angular grid $\Psi = \sin(\varphi) \cos(\theta)$ and the elevation angular grid $\Omega = \sin(\theta)$ are sampled with a period of $1/40$, ensuring that all dictionary angle pairs satisfy the condition $\Psi^2 + \Omega^2 \leq 1$. This leads to a dictionary size of 5019.

In this simulation, the paths are assumed equally strong on the average, and the pilot SNR is 10 dB. Multiple random channel realizations are considered for each pilot length, with azimuth and elevation angles chosen from the range $[-\frac{0.9\pi}{2}, \frac{0.9\pi}{2}]$. The classical orthogonal matching pursuit (OMP) algorithm [69] is utilized for the compressed-sensing-based estimation. Fig. 12 shows the resulting estimation performance, compared with the LS and RS-LS estimators, for which all inverses are replaced by pseudo-inverses when matrix inversion issues arise. The results show that when $\tau_p \geq 10$, the OMP-based estimator significantly outperforms the conventional LS estimator and provides the lowest NMSE. $\tau_p = \bar{\tau} = 44$ is the point where the RS-LS estimator has sufficient pilots to explore all the dimensions of the reduced subspace, after which there is a slight performance gap between the RS-LS and OMP algorithms. The OMP algorithm experiences a performance floor for $\tau_p > 50$, which stems from the limited dictionary size and on-grid channel estimation. More advanced compressed sensing-based algorithms may provide better performance, but they come with a significantly higher computational complexity (e.g., caused by increasing the dictionary size). Hence, compressed-sensing-based estimation methods are attractive to reduce the pilot length when estimating channels that are known to feature sparse scattering, but where the exact channel statistics are unknown. However, for less sparse channels or when we can afford longer pilots, the more computationally-friendly RS-LS estimator might be a better choice.

The goal of the dictionary design is that the channel path to any scattering cluster should be well represented by one of the dictionary vectors, in terms of having a large inner product (in the magnitude sense). The simulation example considered far-field channels because we considered the channel between a single antenna of the BS and a relatively small array of the UE. In this scenario, the assumption of having a dictionary of far-field array response vectors chosen with uniform sampling of the azimuth and elevation angular domains makes good sense, because all scattering clusters will be in the far-field of the UE. However, when considering the complete MIMO channel between the BS and UE, it becomes imperative to also account for radiative near-field effects. All the considered estimation algorithms are also applicable in this situation. When it comes to compressed sensing-based channel estimation, we must revise the dictionary design to also account for scattering clusters located in the radiative near-field of transmitter or receiver

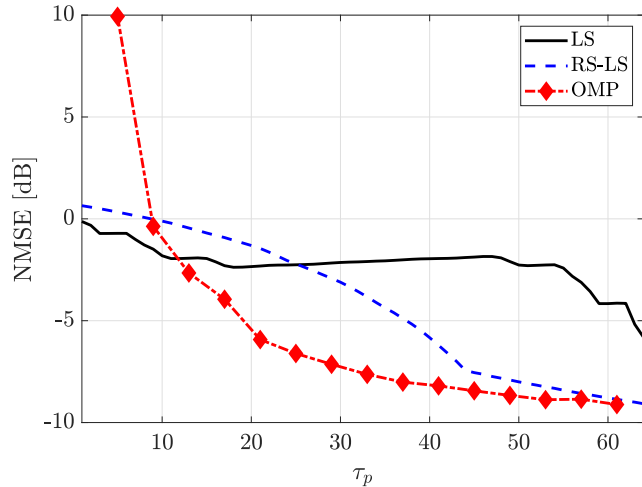


Fig. 12. The NMSE versus the pilot length (τ_p) for the LS, RS-LS, and OMP-based channel estimators in a propagation environment with sparse scattering.

[70], [71]. Recent proposals advocate for using a polar-domain grid where the dictionary vectors represent channels to points in the different angles and depths [23], [71], [72]. This design is more challenging compared to the far-field case due to the new depth dimension, and the fact that the shape of a beam depends on both depth and angle. One basically needs to find a reasonably small collection of finite-depth and far-field beams that jointly cover all possible locations in the coverage area, in the sense of guaranteeing a large array gain anywhere. Once this dictionary has been designed, similar compressed sensing-based estimation algorithms can be applied as in the far-field [23].

V. A LINEAR SYSTEM APPROACH TO ELECTROMAGNETIC THEORY

Highly simplified models for the function of antennas and the propagation of EM waves have taken us a long way in designing wireless communication systems, culminating in Massive MIMO—the most spectrally efficient scheme yet devised. However, to further push the limits of the MIMO technology, we must not only consider radiative near-field effects (as earlier in this paper) but also thoroughly model the interactions among the antennas in an array. It seems necessary to create a close union of EM theory and communication theory. Important characteristics such as polarization must be taken into account, whose analysis is required to rigorously assess the spatial DoF [73]. We have to abandon the common assumptions of i.i.d. Rayleigh fading, independence of interference among receivers, and the neglect of mutual coupling [74]. None of these classic assumptions hold when considering large and densely spaced arrays. Revising the basic system models is a daunting prospect for most communication theorists and signal processors. The very mention of EM theory evokes non-physical scalar and vector potentials, black-box finite-element simulations, complicated partial differential equations in cylindrical or spherical coordinates, and unfamiliar Bessel and Hankel functions.

In fact, tractable physics-based communication-theoretic models are obtainable quite simply from three fundamental principles [75]: 1) Any system of antennas, operating in a fixed propagation medium, is completely characterized by an impedance matrix, which quantifies the interactions among all antennas in the system. 2) Maxwell's equations describe a linear space/time-invariant system. 3) The external EM field due to any space/time distribution of electrical currents can be exactly represented as a superposition of outgoing plane waves, both ordinary and evanescent.

In this section, we will provide a linear system approach to EM theory that is distinct from the conventional physicist's approach, which is based on vector and scalar potentials and the method of separation of variables. The linear system approach is ideally suited for the needs of the communications and signal processing communities when developing future wireless technologies.

A. Impedance matrix description of antennas

A resistor, capacitor, inductor, or an antenna is a *ported device*, having a pair of wires carrying equal and opposite currents, across which is a voltage. A system of M antennas is an M -port network [76], [77]. If operating in a linear, time-invariant medium, the relation between the M voltages and the M currents constitutes a linear time-invariant (LTI) system having a real-valued causal $M \times M$ matrix impulse response. The Fourier transform of the matrix impulse response is equal to the impedance matrix, which leads to the relation

$$\mathbf{V}(\omega) = \mathbf{Z}(\omega)\mathbf{I}(\omega), \quad (64)$$

where $\mathbf{V}(\omega) \in \mathbb{C}^M$ denotes the vector of voltages, $\mathbf{I}(\omega) \in \mathbb{C}^M$ denotes the vector of currents, and ω is the angular frequency. The diagonal elements of the impedance matrix $\mathbf{Z}(\omega) \in \mathbb{C}^{M \times M}$ are the self-impedances, while the off-diagonal elements are the mutual impedances. A valid impedance matrix must satisfy four properties:

- Conjugate-symmetry in frequency, $\mathbf{Z}(-\omega) = \mathbf{Z}^*(\omega)$;
- Causality: the real and imaginary parts of the (n, m) th entry $Z_{nm}(\omega)$ satisfy the Kramers-Kronig relations;
- Reciprocity: non-conjugate transpose symmetry, $\mathbf{Z}^T(\omega) = \mathbf{Z}(\omega)$;
- Conservation of energy: $\text{Re}(\mathbf{Z}(\omega))$ is nonnegative-definite.

A system comprising a set of transmit antennas and a set of receive antennas has a partitioned impedance matrix, such that

$$\begin{bmatrix} \mathbf{V}_T(\omega) \\ \mathbf{V}_R(\omega) \end{bmatrix} = \begin{bmatrix} \mathbf{Z}_T(\omega) & \mathbf{Z}_{TR}(\omega) \\ \mathbf{Z}_{RT}(\omega) & \mathbf{Z}_R(\omega) \end{bmatrix} \begin{bmatrix} \mathbf{I}_T(\omega) \\ \mathbf{I}_R(\omega) \end{bmatrix}, \quad (65)$$

where $\mathbf{Z}_{TR}(\omega) = \mathbf{Z}_{RT}^T(\omega)$ due to the reciprocity property. The instantaneous transmitted sum-power is

$$p_T(t) = \mathbf{i}_T^T(t)\mathbf{v}_T(t). \quad (66)$$

For a time-harmonic source, $\mathbf{i}_T(t) = \text{Re}(\mathbf{I}_T(\omega)e^{-j\omega t})$ ² and the time-average transmitted power is

$$\begin{aligned} \bar{P}_T &= \frac{1}{2}\mathbf{I}_T^H(\omega)\text{Re}(\mathbf{Z}_T(\omega))\mathbf{I}_T(\omega) \\ &\quad + \frac{1}{2}\text{Re}(\mathbf{I}_T^H(\omega)\mathbf{Z}_{TR}(\omega)\mathbf{I}_R(\omega)). \end{aligned} \quad (67)$$

If the receive currents are equal to zero (implying that open-circuit voltages are measured) or if the receive antennas are sufficiently decoupled from the transmit antennas (the typical cellular scenario) then the receive currents do not contribute to the transmitted power. In the communication literature, the power is typically computed by summing the squared magnitudes of the currents, but this is an inaccurate procedure whenever there is significant mutual coupling, i.e., the off-diagonal entries of $\mathbf{Z}_T(\omega)$ are non-negligible.

We again stress that the impedance matrix is an *exact* description of any system of antennas operating in an LTI propagation medium. In the context of wireless communication theory, the sole purpose of EM theory is to populate the entries of the impedance matrix.

B. Space/Time Fourier solutions to Maxwell's equations

We will now derive general expressions for the EM fields that can appear in communication systems, and thereby be used for carrying data. The electric and magnetic fields are obtained as solutions to Maxwell's equations. In a homogeneous and isotropic medium, Maxwell's equations become

$$\begin{aligned} \nabla \times \mathbf{E}(t, \mathbf{p}) &= -\mu_0 \frac{\partial \mathbf{H}(t, \mathbf{p})}{\partial t} \\ \nabla \times \mathbf{H}(t, \mathbf{p}) &= \epsilon_0 \frac{\partial \mathbf{E}(t, \mathbf{p})}{\partial t} + \mathbf{J}(t, \mathbf{p}) \\ \epsilon_0 \nabla \cdot \mathbf{E}(t, \mathbf{p}) &= \rho(t, \mathbf{p}) \\ \mu_0 \nabla \cdot \mathbf{H}(t, \mathbf{p}) &= 0, \end{aligned} \quad (68)$$

where $\mathbf{p} = [x, y, z]^T$ contains the Cartesian coordinates. The EM medium is characterized by the dielectric permittivity ϵ_0 ($\frac{\text{A}\cdot\text{s}}{\text{V}\cdot\text{m}}$) and the magnetic permeability μ_0 ($\frac{\text{V}\cdot\text{s}}{\text{A}\cdot\text{m}}$). The field quantities are defined as follows: the electric field intensity, \mathbf{E} (V/m), the magnetic field intensity, \mathbf{H} (A/m), and the electric charge density, ρ ($\text{A}\cdot\text{s}/\text{m}^3$). Maxwell's equations are driven, in general, by a space/time distributed electric current density, \mathbf{J} (A/m^2).

1) *Linear space/time-invariant system*: Maxwell's equations describe a system whose inputs are the three components of the vector-valued electric current density, and whose outputs are the six components of the vector-valued electric and magnetic fields.³ If the electric current density is displaced in time and space, the corresponding electric and magnetic fields are displaced in the same way. We are dealing with a linear space/time-invariant system for which there is a 6×3

²A slight abuse of notation: elsewhere we represent the Fourier transform by $\mathbf{I}_T(\omega)$.

³The divergence of the second Maxwell equation, combined with the third, yields the charge density in terms of the current density, $\frac{\partial \rho}{\partial t} = -\nabla \cdot \mathbf{J}$.

space/time impulse response (called the *Green's function*), $\mathbf{G}(t, \mathbf{p})$, such that

$$\begin{bmatrix} \mathbf{E}(t, \mathbf{p}) \\ \mathbf{H}(t, \mathbf{p}) \end{bmatrix} = \mathbf{G}(t, \mathbf{p}) * \mathbf{J}(t, \mathbf{p}), \quad (69)$$

where $*$ denotes space/time convolution. Applying the space/time Fourier transform, $\int \int \int \int (\cdot) e^{j(\omega t - \mathbf{k}^T \mathbf{p})} dt dx dy dz$ to (69) yields⁴

$$\begin{bmatrix} \mathbf{E}(\omega, \mathbf{k}) \\ \mathbf{H}(\omega, \mathbf{k}) \end{bmatrix} = \mathbf{G}(\omega, \mathbf{k}) \mathbf{J}(\omega, \mathbf{k}), \quad (70)$$

where the convolution becomes a multiplication, and $\mathbf{k} = [k_x, k_y, k_z]^T$ is the wavenumber. As previously,

$$\kappa = \omega \sqrt{\epsilon_0 \mu_0} = \frac{\omega}{c} = \frac{2\pi}{\lambda}. \quad (71)$$

We can take the space/time Fourier transform of both sides of the four Maxwell's equations (68) and then algebraically obtain a remarkably simple analytical solution for the electric and magnetic fields due to the source distribution:

$$\begin{bmatrix} \mathbf{E}(\omega, \mathbf{k}) \\ \mathbf{H}(\omega, \mathbf{k}) \end{bmatrix} = \frac{1}{\mathbf{k}^T \mathbf{k} - \kappa^2} \begin{bmatrix} (\kappa^2 \mathbf{I}_3 - \mathbf{k} \mathbf{k}^T) \\ -j\omega \epsilon_0 \\ (\mathbf{j} \mathbf{k} \times) \end{bmatrix} \mathbf{J}(\omega, \mathbf{k}), \quad (72)$$

where

$$\mathbf{k} \times = \begin{bmatrix} 0 & -k_z & k_y \\ k_z & 0 & -k_x \\ -k_y & k_x & 0 \end{bmatrix}. \quad (73)$$

2) *Plane-wave solution*: All the action of the wave equation is embodied in the denominator polynomial of (72) which constitutes two simple poles in any of the three wavenumbers:

$$\mathbf{k}^T \mathbf{k} - \kappa^2 = (k_z - \gamma)(k_z + \gamma), \quad (74)$$

where

$$\gamma(\omega, k_x, k_y) = \sqrt{\kappa^2 - k_x^2 - k_y^2}. \quad (75)$$

The Sommerfeld rule for choosing the sign of the square root is that the imaginary part of the square root should be non-negative, and when the imaginary part is zero, the real part should be non-negative. Expressed as functions of $\{\omega, k_x, k_y, z\}$, the electric and magnetic fields satisfy an ordinary second-order differential equation in z . Suppose that the electric current distribution is confined to the slab $|z| \leq z_0$ (e.g., $\mathbf{J}(\omega, \mathbf{p}) = \mathbf{0}$, $\forall |z| > z_0$). Then for $|z| > z_0$, the field has to satisfy the 1D homogeneous Helmholtz equation,

$$\left[\frac{\partial^2}{\partial z^2} + (\kappa^2 - k_x^2 - k_y^2) \right] \begin{bmatrix} \mathbf{E}(\omega, k_x, k_y, z) \\ \mathbf{H}(\omega, k_x, k_y, z) \end{bmatrix} = 0, \quad (76)$$

with z -dependence given by

$$\begin{bmatrix} \mathbf{E}(\omega, k_x, k_y, z) \\ \mathbf{H}(\omega, k_x, k_y, z) \end{bmatrix} \propto e^{j\gamma(\omega, k_x, k_y)|z|}, \quad |z| > z_0, \quad (77)$$

which implies that the field constitutes outgoing plane-waves on either side of the source distribution.

To obtain the electric and magnetic fields as functions of (ω, \mathbf{p}) in the source-free region (i.e., outside the antenna), we

⁴We use *mixed notation*: the same symbol is used for the space/time function and its transform.

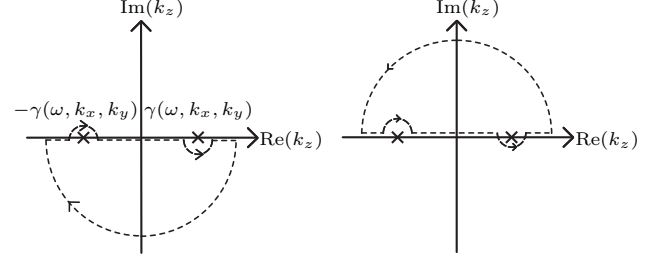


Fig. 13. Integration contours for plane-wave expansion of the outgoing field. a) $z < -z_0$; b) $z > z_0$.

take the inverse wavenumber transform of (72). We perform the k_z integral by evaluating the residues of the two poles, closing the contour in the upper-half plane for $z > z_0$, and in the lower-half plane for $z < z_0$, illustrated in Fig. 13:

$$\begin{bmatrix} \mathbf{E}(\omega, \mathbf{p}) \\ \mathbf{H}(\omega, \mathbf{p}) \end{bmatrix} = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \frac{dk_x dk_y}{(2\pi)^2} \frac{j}{2\gamma(\omega, k_x, k_y)} \cdot \begin{bmatrix} (\kappa^2 \mathbf{I}_3 - \mathbf{k} \mathbf{k}^T) \\ -j\omega \epsilon_0 \\ (\mathbf{j} \mathbf{k} \times) \end{bmatrix} \mathbf{J}(\omega, \mathbf{k}) e^{j\mathbf{k}^T \mathbf{p}} \Big|_{k_z = \text{sgn}(z) \cdot \gamma(\omega, k_x, k_y)}, \quad (78)$$

$|z| > z_0.$

This formula embodies the most important principle of wave propagation: For *any* compact space-time distribution of electric current, the resulting external EM field comprises a superposition of outgoing plane-waves [78]. The plane-waves are of two types: *ordinary (propagating)* for $k_x^2 + k_y^2 < \kappa^2$ where k_z is positive-real and *evanescent (inhomogeneous)* for $k_x^2 + k_y^2 > \kappa^2$ where k_z is positive-imaginary and the wave decays exponentially fast in $|z|$. The evanescent waves carry only reactive power in the z -direction. The difference between these types was previously illustrated in Fig. 8, where we concluded that only the ordinary type provides spatial DoF that can be used to carry data to the radiative near-field and far-field of the source.

Every (k_x, k_y) represents a plane-wave $e^{j(k_x x + k_y y \pm \gamma z)}$ propagating in the $\pm z$ directions. For each of these plane-waves, the wavenumber vector, \mathbf{k} , and the electric and magnetic fields are mutually orthogonal.⁵

The external outgoing field depends only on $\mathbf{J}(\omega, k_x, k_y, \pm \gamma(\omega, k_x, k_y))$, implying that more than one current source distribution can generate the *same* external field. It is the variability in the external field that can be used to carry data, and thereby determine the spatial DoF.

3) *The plane-wave representation of the spherical wave*: Consider the Helmholtz equation (i.e., the temporal Fourier transform of the wave equation), driven by a spatial impulse,

$$(\nabla^2 + \kappa^2) F(\omega, \mathbf{p}) = -4\pi \delta(x) \delta(y) \delta(z). \quad (79)$$

It is not yet obvious, but the solution is the spherical wave, $F(\omega, \mathbf{p}) = \frac{e^{j\kappa|\mathbf{p}|}}{|\mathbf{p}|}$. As before, we take spatial Fourier transforms of both sides of (79), algebraically solve for $F(\omega, \mathbf{k})$, and then

⁵Here we mean “complex-orthogonal”, i.e., $\mathbf{a}^T \mathbf{b} = 0$, not $\mathbf{a}^H \mathbf{b} = 0$.

take inverse wavenumber transforms to obtain the plane-wave representation called the *Weyl integral* [79]:

$$F(\omega, \mathbf{p}) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \frac{dk_x dk_y}{(2\pi)^2} \frac{j2\pi}{\gamma(\omega, k_x, k_y)} \cdot e^{j(k_x x + k_y y + \gamma(\omega, k_x, k_y)|z|)}. \quad (80)$$

To show that (80) is indeed representing a spherical wave, we switch to cylindrical coordinates, and use the fact that (79) is spherically symmetric, implying that $F(\omega, \mathbf{p})$ is also spherically symmetric:

$$\begin{aligned} F(\omega, \mathbf{p}) &= F(\omega, 0, 0, |\mathbf{p}|) \\ &= \int_0^{\infty} \int_0^{2\pi} \frac{k_r dk_r d\phi}{2\pi} \frac{j}{\sqrt{\kappa^2 - k_r^2}} \cdot e^{j\sqrt{\kappa^2 - k_r^2}|\mathbf{p}|} \\ &= - \left. \frac{e^{j\sqrt{\kappa^2 - k_r^2}|\mathbf{p}|}}{|\mathbf{p}|} \right|_0^{\infty} = \frac{e^{j\kappa|\mathbf{p}|}}{|\mathbf{p}|}, \end{aligned} \quad (81)$$

where the choice of integration contour is an exercise in the application of the Cauchy integral theorem: the contour cannot cross the branch cut associated with the square-root singularity. Hence, we have established that the spherical wave $\frac{e^{j\kappa|\mathbf{p}|}}{|\mathbf{p}|}$ can be expanded as the integral in (80) over all plane-waves, both propagating and evanescent.

It is well-known that many wavelengths away, a spherical wave looks like a plane-wave *locally*. When the observer is sufficiently small to observe a plane-wave, we say that it is in the far-field. It is much less obvious that the spherical wave can exactly be represented by a superposition of plane-waves. The local behavior of the spherical wave in the far-field can be inferred by application of the method of stationary phase: for $|\mathbf{p}| \gg \lambda$, the phase of the integrand of (80) oscillates violently, and the only significant contribution to the integral occurs in the vicinity of $\mathbf{k} = \kappa\mathbf{p}/|\mathbf{p}|$.

We note that (72) is a product, in the frequency/wavenumber domain, of three terms, equivalent to convolutions in the space/time domain. The first two terms comprise the Green's function

$$\mathbf{G}(\omega, \mathbf{k}) = \frac{1}{\mathbf{k}^T \mathbf{k} - \kappa^2} \cdot \begin{bmatrix} \left(\frac{\kappa^2 \mathbf{I}_3 - \mathbf{k} \mathbf{k}^T}{-j\omega\epsilon_0} \right) \\ (\mathbf{j} \mathbf{k} \times) \end{bmatrix}. \quad (82)$$

Given the Fourier transform relation, $\frac{1}{\mathbf{k}^T \mathbf{k} - \kappa^2} \leftrightarrow \frac{e^{j\kappa|\mathbf{p}|}}{4\pi|\mathbf{p}|}$, we directly obtain the Green's function in the space/frequency domain as

$$\mathbf{G}(\omega, \mathbf{p}) = \begin{bmatrix} \mathbf{G}_E(\omega, \mathbf{p}) \\ \mathbf{G}_H(\omega, \mathbf{p}) \end{bmatrix}, \quad (83)$$

where

$$\begin{aligned} \mathbf{G}_E(\omega, \mathbf{p}) &= \frac{1}{-j\omega\epsilon_0} \left(\kappa^2 \mathbf{I}_3 + \nabla \nabla^T \right) \frac{e^{j\kappa|\mathbf{p}|}}{4\pi|\mathbf{p}|} \\ &= \frac{1}{-j\omega\epsilon_0} \begin{bmatrix} \frac{\partial^2}{\partial x^2} + \kappa^2 & \frac{\partial^2}{\partial x \partial y} & \frac{\partial^2}{\partial x \partial z} \\ \frac{\partial^2}{\partial y \partial x} & \frac{\partial^2}{\partial y^2} + \kappa^2 & \frac{\partial^2}{\partial y \partial z} \\ \frac{\partial^2}{\partial z \partial x} & \frac{\partial^2}{\partial z \partial y} & \frac{\partial^2}{\partial z^2} + \kappa^2 \end{bmatrix} \frac{e^{j\kappa|\mathbf{p}|}}{4\pi|\mathbf{p}|}, \end{aligned}$$

and

$$\mathbf{G}_H(\omega, \mathbf{p}) = [\nabla \times] \frac{e^{j\kappa|\mathbf{p}|}}{4\pi|\mathbf{p}|} = \begin{bmatrix} 0 & -\frac{\partial}{\partial z} & \frac{\partial}{\partial y} \\ \frac{\partial}{\partial z} & 0 & -\frac{\partial}{\partial x} \\ -\frac{\partial}{\partial y} & \frac{\partial}{\partial x} & 0 \end{bmatrix} \frac{e^{j\kappa|\mathbf{p}|}}{4\pi|\mathbf{p}|}.$$

For a particular electric current distribution, $\mathbf{J}(\omega, \mathbf{p})$, we can perform a spatial convolution with the Green's function to obtain the electric and magnetic fields as superpositions of spherical waves and their first and second spatial derivatives.

We have now described two solution techniques for Maxwell's equations. The first one in (78) represents the electric and magnetic fields as a superposition of plane-waves. The second one in (69) and (83) represents the fields as a superposition of spherical waves. The spherical-wave representation is useful for computing antenna self-impedances. In general, the plane-wave representation is vastly superior. Firstly, the outgoing fields can be efficiently computed via 2D discrete Fourier transform (DFTs). Secondly, plane-waves, unlike spherical waves, are origin-free so we must only characterize their respective strength. Thirdly, there is a particularly simple description of the generation of plane-waves by Cartesian-grid antenna arrays. Finally, propagation in a horizontally-stratified medium can be solved by expanding the source distribution in plane-waves, propagating the constituent plane-waves through the parallel layers of the medium, and finally integrating over horizontal wavenumber to obtain the EM field within each layer.

C. Mutual- and self-impedance

The creation of an electric current distribution entails the exertion of power (both real and reactive) to drive the current against an electric field which itself arises from the combination of the same current distribution and other current distributions. The computation of this power yields either the self-impedance of an antenna or the mutual impedance between two antennas. The instantaneous power associated with the interaction of a current density and an electric field is

$$p(t) = - \int \mathbf{J}^T(t, \mathbf{p}) \mathbf{E}(t, \mathbf{p}) d\mathbf{p}. \quad (84)$$

An antenna, located at the position \mathbf{p}_1 , has the associated electric current density

$$\mathbf{J}_1(t, \mathbf{p}) = \text{Re} \left(I_1(\omega) \mathbf{s}(\mathbf{p} - \mathbf{p}_1) e^{-j\omega t} \right), \quad (85)$$

where $\mathbf{s}(\mathbf{p})$ is a real-valued function of space that describes the shape of the antenna's current distribution. The convolution of the current distribution with the Green's function is equal to an electric field, $\mathbf{E}_1(t, \mathbf{p})$. A second antenna, located at the position \mathbf{p}_2 , creates a current distribution, $\mathbf{J}_2(t, \mathbf{p})$, with an expenditure of instantaneous power due to the mutual coupling, $p_{21}(t) = - \int \mathbf{J}_2^T(t, \mathbf{p}) \mathbf{E}_1(t, \mathbf{p}) d\mathbf{p}$. This integral yields the mutual impedance between the two antennas:

$$\begin{aligned} Z(\omega, \mathbf{p}_1 - \mathbf{p}_2) &= - \int \int \mathbf{s}^T(\mathbf{p}) \mathbf{G}_E(\omega, \mathbf{p} - \mathbf{p}' + \mathbf{p}_1 - \mathbf{p}_2) \mathbf{s}(\mathbf{p}') d\mathbf{p} d\mathbf{p}'. \end{aligned} \quad (86)$$

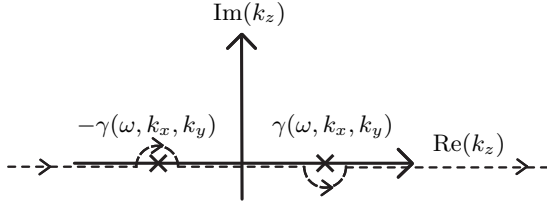


Fig. 14. The only contribution to the real part of the self-impedance are the half-residues associated with the poles for $k_x^2 + k_y^2 < \kappa^2$.

For example, an incremental vertical electric dipole of length L_0 is characterized by $\mathbf{s}(\mathbf{p}) = L_0 \mathbf{e}_z \delta(x) \delta(y) \delta(z)$, and the mutual impedance is

$$\begin{aligned} Z(\omega, \mathbf{p}_1 - \mathbf{p}_2) &= -L_0^2 \mathbf{e}_z^T \mathbf{G}_E(\omega, \mathbf{p}_1 - \mathbf{p}_2) \mathbf{e}_z \\ &= \frac{L_0^2}{j\omega\epsilon_0} \left[\frac{\partial^2}{\partial z^2} + \kappa^2 \right] \frac{e^{j\kappa|\mathbf{p}|}}{4\pi|\mathbf{p}|}. \end{aligned} \quad (87)$$

For $\mathbf{p}_1 = \mathbf{p}_2$, (87) yields the self-impedance; the imaginary part is infinite, but the real part (that figures in the computation of transmitted power) is finite.

An incremental current loop (magnetic dipole) with area A_0 , oriented in the z -direction, is characterized by $\mathbf{s}(\mathbf{p}) = A_0 \delta(x) \delta(y) \delta(z) \left[\frac{\partial}{\partial y} - \frac{\partial}{\partial x} \ 0 \right]^T$, for which the mutual impedance is

$$Z(\omega, \mathbf{p}_1 - \mathbf{p}_2) = \frac{A_0^2}{j\omega\epsilon_0} \left[\frac{\partial^2}{\partial x^2} + \frac{\partial^2}{\partial y^2} \right] \frac{e^{j\kappa|\mathbf{p}|}}{4\pi|\mathbf{p}|}. \quad (88)$$

The expression for the mutual impedance in (86) can be written in the wavenumber domain as

$$Z(\omega, \mathbf{p}) = \int \frac{d\mathbf{k}}{(2\pi)^3} \mathbf{S}^H(\mathbf{k}) \mathbf{G}_E(\omega, \mathbf{k}) \mathbf{S}(\mathbf{k}) e^{j\mathbf{k}^H \mathbf{p}}. \quad (89)$$

For $|\mathbf{p}|$ greater than the diameter of the antenna, we can extract the two residues to obtain the plane-wave representation, which can be computed via 2D DFT, as

$$\begin{aligned} Z(\omega, \mathbf{p}) &= \int \int \frac{dk_x dk_y}{(2\pi)^2} \frac{1}{2\gamma} \mathbf{S}^H(\mathbf{k}) \mathbf{G}_E(\omega, \mathbf{k}) \\ &\quad \cdot \mathbf{S}(\mathbf{k}) e^{j(k_x x + k_y y + k_z |z|)} \Big|_{k_z = \gamma(\omega, k_x, k_y)}. \end{aligned} \quad (90)$$

For $\mathbf{p} = \mathbf{0}$, (89) appears to give a purely imaginary-valued self-impedance, because \mathbf{G}_E (82) is imaginary-valued for real-valued \mathbf{k} . In fact, this does not happen, because the k_z contour has to be indented; the sole contributions to the real part of self-impedance are the half-residues associated with the propagating plane-waves as shown in Fig. 14, which becomes

$$\begin{aligned} \text{Re}(Z(\omega, \mathbf{0})) &= \int \int_{k_x^2 + k_y^2 < \kappa^2} \frac{dk_x dk_y}{(2\pi)^2} \\ &\quad \cdot \frac{\mathbf{S}^H(\mathbf{k}) (\kappa^2 \mathbf{I}_3 - \mathbf{k} \mathbf{k}^T) \mathbf{S}(\mathbf{k})}{2\omega\epsilon_0\gamma} \Big|_{k_z = \gamma(\omega, k_x, k_y)}. \end{aligned} \quad (91)$$

D. Applications of EM theory to communications

By utilizing the exact EM models, we can gain insights into the design of UM-MIMO communication systems. We will revisit the DoF concept and its connection to polarization, take a closer look at the evanescent waves, and finally discuss how to model thermal noise.

1) *Polarization and degrees-of-freedom:* The plane-wave expansion in (78) implies that, for every (ω, k_x, k_y) , there are two plane-waves having wavenumber vector $\mathbf{k} = [k_x \ k_y \ \pm \gamma]^T$. Recall that the wavenumber vector and the electric and magnetic field vectors are mutually orthogonal. Consequently, for each of the $\pm z$ -waves, the electric and magnetic field vectors are confined to a two-dimensional subspace. One way to characterize the subspace is a set of three mutually orthogonal unit vectors, which form a unitary matrix

$$\begin{aligned} \Psi(\mathbf{k}) &= \begin{bmatrix} \frac{\mathbf{k}}{\kappa} & \boldsymbol{\psi}_v & \boldsymbol{\psi}_h \end{bmatrix} \\ &= \begin{bmatrix} +\frac{k_x}{\kappa} & +\frac{k_x k_z}{k_r \kappa} & -\frac{k_y}{k_r} \\ +\frac{k_y}{\kappa} & +\frac{k_y k_z}{k_r \kappa} & +\frac{k_x}{k_r} \\ +\frac{k_z}{\kappa} & -\frac{k_x}{\kappa} & 0 \end{bmatrix}, \end{aligned} \quad (92)$$

where $k_r = \sqrt{k_x^2 + k_y^2}$ and $k_z = \pm \sqrt{\kappa^2 - k_r^2}$. The first column vector is the normalized wavenumber vector, the second lies in a vertical plane, and the third lies in the horizontal, (x, y) , plane. We can re-write the plane-wave representation in (78) in terms of horizontally and vertically polarized plane-waves by projecting both the current distribution and the E and H fields onto the orthogonal unit vectors (note that \mathbf{E} and $Z_0 \mathbf{H}$ have the same physical units, where $Z_0 = \sqrt{\frac{\mu_0}{\epsilon_0}}$ represents the characteristic impedance, thereby simplifying the expression):

$$\begin{aligned} \begin{bmatrix} \mathbf{E}(\omega, \mathbf{p}) \\ Z_0 \mathbf{H}(\omega, \mathbf{p}) \end{bmatrix} &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \frac{dk_x dk_y}{(2\pi)^2} \frac{-\mu_0}{2\gamma} e^{j(k_x x + k_y y + \gamma |z|)} \\ &\quad \cdot \begin{cases} \left(A_h^+ \begin{bmatrix} \boldsymbol{\psi}_h \\ -\boldsymbol{\psi}_v \end{bmatrix} + A_v^+ \begin{bmatrix} \boldsymbol{\psi}_v \\ \boldsymbol{\psi}_h \end{bmatrix} \right) \Big|_{k_z = +\gamma}, & z > +z_0 \\ \left(A_h^- \begin{bmatrix} \boldsymbol{\psi}_h \\ -\boldsymbol{\psi}_v \end{bmatrix} + A_v^- \begin{bmatrix} \boldsymbol{\psi}_v \\ \boldsymbol{\psi}_h \end{bmatrix} \right) \Big|_{k_z = -\gamma}, & z < -z_0 \end{cases}, \end{aligned} \quad (93)$$

where the *polarization amplitudes*, A_h^+ , A_v^+ , A_h^- , A_v^- , are related to the current density by

$$\begin{bmatrix} A_h^\pm(\omega, k_x, k_y) \\ A_v^\pm(\omega, k_x, k_y) \end{bmatrix} = \left(\begin{bmatrix} \boldsymbol{\psi}_h^T(\mathbf{k}) \\ \boldsymbol{\psi}_v^T(\mathbf{k}) \end{bmatrix} \mathbf{J}(\omega, \mathbf{k}) \right) \Big|_{k_z = \pm \gamma}. \quad (94)$$

The polarization coordinate system is illustrated in Fig. 15.

When spatial DoF were discussed in Section III, we only considered a single polarization. We will now revisit this concept by considering the extent to which an array of antennas can control the spectrum of polarimetric plane-waves.

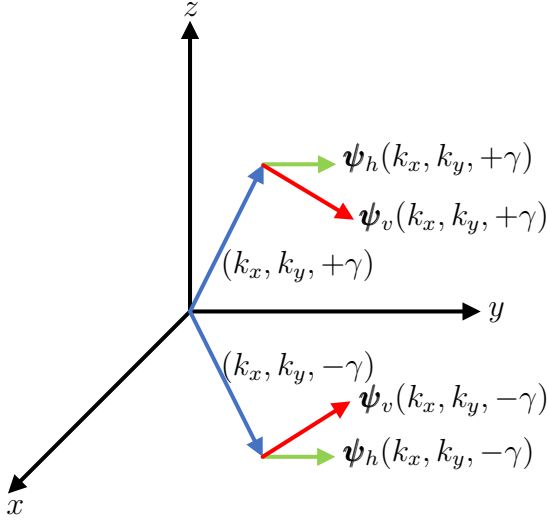


Fig. 15. For every pair of horizontal wavenumbers (k_x, k_y) there are four linearly independent plane-waves: in either the plus- or minus- z directions an h -wave whose E -field lies in the horizontal plane, ψ_h , and a v -wave whose E -field lies in a vertical plane, ψ_v .

Consider a $N \times N$ UPA in the (x, y) -plane with antenna spacing Δ . The resulting current density is

$$\mathbf{J}(\omega, \mathbf{p}) = \sum_{n_x=0}^{N-1} \sum_{n_y=0}^{N-1} I_{n_x, n_y}(\omega) \cdot \mathbf{s} \left(x - \Delta \left[n_x - \frac{N-1}{2} \right], y - \Delta \left[n_y - \frac{N-1}{2} \right], z \right). \quad (95)$$

The currents that drive the antennas are made equal to an inverse 2D DFT:

$$I_{n_x, n_y}(\omega) = \frac{1}{N} \sum_{m_x=0}^{N-1} \sum_{m_y=0}^{N-1} \hat{I}_{m_x, m_y}(\omega) e^{-\frac{j\pi(N-1)(m_x+m_y)}{N}} \cdot e^{\frac{j2\pi(m_x n_x + m_y n_y)}{N}}. \quad (96)$$

We substitute (96) into (95), and then take the wavenumber Fourier transforms to obtain

$$\mathbf{J}(\omega, \mathbf{k}) = \frac{\mathbf{S}(\mathbf{k})}{N} \sum_{m_x=0}^{N-1} \sum_{m_y=0}^{N-1} \hat{I}_{m_x, m_y}(\omega) \cdot \frac{\sin \left[\frac{(N-1)\Delta}{2} \cdot \left(k_x - \frac{2\pi m_x}{N\Delta} \right) \right]}{\sin \left[\frac{\Delta}{2} \cdot \left(k_x - \frac{2\pi m_x}{N\Delta} \right) \right]} \cdot \frac{\sin \left[\frac{(N-1)\Delta}{2} \cdot \left(k_y - \frac{2\pi m_y}{N\Delta} \right) \right]}{\sin \left[\frac{\Delta}{2} \cdot \left(k_y - \frac{2\pi m_y}{N\Delta} \right) \right]}. \quad (97)$$

We note that these Dirichlet kernel functions are periodic in the wavenumber, with a period equal to $\frac{2\pi}{\Delta}$. Furthermore, the N^2 Dirichlet-products are mutually orthogonal as well. To avoid aliasing the propagating plane-waves, the spacing must be no more than half of a wavelength, so $\Delta \leq \frac{\lambda}{2}$.

Associated with each DFT current coefficient, $\hat{I}_{m_x, m_y}(\omega)$ is a bundle of plane-waves, centered at $(k_x, k_y) = \left(\frac{2\pi m_x}{N\Delta}, \frac{2\pi m_y}{N\Delta} \right)$, occupying a square region $\frac{2\pi}{N\Delta}$ on a side.⁶ For sufficiently large $|\mathbf{p}|$ (i.e., in the far-field), the dominant contribution to the field is a single DFT coefficient such that $\left(\frac{x}{|\mathbf{p}|}, \frac{y}{|\mathbf{p}|} \right) = \left(\frac{2\pi m_x}{\kappa N \Delta}, \frac{2\pi m_y}{\kappa N \Delta} \right)$.

Recall that the propagating plane-waves correspond to $k_x^2 + k_y^2 \leq \kappa^2$, or $m_x^2 + m_y^2 \leq \left(\frac{N\Delta}{\lambda} \right)^2$, while the others are evanescent. The number of propagating plane-waves is approximately $\pi \left(\frac{N\Delta}{\lambda} \right)^2$, which was stated in (41) as the spatial DoF of a UPA.

The preceding calculation has profound implications for planar and volumetric arrays. Firstly, the activity of halving the spacing between the antennas while simultaneously quadrupling the number of antennas does not enable the UPA to create any additional propagating plane-waves. A UPA with $\lambda/2$ -spacing is sufficient to control all spatial DoF, in the sense of giving full control of all possible propagating plane-waves, subject to the wavenumber resolution of the array. This is consistent with the discussion in Section III-D.

Secondly, although the considered UPA can create $\pi \left(\frac{N\Delta}{\lambda} \right)^2$ propagating plane-waves, there is only one DoF per plane wave so we cannot distinguish between the four waves sharing the same horizontal wavenumbers k_x, k_y (i.e., $\pm z$, and the two polarizations). A different kind of array design is required to obtain four DoF per plane-wave, thereby quadrupling the total spatial DoF. Two possible designs are:

- Four parallel planar arrays, separated in z by at least $\lambda/2$, with two arrays employing, for example, vertical electric dipoles, and two arrays employing vertical magnetic dipoles.
- A single polarimetric planar array that simultaneously employs, for example, x - and y -electric dipoles, and x - and y -magnetic dipoles. Not all combinations of the six types of antennas are admissible. For example, a vertical magnetic dipole combined with two horizontal electric dipoles is linearly redundant in view of the identity $\mathbf{H} \propto \nabla \times \mathbf{E}$ [80].

Except as noted above, the expansion of a planar array into a volume array does not yield additional DoF.

Antenna polarization features have been exploited in wireless communication systems for many years [81], [82], and even predate spatial multiplexing [83].

2) *Are evanescent waves of any use?*: Conventionally, DoFs are counted as the number of linearly independent propagating plane-waves that can be created. This excludes evanescent waves, based on the fact that they decay exponentially fast in the z -direction, and carry only reactive power in the z -direction. There are two scenarios, however, which may contradict this popular notion.

The first scenario is an extreme near-field operation (i.e., in the reactive near-field region), where the evanescent wave could be a significant component of the EM field. Just as an

⁶When the range is comparable to the size of the array, the Dirichlet kernels behave, in the distributional sense, as Dirac delta functions, and every DFT current coefficient is associated with a single discrete plane-wave.

increased array aperture extends the Fraunhofer distance so that many practical communication situations can take place in the radiative-near-field, it also expands the reactive near-field so it can be used for some short-range systems.

The second scenario builds on super-directivity. A super-directive array (for example a planar xy -array) has sub-wavelength spacing [84], [85], deliberately to create strong mutual coupling. The evanescent waves decay exponentially in the z -direction, but not in the $x - y$ -directions [86]. Moreover, the evanescent waves can carry real power in the $x - y$ directions. A linear array of M antennas (say, along the x -axis), operating in end-fire mode, has a limiting gain of M^2 as $\Delta \rightarrow 0$, compared with the gain of M for $\Delta = \lambda/2$. The super-directive array is an old concept but has never been realized on a large scale because the antennas have to be driven by numerically large currents which create unacceptable I^2R losses (unless super-conductive antennas are used), and an enormous local reactive field.

3) *Dense scattering propagation*: The plane-wave expansion turns out to be the ideal theoretical tool for handling scattering propagation [80], [87]. Consider a transmit array and a receive array embedded in a scattering environment. As before, we can expand the transmitted field into outgoing plane-waves, characterized for every horizontal wavenumber pair, (k_{Tx}, k_{Ty}) , by four polarization amplitudes, $\{A_{Th}^\pm(\omega, k_{Tx}, k_{Ty}), A_{Tv}^\pm(\omega, k_{Tx}, k_{Ty})\}$. The transmit plane-waves interact with the scatterers to produce the incoming field that is measured by the receive array, which itself comprises a superposition of plane-waves. The received plane-waves are characterized by their own polarization amplitudes, $\{A_{Rh}^\pm(\omega, k_{Rx}, k_{Ry}), A_{Rv}^\pm(\omega, k_{Rx}, k_{Ry})\}$. The most general linear relation between the transmitted field and the received field is through the action of a 4×4 scattering kernel, $\mathbf{K}(\omega, k_{Rx}, k_{Ry}, k_{Tx}, k_{Ty})$,

$$\begin{aligned} \mathbf{A}_R(\omega, k_{Rx}, k_{Ry}) &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} dk_{Tx} dk_{Ty} \\ \mathbf{K}(\omega, k_{Rx}, k_{Ry}, k_{Tx}, k_{Ty}) \mathbf{A}_T(\omega, k_{Tx}, k_{Ty}), \end{aligned} \quad (98)$$

where

$$\begin{aligned} \mathbf{A}_T(\omega, k_{Tx}, k_{Ty}) &= \begin{bmatrix} A_{Th}^+(\omega, k_{Tx}, k_{Ty}) \\ A_{Tv}^+(\omega, k_{Tx}, k_{Ty}) \\ A_{Th}^-(\omega, k_{Tx}, k_{Ty}) \\ A_{Tv}^-(\omega, k_{Tx}, k_{Ty}) \end{bmatrix}, \\ \mathbf{A}_R(\omega, k_{Rx}, k_{Ry}) &= \begin{bmatrix} A_{Rh}^+(\omega, k_{Rx}, k_{Ry}) \\ A_{Rv}^+(\omega, k_{Rx}, k_{Ry}) \\ A_{Rh}^-(\omega, k_{Rx}, k_{Ry}) \\ A_{Rv}^-(\omega, k_{Rx}, k_{Ry}) \end{bmatrix}. \end{aligned} \quad (99)$$

This formulation is both physically and mathematically exact.

Particular assumptions concerning the scattering kernel result in a spatially-stationary stochastic model for the propagation. First, assume that neither transmitted nor received evanescent waves contribute significantly to the propagation, i.e., the scattering kernel vanishes if either $k_{Tx}^2 + k_{Ty}^2 > \kappa^2$ or $k_{Rx}^2 + k_{Ry}^2 > \kappa^2$. Second, assume that the sixteen elements of the scattering kernel are complex Gaussian distributed, and

independent among themselves and independent over the four wavenumbers. Then the Green's function for the propagation is spatially stationary in both transmit and receive coordinates. This is the generalization of the spectral representation for a stationary Gaussian random process to spatial random fields. This formulation results in the most general physically tenable (e.g., satisfying Maxwell's equations in free-space) spatially-stationary stochastic model for propagation.

Isotropic scattering results in a spatial autocorrelation function for the Green's function which is proportional to $\text{sinc}(\kappa|\mathbf{p}_R|) \cdot \text{sinc}(\kappa|\mathbf{p}_T|)$ which is equivalent to the Clarke model [88]. This is as close to i.i.d. Rayleigh fading as could ever exist, and features the spatial correlation behaviors shown in Fig. 9 when the antenna spacing is smaller than $\lambda/2$ or the array is planar.

4) *Rayleigh-Jeans-Clarke model for ambient thermal noise*:

The additive noise in the receiver array can also be spatially correlated. Classical statistical mechanics, when combined with elementary wave propagation theory, yields a space/time model for thermally induced ambient noise [89]. We begin by considering a resonant chamber, such as a lossless copper box. The interior of the box supports a countably infinite number of EM standing-wave (Sturm-Liouville) normal modes whose tangential electric fields vanish on the boundaries of the box [90]. According to the Equipartition Theorem, every energy-storage mode has expected thermal energy equal to $\frac{k_B T}{2}$, where k_B is Boltzmann's constant and T is absolute temperature [91]. The superposition of the modes results in a Gaussian space/time stochastic process which is stationary in time, but non-stationary in space due to the boundary conditions. As the size of the box grows large compared with the wavelength, the random EM field becomes stationary in space as well as time. At a given point in space, the temporal power-density spectrum is proportional to $k_B T \omega^2$, called the *Rayleigh-Jeans spectrum* [91]. At a particular temporal frequency, the spatial autocorrelation function is proportional to $\text{sinc}(\kappa|\mathbf{p}|)$, which is the Clarke autocorrelation [88]. This *Rayleigh-Jeans-Clarke* spectrum describes dark-sky noise, and it represents the minimum EM noise that a receive antenna would be subject to.

5) *Shannon capacity of a wireless link in a resonant chamber*: As a case study, consider a transmit antenna and a receive antenna operating inside a lossless resonant chamber [92]. The 2×2 impedance matrix (whose entries are obtained from EM theory) is purely imaginary, and it has a countably infinite number of simple poles on the real- ω axis. The transmit antenna is driven by a current source, subject to bandwidth and power constraints. The receive antenna is connected to a load resistor, R_L , which in turn is connected to the infinite impedance input of a voltage amplifier. There are two sources of noise in the receiver: the additive amplifier noise, and the Johnson noise of the load resistor, whose voltage spectral density is equal to $2k_B T R_L$.⁷

The exercise of computing the channel capacity yields two surprises: 1) For a fixed transmit power, the capacity increases

⁷The thermally-induced standing-waves in the resonant chamber do not constitute noise in addition to the resistor Johnson noise. To include them would constitute a form of double-counting.

without bound as the load resistance, R_L , approaches infinity, despite the Johnson noise spectral density becoming infinite. 2) The Shannon-optimum allocation of transmit power versus frequency avoids placing power in the vicinity of the resonant frequencies (e.g., the system poles).

The conclusion is that EM theory governs the channel modeling in wireless communications, including the available spatial DoF, interaction between the antennas and scattering environment, array design, and fundamental thermal noise. We will take a closer look at several of these aspects in the remainder of this paper. The general way to expand the DoF is to increase the array size, but one can also exploit polarization and adapt the arrays to a specific scattering environment.

VI. CIRCUIT THEORY FOR PHYSICAL CHANNEL MODELING

The last section provided a comprehensive overview of the physics governing wireless transmission through Maxwell's equations, along with an introduction to the lumped impedance representation of antennas. Nevertheless, a fundamental challenge in designing advanced transmission techniques, such as Massive MIMO systems, stems from the intricate relationship between array signal processing quantities (e.g., signals from the analog-to-digital converters (ADCs) or to the digital-to-analog converters (DACs)) and the resultant fields. This complexity is heightened by mutual coupling effects among the antenna elements, making the connection between these processing quantities and the realized fields complicated [76].

The objective of this section is to describe the analytical tools necessary to establish this connection and formulate an end-to-end model of such a communication link. This model from [93] encompasses both the antenna and radio-frequency (RF) frontend and leverages fundamental physics, such as the superposition principle, the conservation of power, and reciprocity. Throughout this section, for the sake of simplicity in notation, we may occasionally omit the frequency variable ω .

A complete MIMO transceiver system can be modeled as a noisy multiport circuit with the mutual MIMO impedance matrix \mathbf{Z}_{MIMO} from (65) as illustrated in Fig. 16. In most practical systems, the back-scattering effects between the transmit and receive antennas can be neglected and the unilateral approximation, where only the forward channel impedance is taken into account, can be made:

$$\mathbf{Z}_{\text{MIMO}}(\omega) \approx \begin{bmatrix} \mathbf{Z}_{\text{T}}(\omega) & \mathbf{0} \\ \mathbf{Z}_{\text{RT}}(\omega) & \mathbf{Z}_{\text{R}}(\omega) \end{bmatrix}, \quad (100)$$

where $\mathbf{Z}_{\text{T}}(\omega)$ is the mutual impedance matrix of the transmitting array and $\mathbf{Z}_{\text{R}}(\omega)$ is the mutual impedance matrix of the receiving array, while the propagation channel is represented by the transimpedance $\mathbf{Z}_{\text{RT}}(\omega)$. In the far-field, the latter matrix can be related to the open circuit radiation responses of both arrays and the coefficients of the channel directions

$$\mathbf{Z}_{\text{RT}}(\omega) = \sum_k g(\varphi_k, \theta_k, \omega) \mathbf{s}_{\text{R}}^{\text{OC}}(\varphi_k, \theta_k, \omega) \left(\mathbf{s}_{\text{T}}^{\text{OC}}(\varphi_k, \theta_k, \omega) \right)^{\text{T}}. \quad (101)$$

In Fig. 16, the ambient thermal noise as well as the thermal noise due to the antenna losses are represented by the vector

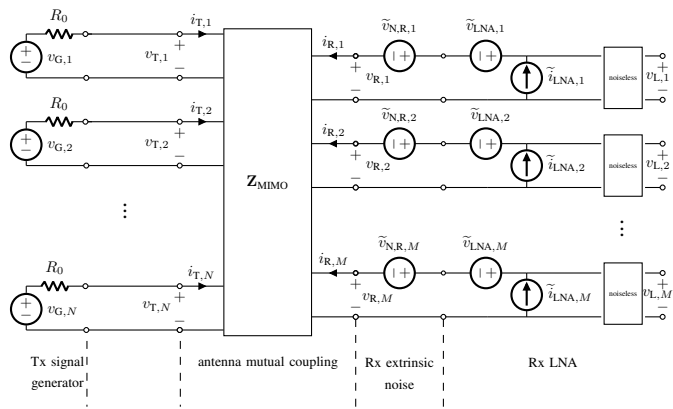


Fig. 16. A linear multiport model of a MIMO communication system showing the signal generation, antenna mutual coupling, and noise from both extrinsic sources (i.e., picked up by the antennas) and intrinsic sources (i.e., generated by low-noise amplifiers and local circuitry).

$\tilde{\mathbf{v}}_{\text{N,R}}$. Each transmit amplifier is represented as a Thévenin equivalent source with a generator open circuit voltage $v_{\text{G},n}$ and internal resistance R_0 . Each low-noise amplifier (LNA) is modeled as a noisy linear two-port network with two equivalent input noise sources. The noiseless part of the LNA can be assumed to have open-circuit input and unit voltage gain, without affecting the communication performance.

Using basic circuit analysis [76], we can establish the input-output relationship of the MIMO communication system between the input voltage $\mathbf{v}_{\text{G}}(\omega)$ and the output voltage $\mathbf{v}_{\text{L}}(\omega)$ as

$$\mathbf{v}_{\text{L}}(\omega) = \mathbf{H}(\omega) \mathbf{v}_{\text{G}}(\omega) + \mathbf{n}(\omega), \quad (102)$$

where

$$\mathbf{H}(\omega) = \mathbf{Z}_{\text{RT}}(\omega) \left(\mathbf{Z}_{\text{T}}(\omega) + R_0 \mathbf{I}_N \right)^{-1}, \quad (103a)$$

$$\mathbf{n}(\omega) = \tilde{\mathbf{v}}_{\text{N,R}}(\omega) + \tilde{\mathbf{v}}_{\text{LNA}}(\omega) + \mathbf{Z}_{\text{R}} \tilde{\mathbf{i}}_{\text{LNA}}(\omega). \quad (103b)$$

The characterization of the circuit model in Fig. 16 requires the specification of the circuit structure of the joint impedance matrix \mathbf{Z}_{MIMO} along with the statistical signal and noise properties.

A. Antenna circuit models and their key properties

An antenna can be viewed as a wave transformer that converts (single-mode) guided waves at the terminal ports into EM fields that propagate in free space and vice versa. The EM properties of an antenna array are thus characterized by its radiating/receiving patterns, the space-side scattering pattern, and the electrical multi-port behavior of its terminals [93], [94] as depicted in Fig. 17. For simplicity, we simplify the EM field, typically a complex vector quantity, as a complex-valued scalar quantity (e.g., considering a single polarization as we did in Section III). At each accessible port (m th port), the forward and backward traveling wave phasors along the antenna feed line are denoted by $a_{\alpha,m}$ and $b_{\alpha,m}$, respectively, which are related to the port current and voltages as

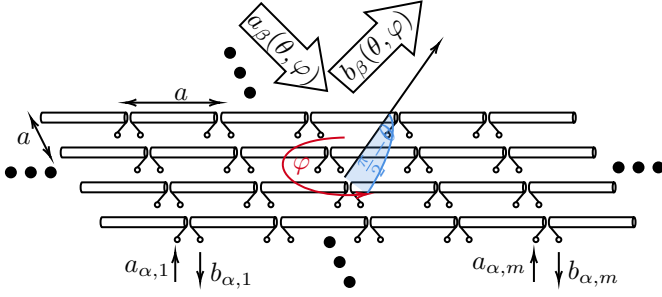


Fig. 17. Antenna scattering with m accessible ports.

$$\begin{aligned} \mathbf{a}_\alpha &= \frac{\mathbf{v} + R_0 \mathbf{i}}{2\sqrt{R_0}} = \frac{\mathbf{v}_G}{2\sqrt{R_0}}, \\ \mathbf{b}_\alpha &= \frac{\mathbf{v} - R_0 \mathbf{i}}{2\sqrt{R_0}}. \end{aligned} \quad (104)$$

In the far-field, the angular transmission characteristic of a specific polarization (i.e., the complex embedded pattern under reference resistance termination) is defined when the antenna port m is connected to a wave-source amplitude $a_{\alpha,m}$ (in \sqrt{W}) at port m while all other ports are terminated with R_0 :

$$s_m(\varphi, \theta) = \lim_{r \rightarrow \infty} -j r e^{jkr} \frac{E^{(m)}(\varphi, \theta, r)}{a_{\alpha,m}} \sqrt{\frac{1}{Z_0}}, \quad (105)$$

where $Z_0 = \sqrt{\mu_0/\epsilon_0}$ is the characteristic impedance of free space and $E^{(m)}(\varphi, \theta, r)$ is the corresponding generated electrical far-field. The complex patterns of the embedded elements, denoted as $s_m(\varphi, \theta)$, are aggregated to construct the characteristic response vector $\mathbf{s}_{\alpha\beta}(\varphi, \theta)$ for receiving and transmitting. Considering a single polarization, for simplicity, the overall description of the antenna array can be expressed through the linear scattering representation on the terminal side

$$\mathbf{b}_\alpha = \mathbf{S}_{\alpha\alpha} \mathbf{a}_\alpha + \int_0^\pi \int_{-\pi}^\pi \mathbf{s}_{\alpha\beta}(\varphi, \theta) a_\beta(\varphi, \theta) \sin(\theta) d\varphi d\theta, \quad (106)$$

where $\mathbf{a}_\alpha \triangleq [a_{\alpha,1}, \dots, a_{\alpha,M}]^T$ and $\mathbf{S}_{\alpha\alpha}$ is the scattering matrix. For the space-side, the angular spectra of the outgoing propagating wave phasors $b_\beta(\varphi, \theta)$ are expressed as linear functions of the incoming wave $a_\beta(\varphi, \theta)$ ($= g(\varphi, \theta) \cdot s$ for a single source in space) as well as the port incident phasor

$$\begin{aligned} b_\beta(\varphi, \theta) &= \mathbf{s}_{\alpha\beta}^T(\varphi, \theta) \mathbf{a}_\alpha + \\ &\int_0^\pi \int_{-\pi}^\pi s_{\beta\beta}(\varphi, \theta, \varphi', \theta') a_\beta(\varphi', \theta') \sin(\theta') d\varphi' d\theta', \end{aligned} \quad (107)$$

where $s_{\beta\beta}(\varphi, \theta, \varphi', \theta')$ is the wave back-scattering characteristic. Ideally, $s_{\beta\beta}(\varphi, \theta, \varphi', \theta')$ is effectively negligible or is just considered as part of the fixed environment. Alternatively, one can use port current and voltage phasors to describe the antenna instead of wave phasors. Substituting (104) into (106), we get the alternative impedance-based representation

$$\mathbf{v} = \mathbf{Z} \mathbf{i} + \int_0^\pi \int_{-\pi}^\pi \mathbf{s}_{\alpha\beta}^{\text{OC}}(\varphi, \theta) a_\beta(\varphi, \theta) \sin(\theta) d\varphi d\theta, \quad (108)$$

where the impedance matrix is then related to the S-matrix via the relation

$$\mathbf{Z} = R_0(\mathbf{I}_M - \mathbf{S}_{\alpha\alpha})^{-1}(\mathbf{I}_M + \mathbf{S}_{\alpha\alpha}), \quad (109)$$

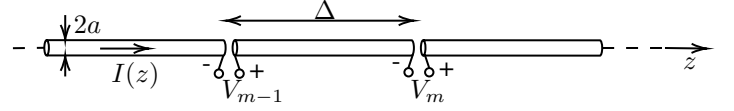


Fig. 18. Infinite uniform connected co-linear array.

and the open circuit embedded radiation pattern is related to the terminated radiation pattern as

$$\mathbf{a}_{\alpha\beta}^{\text{OC}}(\varphi, \theta) = \frac{1}{\sqrt{R_0}}(\mathbf{Z} + R_0 \mathbf{I}_M) \mathbf{s}_{\alpha\beta}(\varphi, \theta). \quad (110)$$

Assuming lossless antennas with the property

$$\mathbf{S}_{\alpha\alpha} \mathbf{S}_{\alpha\alpha}^H + \underbrace{\int_0^\pi \int_{-\pi}^\pi \mathbf{s}_{\alpha\beta}(\varphi, \theta) \mathbf{s}_{\alpha\beta}^H(\varphi, \theta) \sin(\theta) d\varphi d\theta}_{\triangleq \mathbf{B}} = \mathbf{I}_M, \quad (111)$$

then a relationship between the embedded pattern coupling matrix \mathbf{B} (also called radiation matrix) and the S-matrix $\mathbf{S}_{\alpha\alpha}$ is given by

$$\mathbf{B} = \mathbf{I}_M - \mathbf{S}_{\alpha\alpha} \mathbf{S}_{\alpha\alpha}^H = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^H. \quad (112)$$

Due to reciprocity (i.e., $\mathbf{S}_{\alpha\alpha} = \mathbf{S}_{\alpha\alpha}^T$), we can obtain when all eigenvalues of \mathbf{B} are distinct

$$\mathbf{S}_{\alpha\alpha} = \mathbf{U} \text{diag}(e^{j\alpha_1}, \dots, e^{j\alpha_M})(\mathbf{I}_M - \mathbf{\Lambda})^{1/2} \mathbf{U}^T, \quad (113)$$

with arbitrary phases $\alpha_1, \dots, \alpha_M$ (which are equal if the antennas are assumed to have identical embedded patterns). It is worth noting that the port scattering matrix $\mathbf{S}_{\alpha\alpha}$, the impedance matrix \mathbf{Z} , and the pattern coupling matrix \mathbf{B} have the same eigenvectors but with different eigenvalues.

Antenna mutual coupling, primarily occurring in the reactive near-field, can be characterized by considering the complex far-field pattern, up to a diagonal complex rotation. While this observation seems to be counter-intuitive, it follows from the uniqueness theorem in electromagnetism, wherein boundary conditions are specified at infinity. It is crucial to emphasize that the required patterns for this derivation are the embedded patterns as opposed to the isolated patterns which relativizes the usefulness of this observation.

If the antenna is lossy, then (111) becomes an inequality. A simple way to account for losses is to introduce the loss factor $\eta \leq 1$ which defines the antenna efficiency. Accordingly, (113) becomes

$$\mathbf{S}_{\alpha\alpha} = \mathbf{U} \text{diag}(e^{j\alpha_1}, \dots, e^{j\alpha_M})(\eta \mathbf{I}_M - \mathbf{\Lambda})^{1/2} \mathbf{U}^T. \quad (114)$$

Generally, characterizing the circuit behavior of large antenna arrays across different frequencies through the embedded pattern is challenging computationally as well as experimentally. In addition, the current distribution is not known in advance to solve (78) in a straightforward manner. Instead, some boundary conditions are to be imposed depending on the antenna structure. In the following, two main analytical techniques for characterizing the mutual coupling effects are discussed.

1) *Infinite arrays*: Infinite periodic arrays are generally more tractable to analyze due to the identical radiation properties of the elements and the absence of edge effects. As an example, consider the infinite co-linear one-dimensional array in Fig. 18, which can be treated as an infinitely long dipole with multiple feeds. This setup serves as a valuable approximation for large arrays and is mathematically manageable in numerous instances, thanks to the periodic structure. The feed points are periodic infinitesimal gaps distributed along the thin wire. To find the corresponding admittance matrix, we first examine the current distribution on a linear cylindrical antenna center-driven by a delta function generator. For a hollow-cylindrical antenna with a radius a and infinite length, the boundary condition in the polar coordinates system along the center-driven wire is

$$E_z(\rho = a, z) = -V\delta(z), \quad (115)$$

where V is the voltage maintained at the driving point $z = 0$. By solving (72) in the Fourier domain of z under the above boundary condition, we obtain the solution (which is a particular solution of the Pocklington's integral equation [95])

$$y(z) = \frac{I(z)}{V} = \frac{2\kappa}{\pi Z_0} \int_{-\infty}^{\infty} \frac{e^{j\alpha z}}{\beta^2 J_0(\beta a) H_0^{(2)}(\beta a)} d\alpha, \quad (116)$$

where $\beta = \sqrt{\kappa^2 - \alpha^2}$, while $J_0(\cdot)$ and $H_0^{(2)}(\cdot)$ are the zero-order Bessel and Hankel functions. The resulting magnetic field reads as (c.f. (69))

$$\begin{aligned} H_{\varphi, \text{far-field}}(\varphi, \theta, r) &= \\ \sin(\theta) e^{-j\kappa r} \frac{j\kappa}{4\pi r} \int_{-\infty}^{\infty} \frac{I(z)}{2\pi} \int_0^{2\pi} e^{-j\kappa(z \cos(\theta) + a \sin(\theta) \cos(\varphi))} d\varphi dz & \\ = \frac{jV e^{-j\kappa r}}{\pi r Z_0 \sin(\theta) H_0^{(2)}(a\kappa \sin(\theta))}. & \end{aligned} \quad (117)$$

Hence, the short-circuit pattern for an excitation port at $z = 0$ is expressed as

$$\begin{aligned} s_{\alpha\beta}^{\text{SC}}(\varphi, \theta) &= -j r e^{j\kappa r} \frac{H_{\varphi, \text{far-field}}(\varphi, \theta, r)}{V} \sqrt{Z_0} \\ &= \frac{1}{\pi \sqrt{Z_0} \sin(\theta) H_0^{(2)}(a\kappa \sin(\theta))}. \end{aligned} \quad (118)$$

To derive the impedance description of the infinite co-linear array, we uniformly discretize the admittance function in (116) with a sampling spacing of Δ . The periodic admittance spectrum can subsequently be inverted to form the impedance function

$$z[m\Delta] = \frac{Z_0 \Delta^2}{8\pi\kappa} \int_0^{\frac{2\pi}{\Delta}} \frac{e^{j\alpha m\Delta}}{\sum_{\ell=-\infty}^{\infty} \frac{1}{\beta_\ell^2 J_0(\beta_\ell a) H_0^{(2)}(\beta_\ell a)}} d\alpha, \quad (119)$$

where $\beta_\ell = \sqrt{\kappa^2 - (\alpha - \frac{2\pi\ell}{\Delta})^2}$. In addition, the open-circuit embedded pattern is obtained similarly as

$$s_{\alpha\beta}^{\text{OC}}(\varphi, \theta) = \frac{s_{\alpha\beta}^{\text{SC}}(\varphi, \theta) Z_0 \Delta}{\sum_{\ell=-\infty}^{\infty} \frac{4\kappa}{\beta_\ell^2 J_0(\beta_\ell a) H_0^{(2)}(\beta_\ell a)}}, \quad \beta_\ell = \sqrt{\kappa^2 - (\kappa \cos(\theta) - \frac{2\pi\ell}{\Delta})^2}. \quad (120)$$

The impedance description of the infinite array is more appropriate than the admittance one to approximate the actual finite case since open-circuiting the edge portions of the infinite array forces the edge currents to decay rapidly to zero rather than slowly in the shorted array case.

2) *Array of identical minimum scattering antennas*: The radiation patterns of embedded elements usually differ from the patterns emitted by an element when other elements in the array are absent, known as isolated element patterns, due to scattering effects from nearby elements. Finding the impedance matrix \mathbf{Z} directly from the isolated far-field radiation pattern is however possible in the case of minimum scattering antennas [96] based on the result in (90). Minimum scattering antennas are invisible (or radio-transparent) under certain reactive termination. In the canonical case, this happens when the antenna port is open-circuited. This generally means that the antenna elements, embedded into the array, induce in the structure individual current distributions within non-overlapping spheres when fed with current sources. In such a case, the embedded open circuit patterns are equivalent to the isolated pattern. Minimum scattering antennas also have other properties such as identical radiation and scattering patterns that are symmetrical with respect to the origin [96]. By changing the variables k_x and k_y to the angular coordinates, (90) can be rewritten in terms of normalized isolated radiation pattern (assumed to be known for complex angles) [96]

$$\frac{Z(\omega, \mathbf{p})}{\text{Re}(Z(\omega, \mathbf{0}))} = 2 \int_0^{2\pi} \int_0^{\frac{\pi}{2} + j\infty} e^{-j\kappa \mathbf{e}_r^T \mathbf{p}} |s^{\text{iso}}(\varphi, \theta)|^2 \sin(\theta) d\theta d\varphi, \quad (121)$$

where the θ integration is taken along the real axis, $0 \leq \theta \leq \frac{\pi}{2}$, and then along the line $\text{Re}(\theta) = \frac{\pi}{2}$, $0 \leq \text{Im}(\theta) < \infty$, with the normalized isolated angular pattern $s^{\text{iso}}(\varphi, \theta)$

$$\int_0^{2\pi} \int_0^{\pi} |s^{\text{iso}}(\varphi, \theta)|^2 \sin(\theta) d\theta d\varphi = 1. \quad (122)$$

It is worth mentioning that the result in (121) is also valid in the non-minimum scattering case if the elements are sufficiently spaced ($> \lambda$) [97]. It might be, however, inaccurate for closely spaced elements.

B. Signal power

The characterizing of signal power is a fundamental aspect in the analysis and optimization of communication systems. While in communication theory, mainly the Euclidean norm of signals is emphasized, reflecting their abstract representations and mathematical properties, UM-MIMO systems necessitate an understanding of physical power involving different physical quantities such as voltages and currents. Here, two different power metrics can be relevant in the realm of physical-consistent modeling and analysis: the radiated power and the maximum available power spectral densities obtained by taking the signal duration T_0 to infinity:

$$\begin{aligned} p_{\text{avail}}(\omega) &= \frac{1}{2} \lim_{T_0 \rightarrow \infty} \frac{1}{T_0} \mathbb{E} \left\{ \|\mathbf{a}_\alpha(\omega)\|_2^2 \right\} \\ &= \lim_{T_0 \rightarrow \infty} \frac{1}{T_0} \frac{\mathbb{E} \left\{ \|\mathbf{v}_G(\omega)\|_2^2 \right\}}{8R_0}, \end{aligned} \quad (123)$$

$$\begin{aligned}
p_T(\omega) &= \frac{1}{2} \lim_{T_0 \rightarrow \infty} \frac{1}{T_0} \mathbb{E} \{ \mathbf{a}_\alpha^H(\omega) \mathbf{B}(\omega) \mathbf{a}_\alpha(\omega) \} \\
&\leq \underset{\text{lossy antenna}}{\frac{1}{2} \lim_{T_0 \rightarrow \infty} \frac{1}{T_0} \mathbb{E} \{ \mathbf{i}_T^H(\omega) \text{Re}(\mathbf{Z}_T(\omega)) \mathbf{i}_T(\omega) \}} \\
&\leq \underset{\substack{\text{Imped.} \\ \text{Mismatch}}}{p_{\text{avail}}(\omega)}. \tag{124}
\end{aligned}$$

The radiated power density $p_T(\omega)$ is usually restricted by regulatory measures and interference policies, whereas the available power $p_{\text{avail}}(\omega)$ is constrained by device/technology limitations. In the case of perfect matching between the generator and antenna impedances ($\mathbf{Z}_T = R_0 \mathbf{I}$) along with a lossless array structure, these two powers are equivalent. In advanced antenna systems, however, impedance mismatching, structural losses, and technology limitations introduce additional model complexities that require careful characterization of the different power limitations and their relevance. At the same time, it enables flexible signal design as the antenna structure can serve as spatial filters for blocking unwanted distortion [98].

C. Noise modeling

1) *Background noise of antennas:* The multi-port network \mathbf{Z}_{MIMO} is only composed of passive components which can be assumed to have the same absolute temperature T of the environment [76]. Therefore, the noise of the joint impedance matrix \mathbf{Z}_{MIMO} originates solely from the thermal agitation of the electrons flowing inside its all passive components, a.k.a. thermal noise at the equilibrium temperature T [99]. Due to the unilateral approximation, the transmit side noise is neglected. In Fig. 16, the receive noise voltages $\tilde{v}_{\text{N,R}}(\omega)$ model the background noise of receive antennas as well as the resistive losses. When the mutual coupling is taken into account within the transmit/receive arrays, the correlation between the m th and m' th receive noise voltages $\tilde{v}_{\text{N,R},m}(\omega)$ and $\tilde{v}_{\text{N,R},m'}(\omega)$ taking the signal duration T_0 to infinity [99] is

$$\begin{aligned}
\lim_{T_0 \rightarrow \infty} \frac{1}{T_0} \mathbb{E} \{ \tilde{v}_{\text{N,R},m}(\omega) \tilde{v}_{\text{N,R},m'}^*(\omega) \} &= 4 k_B T \text{Re} (Z_{R,mm'}) \\
&\quad \forall m, m' \in [1, \dots, M]. \tag{125}
\end{aligned}$$

2) *The receive LNA model:* The LNAs are modeled as independently noisy frequency flat devices with unit gain. For the m th amplifier, the spectral second-order statistics of the noise voltage, $\tilde{v}_{\text{LNA},m}(\omega)$, and current, $\tilde{i}_{\text{LNA},m}(\omega)$, generated inside the LNA are determined using the truncated Fourier transform:

$$\lim_{T_0 \rightarrow \infty} \frac{1}{T_0} \mathbb{E} \{ |\tilde{v}_{\text{LNA},m}(\omega)|^2 \} = 4 k_B T R_{v,\text{LNA}}, \tag{126a}$$

$$\lim_{T_0 \rightarrow \infty} \frac{1}{T_0} \mathbb{E} \{ |\tilde{i}_{\text{LNA},m}(\omega)|^2 \} = 4 k_B T G_{i,\text{LNA}}, \tag{126b}$$

$$\lim_{T_0 \rightarrow \infty} \frac{1}{T_0} \mathbb{E} \{ \tilde{i}_{\text{LNA},m}(\omega) \tilde{v}_{\text{LNA},m}(\omega)^* \} = 4 k_B T \beta_{\text{LNA}}. \tag{126c}$$

Additionally, at the m th receive antenna, the LNA noise voltage, $\tilde{v}_{\text{N,LNA},m}(\omega)$, is uncorrelated with the receive noise

voltage, $\tilde{v}_{\text{N,R},m}(\omega)$. Using (103b) and (125)–(126), the noise correlation matrix is obtained as:

$$\mathbf{R}_n = 4 k_B T \left[\text{Re} \left((1 + \beta_{\text{LNA}}) \mathbf{Z}_R \right) + R_{v,\text{LNA}} \mathbf{I}_M + G_{i,\text{LNA}} \mathbf{Z}_R \mathbf{Z}_R^H \right]. \tag{127}$$

Generally, noise is spatially correlated due to mutual coupling effects. However, if the noise from the LNAs dominates or the antenna structure is lossy (\mathbf{Z}_R dominated by resistive losses), then the conventional uncorrelated noise assumption becomes more justified.

VII. ANTENNA NEAR-FIELD AND POLARIZATION MODELING

Maxwell's equations govern not only how fields propagate through the environment but also how antennas respond to fields and excitations. In prior works on wireless communication and array processing, the antenna is generally treated as an isotropic point source with linear polarization. While such a point source does not exist in practice, this model is appropriate under far-field and narrowband assumptions. As the array sizes grow larger and antennas grow more complex, however, their near-field characteristics and polarization become increasingly important. In this section, we overview prior work on realistic array modeling and a novel EM-based approach developed in [100].

Physically consistent antenna and channel models are critical in isolating the characteristics of the array from the channel, understanding antenna effects on MIMO communications, and designing systems that achieve high capacity in specific environments. There are several different aspects to this. Firstly, the radiative near-field of an antenna array extends much further than that of an individual antenna, as we elaborated on in Section II. Hence, we need models that account for the amplitude and phase variations observed at different antenna elements due to impinging spherical wavefronts. Secondly, the antenna polarization must be captured since it affects both the channel modeling and spatial DoF. The fundamental EM theory for polarized antennas was provided in Section V-D, but practical antenna hardware has impairments. Each antenna is designed for a specific polarization but reacts to both the intended polarization, known as the co-polarization, and the orthogonal cross-polarization. The ratio between the strength of these two polarizations is known as the cross-polarization discrimination and determines how well a practical antenna can isolate co-polarized signals. The antenna polarization can be modeled by accounting for both the co-polarization and cross-polarization gain patterns in the signal model [101]–[103].

State-of-the-art array models aim to analytically characterize all antenna features simultaneously. The array manifold model in [104] used a generalized array response vector model that accounts for arbitrary locations, spherical propagation, and gain patterns. An EM-based array model proposed in [105] leveraged a Hertzian dipole framework for far-field array propagation. This framework is extended in [100] to capture mutual coupling, polarization, and near-field propagation. We will provide an overview of this model below.

A. Electromagnetic-based array model

Electromagnetically, an antenna is a device that converts electrical currents to fields. Electrical signals are fed into the antenna, inducing a current distribution throughout the structure. Let \mathcal{V} denote the volume containing the antenna. For any point $\mathbf{p} \in \mathcal{V}$, we let $\mathbf{J}(\mathbf{p})$ be the current density along the antenna when it is excited by a unit current (we omitted the time index for simplicity). The radiated electric and magnetic fields can be computed from $\mathbf{J}(\mathbf{p})$ using the magnetic vector potential $\mathbf{A}(\mathbf{p})$ computed as [106]

$$\mathbf{A}(\mathbf{p}) = \mu_0 \int_{\mathcal{V}} \frac{e^{j\kappa|\mathbf{p}-\mathbf{p}'|}}{4\pi|\mathbf{p}-\mathbf{p}'|} \mathbf{J}(\mathbf{p}') d\mathbf{p}'. \quad (128)$$

The electric $\mathbf{E}(\mathbf{p})$ and the magnetic field $\mathbf{H}(\mathbf{p})$ are given by Maxwell's equations in (68). The current distribution along the antenna fully determines the radiated field patterns. Unfortunately, this current distribution can be difficult to obtain in a tractable form that allows easy computation of the magnetic vector potential.

To simplify the calculation of the radiated fields, we can apply a discretization technique to model the antenna as a number of easily characterized segments. EM computational software discretizes antenna volumes into smaller portions to solve complicated problems. We incorporate a similar approach by partitioning \mathcal{V} into K non-overlapping pieces $\{\mathcal{V}_k\}_{k=1}^K$ to obtain

$$\mathbf{A}(\mathbf{p}) = \mu_0 \sum_{k=1}^K \int_{\mathcal{V}_k} \frac{e^{j\kappa|\mathbf{p}-\mathbf{p}'|}}{4\pi|\mathbf{p}-\mathbf{p}'|} \mathbf{J}(\mathbf{p}') d\mathbf{p}'. \quad (129)$$

We assume that the discretization is fine enough such that the following two hold: 1) The current distribution over the k th segment is equal to a constant \mathbf{J}_k (i.e., $\mathbf{J}(\mathbf{p}) = \mathbf{J}_k$ for $\mathbf{p} \in \mathcal{V}_k$); 2) Each segment behaves as a point source (i.e., $\mathbf{p} - \mathbf{p}' = \mathbf{p}_k$ for $\mathbf{p}' \in \mathcal{V}_k$).

Under these assumptions, (129) can be approximated as

$$\mathbf{A}(\mathbf{p}) \approx \mu_0 \sum_{k=1}^K \int_{\mathcal{V}_k} \frac{e^{j\kappa|\mathbf{p}_k|}}{4\pi|\mathbf{p}_k|} \mathbf{J}_k d\mathbf{p}'. \quad (130)$$

$$= \mu_0 \sum_{k=1}^K \frac{e^{j\kappa|\mathbf{p}_k|}}{4\pi|\mathbf{p}_k|} \mathbf{J}_k |\mathcal{V}_k|, \quad (131)$$

where $|\mathcal{V}_k|$ denotes the volume of \mathcal{V}_k . The product of the current distribution and the volume is known as the moment, and will be denoted as $\mathbf{m}_k = \mathbf{J}_k |\mathcal{V}_k|$. If we define the magnetic vector potential of the k th segment as $\mathbf{A}_k(\mathbf{p}) = \frac{\mu_0 \mathbf{m}_k}{4\pi|\mathbf{p}_k|} e^{j\kappa|\mathbf{p}_k|}$, then we can compute (130) as

$$\mathbf{A}(\mathbf{p}) \approx \sum_{k=1}^K \mathbf{A}_k(\mathbf{p}). \quad (132)$$

In essence, the discretization procedure has converted the antenna into an array of point sources, each with magnetic vector potential \mathbf{A}_k . The overall fields can therefore be found from the superposition of the fields from the extended array.

The simplest physically consistent radiating structure is a point source with constant current, known as a Hertzian dipole. Each segment of the antenna can be modeled as a Hertzian

dipole with the current determined from the discretization procedure. The radiated field of a Hertzian dipole can be obtained in closed form following the derivation in [100], [106]. Let \mathbf{p} have a spherical representation with radial distance r , azimuth angle ϕ , and elevation angle θ as

$$\mathbf{p} = r[\cos(\phi) \cos(\theta), \sin(\phi) \cos(\theta), \sin(\theta)]^T. \quad (133)$$

We define the unit vectors of the spherical orthonormal basis at \mathbf{p} as

$$\begin{aligned} \mathbf{u}_r(\mathbf{p}) &= [\cos(\phi) \cos(\theta), \sin(\phi) \cos(\theta), \sin(\theta)]^T, \\ \mathbf{u}_\theta(\mathbf{p}) &= [-\sin(\phi), \cos(\phi), 0]^T, \\ \mathbf{u}_\phi(\mathbf{p}) &= [-\cos(\phi) \cos(\theta), -\sin(\phi) \cos(\theta), \sin(\theta)]^T. \end{aligned} \quad (134)$$

The electric field of a Hertzian dipole decays with distance along each spherical coordinate. We define the decaying amplitudes of the radial and angular components as $\alpha_{\text{rad}}(\mathbf{p})$ and $\alpha_{\text{ang}}(\mathbf{p})$, which are given by

$$\alpha_{\text{rad}}(\mathbf{p}) = \frac{e^{-j\kappa r}}{j\omega\epsilon_0 2\pi} \left(\frac{1}{r^3} + \frac{j\kappa}{r^2} \right), \quad (135)$$

$$\alpha_{\text{ang}}(\mathbf{p}) = -\frac{e^{-j\kappa r}}{j\omega\epsilon_0 4\pi} \left(\frac{1}{r^3} + \frac{j\kappa}{r^2} - \frac{\kappa^2}{r} \right). \quad (136)$$

By further defining the 3×3 matrix $\mathbf{T}(\mathbf{p})$ as the dipole field transform

$$\mathbf{T}(\mathbf{p}) = \begin{bmatrix} \alpha_{\text{rad}}(\mathbf{p}) \mathbf{u}_r(\mathbf{p}), & \alpha_{\text{ang}}(\mathbf{p}) \mathbf{u}_\phi(\mathbf{p}), & \alpha_{\text{ang}}(\mathbf{p}) \mathbf{u}_\theta(\mathbf{p}) \end{bmatrix}^T, \quad (137)$$

the electric field of a Hertzian dipole at the origin with moment \mathbf{m} is given by

$$\mathbf{E}_{\text{dip}}(\mathbf{p}) = \mathbf{T}(\mathbf{p}) \mathbf{m}. \quad (138)$$

Two important properties captured by the Hertzian dipole model that are not seen in an isotropic point source are the polarization and the near-field pattern. The Hertzian dipole has a polarization that is dependent on the orientation of the moment vector. For example, a moment vector that is aligned with the z -axis will exhibit a vertically polarized field. The amplitude variations in the field as a function of distance are also incorporated into the field expressions through $\alpha_{\text{rad}}(\mathbf{p})$ and $\alpha_{\text{ang}}(\mathbf{p})$.

The complete antenna response can be obtained by superposing the fields from Hertzian dipoles located at each segment. To translate the Hertzian dipole response from the origin to the center of the k th segment, we need to account for the dependence of the spherical basis on the position. Letting $\mathbf{Q}(\mathbf{p}) = [\mathbf{u}_r(\mathbf{p}), \mathbf{u}_\phi(\mathbf{p}), \mathbf{u}_\theta(\mathbf{p})]$ be the rotation matrix at \mathbf{p} , the field at \mathbf{p} from a dipole with moment \mathbf{m}_k located at the k th segment is

$$\mathbf{E}_{\text{dip},k}(\mathbf{p}) = \mathbf{Q}(\mathbf{p}) \mathbf{Q}(\mathbf{p}_k) \mathbf{T}(\mathbf{p}_k) \mathbf{m}_k. \quad (139)$$

The two rotations serve to express the radiated field in spherical coordinates with respect to the origin. Letting $\mathbf{R}(\mathbf{p}, \mathbf{p}_k) = \mathbf{Q}(\mathbf{p}) \mathbf{Q}(\mathbf{p}_k)$, the antenna response becomes

$$\mathbf{E}_{\text{ant}}(\mathbf{p}) = \sum_{k=1}^K \mathbf{R}(\mathbf{p}, \mathbf{p}_k) \mathbf{T}(\mathbf{p}_k) \mathbf{m}_k. \quad (140)$$

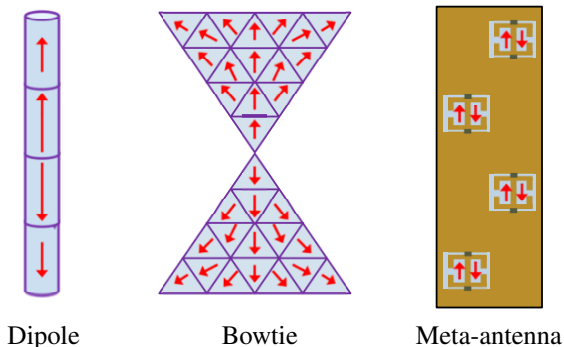


Fig. 19. Examples of the EM characterization of antennas based on discretization. Each antenna is partitioned into a large number of segments, each of which is treated as a Hertzian dipole. The radiated field from each antenna comes from the combined fields of the dipoles.

The polarization and near-field radiation of the antenna are characterized by the combined effect of all of the dipole segments.

The framework can be applied to an antenna array. Consider an array with N antenna elements excited by the length N transmit signal \mathbf{x} , where the n th antenna is discretized into K_n segments. Under the same assumptions regarding the size of the discretization, the k th segment of the n th antenna can be associated with a point $\mathbf{p}_{n,k}$ and a moment vector $\mathbf{m}_{n,k}(\mathbf{x})$. We note that because of mutual coupling, the current distribution in all of the antennas will depend on the excitation signal \mathbf{x} . The array radiated field can be computed as

$$\mathbf{E}_{\text{arr}}(\mathbf{p}) = \sum_{n=1}^N \sum_{k=1}^{K_n} \mathbf{R}(\mathbf{p}, \mathbf{p}_{n,k}) \mathbf{T}(\mathbf{p}_{n,k}) \mathbf{m}_{n,k}(\mathbf{x}). \quad (141)$$

We emphasize that the model still functions in the same way as in the single-antenna case, thus, the total radiated field is calculated by treating the array structure as a larger array of Hertzian dipoles. Linearity can also be applied to further simplify (141). Let \mathbf{e}_n denote the $N \times 1$ vector with 1 on the n th entry and zeros elsewhere and consider the $3 \times N$ matrix $\mathbf{M}_{n,k}$ with ℓ th column given by $\mathbf{m}_{n,k}(\mathbf{e}_\ell)$. It can be shown that

$$\mathbf{E}_{\text{arr}}(\mathbf{p}) = \sum_{n=1}^N \sum_{k=1}^{K_n} \mathbf{R}(\mathbf{p}, \mathbf{p}_{n,k}) \mathbf{T}(\mathbf{p}_{n,k}) \mathbf{M}_{n,k} \mathbf{x}. \quad (142)$$

This isolates the role of the transmit signal by incorporating mutual coupling between the array elements into $\mathbf{M}_{n,k}$.

We end the section with a few notes on the usefulness of this representation for an antenna array. One of the key benefits is the flexibility in characterizing arbitrary arrays since we did not make any assumptions about the structure or shape. Any type of antenna that can be meshed in the considered manner can be approximated by the sum of the Hertzian dipoles, as shown in Fig. 19. In addition, the use of Hertzian dipoles as the fundamental building blocks of each antenna means that properties such as the polarization and near-field radiation patterns of the entire array are natively incorporated into the linear model.

VIII. ANTENNA ARRAY DESIGN FOR UM-MIMO

The beamfocusing and massive spatial multiplexing characteristics of UM-MIMO systems were discussed in Section II, but the ability to control them depends on the signal processing capabilities of the antenna array [107]. The achievable SE varies between hardware architectures that have the same form factor. While half-wavelength antenna spacing is the norm in conventional Massive MIMO systems, a key feature of UM-MIMO is placing an extremely large number of antennas in a small aperture area, using spacings much smaller than $\lambda/2$. Although this feature does not inherently increase spatial DoF, this approach allows for several candidate architecture designs for spatial efficiency and new functionalities to meet the extreme requirements of 6G. These designs range from thousands of discrete antenna elements to methods using metasurfaces, and even approaches that make apertures effectively continuous.

There are many implementation challenges related UM-MIMO arrays. Firstly, large and dense array structures introduce issues related to mutual coupling, polarization discrimination, channel estimation, and near-field propagation, which have been touched upon earlier in this paper. Moreover, the large number of antennas and reduced spacing, coupled with novel signal processing techniques, can significantly increase hardware costs, energy consumption, and complexity in the production of UM-MIMO systems. The efficient realization of such high-density arrays has spurred the development of advanced antenna and device designs for 6G, including using metamaterials for antenna design or developing flexible fluid antenna systems. This section explains key antenna array technologies for 6G UM-MIMO systems, discussing how these technologies can be designed, their implementation methods, and potential application areas.

A. Uniform array-based radiation architecture

One might implement a UM-MIMO array by following the conventional approach with a UPA with discrete elements as shown in Fig. 20(a). While conventional Massive MIMO at the BSs typically involves around 64 antennas, UM-MIMO is envisioned to include many more antennas in the array [40]. This is made possible by making the array aperture larger and possibly by reducing the spacing between antennas. In [41], a BS equipped with 40 000 antennas was proposed, and [108] investigated various antenna spacing values, such as $\lambda/6$ and $\lambda/15$. Additionally, structures composed of thousands of patch antennas have been researched since 5G, including how the antenna size affects the EM modeling in large arrays [109].

Massive MIMO for the 3.5 GHz band has been implemented in 5G using the fully digital beamforming architecture shown in Fig. 20(a). In this structure, each antenna is connected to a dedicated RF chain, allowing the transceivers to generate any desirable superposition of near-field and far-field beams, thus offering high spatial flexibility [110]. This is the ideal implementation from an SE perspective, but not from a complexity viewpoint. A fully digital UM-MIMO architecture requires the use of a large number of RF chains, leading

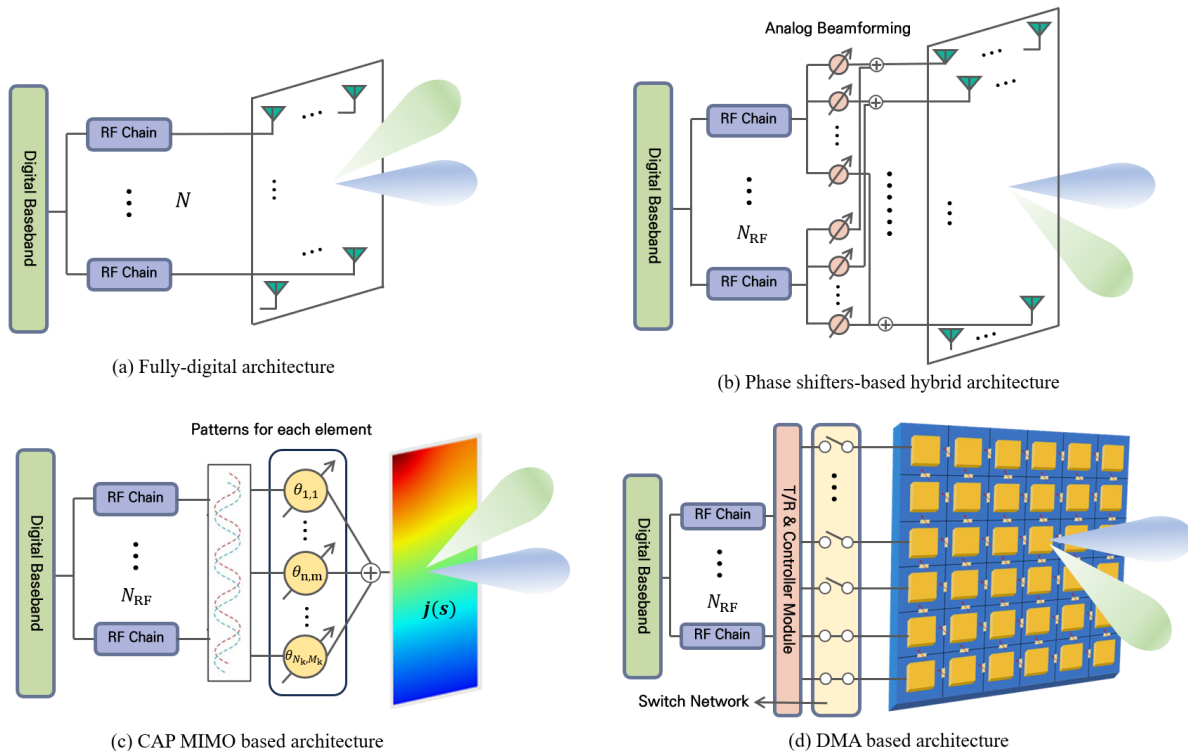


Fig. 20. The four classes of UM-MIMO antenna array architectures. a) Fully-digital architecture, b) Hybrid architecture, c) CAP MIMO-based architecture, d) DMA-based architecture.

to substantially higher design complexity, energy consumption, and cost. These issues might grow faster than linearly with the number of antennas; for example, the computational complexity for interference-suppression beamforming grows cubically with the number of antennas [15]. In principle, a fully analog architecture could be utilized that has a single RF chain that connects to the antennas through phase-shifters, but this implementation lacks the ability of spatial multiplexing—the main motivation behind UM-MIMO. Consequently, intermediate methods have been proposed that divide large arrays into several analog sub-arrays [111]. This is specifically implemented through hybrid processing, which combines the analog and digital architectures and uses fewer RF chains than the number of antennas, realized through the connection of phase shifters and viewed as a promising method to reduce hardware complexity [112], [113]. This architecture is illustrated in Fig. 20(b). These implementation simplifications can have little impact on the SE when there are more RF chains than strong multipath clusters in the propagation environment, which is particularly the case in line-of-sight scenarios.

B. Continuous aperture MIMO system

To achieve full control over the array aperture, continuous-aperture MIMO (CAP MIMO) architectures are gaining attention. CAP MIMO is a novel implementation that modulates information directly in the form of EM waves through a continuous antenna surface, ideally possessing a spatially continuous EM volume composed of an infinite number of infinitesimally small antennas as shown in Fig. 20(c). This is also being researched under the name of holographic MIMO [87], [115]

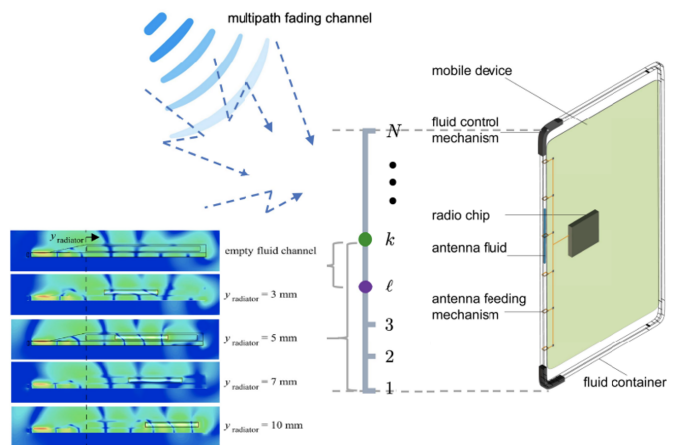


Fig. 21. Fluid antenna systems and E -field distribution of the surface wave when the fluid element changes its location [114].

and differs from traditional discrete antenna arrays by effectively using a continuous antenna aperture, thereby offering transmission efficiency and application flexibility. Although the spatial DoF is the same as for a discrete array of the same size, channel properties such as the singular values of the channel matrix can be improved within those dimensions.

CAP-based UM-MIMO systems have the capability to generate and control all current distributions on a spatially continuous surface, and can directly modulate artificially constructed

EM waves to radiate into space. This can also be seen as a specific case of point-antenna-based UPA UM-MIMO from a mathematical perspective, by considering the limit where the antenna element sizes approach zero. This can be approximately realized through the development of metamaterials and highly flexible reconfigurable antennas, fully leveraging the physical properties of space propagation to achieve high SE and energy efficiency. In [116], an antenna structure with a circular CAP plane of 10 m radius was introduced. The work in [117] presented a lens antenna with a spatially continuous aperture, where hundreds of antenna elements were placed along the focal arc of an EM lens with a radius of 5 m.

The signal processing for discrete MIMO arrays mainly consists of linear algebra operations, but these turn into continuous operations when using CAP MIMO. Designing continuous beam pattern functions for CAP MIMO is typically a non-convex problem, and it is nontrivial to make use of conventional MIMO methods. [118] introduced a pattern division multiplexing (PDM) technique that addresses these challenges by transforming the design of continuous pattern functions into a finite orthogonal basis-based projection length design. Subsequently, the EM performance comparison between CAP MIMO and conventional discrete MIMO systems was demonstrated using a non-asymptotic approach [119]. This technique provides an efficient means to simultaneously serve multiple UEs and demonstrates the potential of CAP MIMO systems in meeting the diverse performance requirements of future networks.

C. Fluid antenna system

Traditional array designs consist of highly conductive and static metal elements at fixed locations, whose joint radiation pattern is controlled by modifying the current distribution over the elements. A fundamentally different approach is fluid antenna systems (FAS) that can dynamically change shapes and positions of the radiating elements [120]–[122]. FAS represents an antenna design capable of reconfiguring characteristics such as shape, position, polarization, and radiation pattern, based on controllable conductive or dielectric elements. The flexibility and benefits over traditional solid materials have led to the emergence of numerous fluid antennas in recent years [123]–[125]. As seen in Fig. 21, the fluid element can move the antenna to one of N fixed positions (referred to as ports) within a predefined space. At any given moment, the FAS can explore the spatial fading characteristics by moving to a spatial location with better signal quality to avoid deep fades. To accomplish this, the FAS must have a fine spatial resolution, allowing the fluid to optimize for each port, thus requiring a large N value. Although the spatial footprint of FAS is small, N can be very large, allowing for diversity across numerous spatially correlated ports. From a communication perspective, FAS behaves as an antenna selection system, and can thus achieve diversity gains that can be comparable to those obtained with maximum ratio processing [126], [127].

To achieve multiplexing gains, the FAS can make use of multiple fluid antenna elements (each with N different ports), which can move independently within a specific non-overlapping space. The number of elements must be equal to or

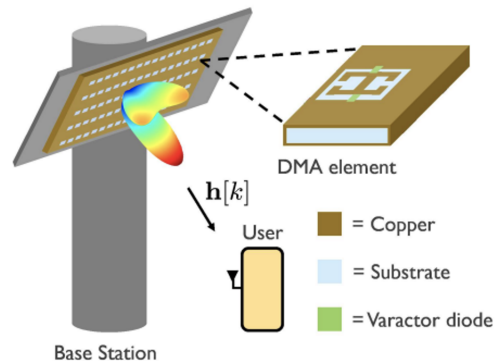


Fig. 22. Various metasurface antenna designs. A DMA-based system that uses metasurface-based antenna arrays [128].

larger than the desired number of layers. The joint optimization of the elements can be utilized to design a MIMO channel matrix with advantageous properties for a multiplexing perspective. The concept of reconfigurable antennas has been considered in the communication literature, which contains algorithms for performance optimization [129]–[132]. FAS is a potent technology for realizing these visions.

From an implementation perspective, FAS can be fabricated using liquid metals like eutectic gallium-indium (eGaIn). The movement of the fluid can be electrically controlled using methods like electrowetting (EW) or RF MEMS, as used in displays, and must have a device structure optimized for the desired radiation pattern [133], [134]. One of the major bottlenecks of FAS is response time, as the time it takes for the fluid to move to the desired position significantly impacts communication performance. While devices with sufficient response speed have not yet been presented, future improvements are expected through creative geometric methods. Additionally, some fluids used in FAS may cause undesirable chemical side effects, which should be noted. While FAS offers significant efficiency in terms of space and cost, stability issues must also be considered.

D. Metasurface-inspired antenna system

The implementation of UM-MIMO using hybrid architectures might lead to lower complexity and energy consumption than a fully digital architecture but still requires plenty of analog circuits due to the addition of numerous phase shifters. In this context, it has been recognized that metasurfaces, capable of effectively controlling EM properties, can replace traditional analog beamforming structures [135]. Metasurfaces, comprising dense arrays of reconfigurable elements smaller than the wavelength, can precisely and dynamically manipulate EM waves, actively participating in improving signal transmission and reducing interference. Research analyzing the Shannon entropy of digitally-coded metasurface communication systems has also been presented [136].

Metasurfaces have been considered in the EM literature for decades, and the development can be divided into generations. Metasurface 1.0 was based on homogeneous periodic

structures, Metasurface 2.0 allows spatial modulation with variable amplitude and phase. Further advancements have led to Metasurface 3.0, capable of real-time modulation in both time and space and featuring electrically large arrays that adapt to environmental changes [137].

Within the latest generation, the term dynamic metasurface antennas (DMA) has emerged as a promising way of implementing UM-MIMO arrays. This architecture is illustrated in Fig. 20(d). Comprising densely packed arrays of reconfigurable metadevices, DMAs utilize the tuning characteristics of each unit cell to provide analog signal processing functions without dedicated analog circuits as shown in Fig. 22. By controlling the amplitude and phase of the DMA elements, they enable low-power precoding structure that does not require phase shifters [138].

Furthermore, by exploiting the properties of metadevices to create an effective medium at much smaller scales than the wavelength, DMAs can design densely packed antenna arrays with spacings smaller than half a wavelength, enhancing beamfocusing with higher efficiency. A major challenge in implementing DMA-based communication is the Lorentzian-limited beamforming weights in the analog precoding process. Studies have been conducted to optimally adjust these limited weights [139], [140], with [141] demonstrating that the proposed approach can reduce the number of RF chains while achieving most of the beamforming gain. The paper [142] mathematically analyzed how large-scale MIMO systems using fully digital, hybrid, and DMA architectures impact beamforming capabilities in the radiative near-field. Additionally, [128] integrated DFT codebooks with DMAs in low-power MIMO systems, showing that transmitter systems using developed DMAs can surpass conventional analog beamforming systems in terms of SE and energy efficiency. An energy-efficiency analysis in [143] also demonstrates the benefits of incorporating DMAs into traditional hybrid precoding architectures. These findings suggest that DMAs, capable of being packaged in a small physical area for a wide operating frequency band, hold significant potential as a design that can surpass conventional beamforming architectures.

IX. INTERPLAY BETWEEN ARRAY DESIGN AND SIGNAL PROCESSING

The last section described four general categories of array architectures, which are summarized in Fig. 20. Architectures (a) and (b), fully digital and hybrid analog/digital, have a long history in the wireless infrastructure industry. On the other hand, (c) and (d), namely CAP and DMA, represent more modern approaches, and only time will reveal whether they will be implemented in BSs during the 6G life-cycle.

We foresee that the first 6G systems will adopt fully digital and hybrid architectures, but such UM-MIMO implementations will face many practical challenges. Improvements in hardware technology have without doubt made it possible to equip a BS with around 64 RF chains [144], but commercial BSs typically have 2-4 times more antenna elements than RF

chains [145].⁸ Viewed from the perspective of the enormous number of antenna elements envisioned for UM-MIMO, it is likely that the number of RF chains will never catch up with the desired number of antenna elements at a BS. If so, a fully digital architecture with one element per RF chain will essentially be relegated to benchmark status only, and hybrid architectures are the practical ones to consider. However, a hybrid architecture is not limited to using analog phase shifters; rather, with the term *hybrid* we refer to an architecture that somehow reduces the signal dimension between the antenna elements and the baseband processor.

It is essential to use the terms *fully digital* and *hybrid* judiciously. As discussed in Section III, the available spatial DoF, denoted as η_{2D} in (41), might be smaller than the number of deployed antenna elements. Consequently, a hybrid architecture can be lossless in terms of communication performance only if the number of RF chains is greater or equal to η_{2D} . Therefore, when referencing a fully digital architecture, it may implicitly refer to a hybrid one with η_{2D} RF chains.

To identify an immediate issue with having η_{2D} RF chains, consider a UM-MIMO array implemented on the jumbotron in a sports venue. Suppose the antenna elements cover an area A of the jumbotron and the carrier frequency is f_c . It follows from (41) that $\eta_{2D} = A \frac{f_c^2 \pi}{c^2}$. For a bandwidth of B Hz, this produces, at the very minimum, $\eta_{2D} B$ samples per second as output from the UM-MIMO array. If we represent each sample with b bits and forward all this information to a baseband unit (BBU), the rate at the input to the BBU is

$$R_{\text{BBU}} = ABb \frac{f_c^2 \pi}{c^2} \text{ bit/s.} \quad (143)$$

To appreciate the scale of this rate, we consider the case of $B = 100$ MHz, $f_c = 3$ GHz, $b = 16$, and $A = 10$ m². This results in $R_{\text{BBU}} \approx 5$ Tbit/s. At mmWave frequencies with 10 times more spectrum, the rate increases to $R_{\text{BBU}} \approx 5000$ Tbit/s. These rates are simply not feasible to implement in the foreseeable future; hence, the design of high-performing but lossy hybrid architectures constitutes a promising research direction.

The remainder of this section will describe ways to co-design antenna arrays with signal processing algorithms. We will describe a modular design of hybrid UM-MIMO antennas and elaborate on the tradeoff between the number of BBU inputs and the processing at the antenna elements.

A. Modular UM-MIMO design

We will now take a closer look at a UM-MIMO array and focus on its backplane; that is, the circuitry behind the antenna elements. Abstractly, we can view the UM-MIMO array as having several outputs, each connected to the BBU. When using a hybrid architecture, these outputs correspond to the RF chains. In a UM-MIMO setup, it is possible to define the density of these RF chains as μ m⁻². Additionally, each RF chain may be linked only to a subset of the antennas within

⁸Each RF chain is connected to a subarray with multiple antenna elements, to achieve a stronger vertical directivity than with a single element. Since no phase shifters are utilized, this implementation is still called fully digital.

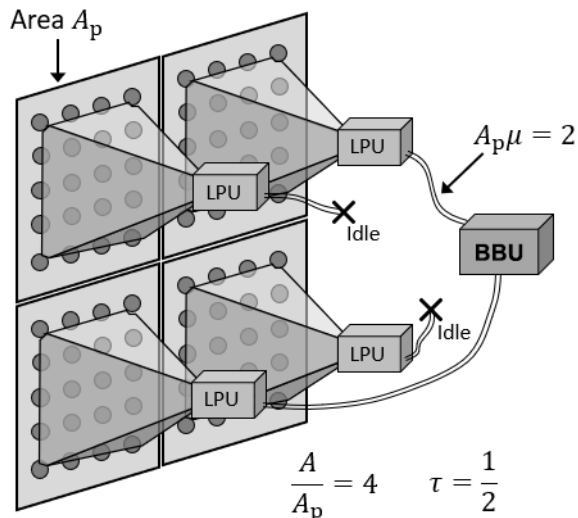


Fig. 23. Example of a modular antenna design where only some panels are connected to the CPU at a given time instance.

the array. We can then define an *antenna panel* as a group of antennas, comprising an area A_p m² to which a given number M_p of RF chains are connected; clearly, $M_p = A_p \mu$. Panel-based implementations of UM-MIMO antennas are thoroughly discussed in [146].

For a given rate limitation at the input interface to the BBU, only a subset of the panels may forward their outputs to the BBU at a given time instance while the remaining are idle. If we let τ denote the fraction of the total antenna area, A , that is active, the total number of RF chains connected to the BBU becomes

$$M = \frac{\tau A}{A_p} M_p = \frac{\tau A}{A_p} A_p \mu = A \tau \mu. \quad (144)$$

A basic example is provided in Fig. 23, where four panels, each of area A_p are equipped with $M_p = 2$ RF chains. We consider $\tau = 0.5$, which implies that at most $\tau A/A_p = 2$ panels may be active (i.e., forward their outputs to the BBU) at any given time. The total number of RF chains connected to the BBU at a given time becomes $M = M_p \tau A/A_p = 4$. The connection between each panel and the BBU is digital and managed by a local processing unit (LPU); thus, there are no physical switches in the array but only a control mechanism that determines which LPUs forward signals to the BBU.

The preferable selection of the parameters τ and μ depends on the deployment scenario. Choosing a very large μ (i.e., letting each panel have a large number of RF chains), implies that τ must be reduced to satisfy a given rate constraint on the BBU input. This may be favorable in case the UEs are believed to appear in hotspots (i.e., small subsets of the coverage area) so that they all have good propagation characteristics to the same panels. However, if the UEs are fairly uniformly distributed over the coverage area, then a small τ (i.e., only a few panels are active), may not be a good choice.

Ultimately, striking a good balance between τ and μ is not a problem that has an analytical solution, but rather requires extensive system simulations—possibly considering a digital twin of the intended deployment area. Let the metric(s)

of interest be denoted by \mathbf{c} , and possibly be vector-valued. This may be the ergodic channel capacity, outage capacity, or any other metric of choice. What quantities/settings do \mathbf{c} depend on? First of all, it strongly depends on the propagation environment and user distribution, which is why a numerical approach is required. Secondly, it depends on various fixed parameters such as the carrier frequency, bandwidth, available deployment area A of the array, etc. Thirdly, it depends on the type of signal processing performed at the LPU. This can range from purely analog beamforming (e.g., using phase-shifters and signal combiners) to a fully digital LPU that can refine/compress the complex baseband samples. Lastly, it depends on the quantities τ , μ , and A_p ; therefore, we express the metric as $\mathbf{c}(\tau, \mu, A_p)$ with the remark that all other properties discussed are held constant. Despite the fact that A_p does not show up in (144), it does impact the overall performance as it impacts the quality of the LPU outputs.

Rule-of-thumbs for selecting τ and μ in a line-of-sight-dominant propagation environment were obtained in [147] by extensive simulations. These show a linear relation between the number of BBU inputs and the number of UEs and active panels. Further work is required to determine the generality of these results.

The signal processing that is performed at each LPU to limit the flow of data to the BBU can also be optimized; for example, to reduce the dimensionality while retaining most of the performance. This problem is studied in [148], [149] with a star-topology between the BBU and LPU, while sequential topologies were considered in [150].

X. FINAL REMARKS AND OPEN CHALLENGES

Many new technology components are envisioned for 6G networks, but one that surely will play a key role is MIMO. An educated guess is that we will first see 6G MIMO in the upper mid-band (from 6-15 GHz), where 512 to 2048 antennas per array are within practical reach. The success of the technology is tightly connected to scientific and hardware maturity, and many grand challenges remain to be tackled. We will end this paper by describing eight such challenges:

- 1) **Beamfocusing in realistic environments:** We showed how near-field effects improve the spatial multiplexing capabilities in line-of-sight scenarios. How prominent are these effects in more realistic environments?
- 2) **Fixed-complexity BBU processing:** The interface between antennas and BBU will be limited in practice, as well as the computational capabilities. How can LPUs reduce the signal dimensionality to harness the benefits of having many antennas with a fixed-complexity BBU?
- 3) **Energy-efficient operation:** The energy consumption of UM-MIMO will likely be higher than that of 5G MIMO, but the energy per bit can be substantially smaller when massive multiplexing is performed. However, intelligent sleep features are required to achieve energy efficiency when the traffic load is small (e.g., at night).
- 4) **Channel estimation:** Channel coefficients must be estimated in every coherence block. Does this impose a fundamental limit on the useful spatial DoF? How does

the limit depend on the available side-information in the estimator?

- 5) **Array topology:** How to optimize the antenna number and spacing for a given area by considering mutual coupling and polarization effects? How to incorporate antenna reconfigurability and different kinds of array architectures?
- 6) **Distortion-aware processing:** Physical-consistent hardware and noise modeling, as well as RF impairments, make the end-to-end system model different than in textbooks. Can these “distortions” be overcome through digital processing?
- 7) **Continuous processing:** The CAP and DMA architectures enable nearly continuous spatial processing over the aperture, but under specific phase and amplitude constraints. How can advanced spatial multiplexing methods (e.g., zero-forcing) be implemented in these cases, where the matrix operations have infinite dimensions?
- 8) **Field trials:** Although near-field effects and their consequences can be simulated using EM-compliant models, true confidence in the technology is obtained through testbeds and measurements. UM-MIMO pushes the limits in terms of the cost and size of such field trials. Antenna designers and communication engineers must come together to overcome these challenges.

REFERENCES

- [1] E. F. W. Alexanderson, “Transatlantic radio communication,” *Trans. American Institute of Electrical Engineers*, vol. 38, no. 2, pp. 1269–1285, 1919.
- [2] P. Bondyopadhyay, “The first application of array antenna,” in *Proc. IEEE Int. Conf. Phased Array Systems and Tech.*, 2000, pp. 29–32.
- [3] H. O. Peterson, H. H. Beverage, and J. B. Moore, “Diversity telephone receiving system of R.C.A. communications, Inc.” *IRE*, vol. 19, no. 4, pp. 562–584, 1931.
- [4] H. T. Friis and C. B. Feldman, “A multiple unit steerable antenna for short-wave reception,” *IRE*, vol. 25, no. 7, pp. 841–917, 1937.
- [5] S. M. Alamouti, “A simple transmit diversity technique for wireless communications,” *IEEE J. Sel. Areas Commun.*, vol. 16, no. 8, pp. 1451–1458, 1998.
- [6] V. Tarokh, H. Jafarkhani, and A. Calderbank, “Space-time block codes from orthogonal designs,” *IEEE Trans. Inf. Theory*, vol. 45, no. 5, pp. 1456–1467, 1999.
- [7] B. M. Hochwald, T. L. Marzetta, and V. Tarokh, “Multiple-antenna channel hardening and its implications for rate feedback and scheduling,” *IEEE Trans. Inf. Theory*, vol. 60, no. 9, pp. 1893–1909, 2004.
- [8] J. H. Winters, “Optimum combining for indoor radio systems with multiple users,” *IEEE Trans. Commun.*, vol. 35, no. 11, pp. 1222–1230, 1987.
- [9] S. C. Swales, M. A. Beach, D. J. Edwards, and J. P. McGeehan, “The performance enhancement of multibeam adaptive base-station antennas for cellular land mobile radio systems,” *IEEE Trans. Veh. Technol.*, vol. 39, no. 1, pp. 56–67, 1990.
- [10] G. J. Foschini and M. J. Gans, “On limits of wireless communications in a fading environment when using multiple antennas,” *Wireless Personal Commun.*, vol. 6, no. 3, pp. 311–335, 1998.
- [11] E. Telatar, “Capacity of multi-antenna Gaussian channels,” *European Trans. Telecom.*, vol. 10, no. 6, pp. 585–595, 1999.
- [12] D. Gesbert, M. Kountouris, R. W. Heath, Jr., C.-B. Chae, and T. Salzer, “Shifting the MIMO paradigm,” *IEEE Signal. Process. Mag.*, vol. 24, no. 5, pp. 36–46, 2008.
- [13] T. L. Marzetta, “Noncooperative cellular wireless with unlimited numbers of base station antennas,” *IEEE Trans. Wireless Commun.*, vol. 9, no. 11, pp. 3590–3600, 2010.
- [14] E. Björnson, E. G. Larsson, and T. L. Marzetta, “Massive MIMO: Ten myths and one critical question,” *IEEE Commun. Mag.*, vol. 54, no. 2, pp. 114–123, Feb. 2016.
- [15] E. Björnson, J. Hoydis, and L. Sanguinetti, “Massive MIMO networks: Spectral, energy, and hardware efficiency,” *Foundations and Trends in Signal Processing*, vol. 11, no. 3–4, pp. 154–655, 2017.
- [16] Ericsson, “Ericsson mobility report,” Nov. 2023. [Online]. Available: <http://www.ericsson.com/mobility-report>
- [17] C. E. Shannon, “Communication in the presence of noise,” *IRE*, vol. 37, no. 1, pp. 10–21, 1949.
- [18] T. L. Marzetta, E. G. Larsson, H. Yang, and H. Q. Ngo, *Fundamentals of Massive MIMO*. Cambridge University Press, 2016.
- [19] E. Björnson, L. Sanguinetti, H. Wymeersch, J. Hoydis, and T. L. Marzetta, “Massive MIMO is a reality—What is next? Five promising research directions for antenna arrays,” *Digital Signal Processing*, vol. 94, pp. 3–20, 2019.
- [20] Z. Wang, J. Zhang, H. Du, W. E. I. Sha, B. Ai, D. Niyato, and M. Debbah, “Extremely large-scale MIMO: Fundamentals, challenges, solutions, and future directions,” *IEEE Wireless Commun.*, 2023.
- [21] M. Cui, Z. Wu, Y. Lu, X. Wei, and L. Dai, “Near-field MIMO communications for 6G: Fundamentals, challenges, potentials, and future directions,” *IEEE Commun. Mag.*, vol. 61, no. 1, pp. 40–46, 2022.
- [22] M. Uusitalo *et al.*, “Rfocus: Beamforming using thousands of passive antennas,” in *Proc. USenix symposium on networked systems design and implementation (NSDI)*, 2020, pp. 1047–1061.
- [23] M. Cui and L. Dai, “Channel estimation for extremely large-scale MIMO: Far-field or near-field?” *IEEE Trans. on Commun.*, vol. 70, no. 4, pp. 2663–2677, 2022.
- [24] D. Dardari, “Communicating with large intelligent surfaces: Fundamental limits and models,” *IEEE J. Sel. Areas Commun.*, vol. 38, no. 11, pp. 2526–2537, 2020.
- [25] H. T. Friis, “A note on a simple transmission formula,” *IRE*, vol. 34, no. 5, pp. 254–256, 1946.
- [26] K. T. Selvan and R. Janaswamy, “Fraunhofer and Fresnel distances: Unified derivation for aperture antennas,” *IEEE Antennas Propag. Mag.*, vol. 59, no. 4, pp. 12–15, 2017.
- [27] C. A. Balanis, *Antenna theory: Analysis and design*. John Wiley & Sons, 2015.
- [28] J. Sherman, “Properties of focused apertures in the Fresnel region,” *IRE Trans. Antennas Propag.*, vol. 10, no. 4, pp. 399–408, 1962.
- [29] E. Björnson and L. Sanguinetti, “Power scaling laws and near-field behaviors of massive MIMO and intelligent reflecting surfaces,” *IEEE Open J. Commun. Soc.*, vol. 1, pp. 1306–1324, 2020.
- [30] A. Kay, “Near-field gain of aperture antennas,” *IRE Trans. Antennas Propag.*, vol. 8, no. 6, pp. 586–593, 1960.
- [31] E. Björnson, Ö. T. Demir, and L. Sanguinetti, “A primer on near-field beamforming for arrays and reconfigurable intelligent surfaces,” in *Proc. Asilomar Conf. on Signals, Systems, and Computers*, 2021, pp. 105–112.
- [32] X. Gao, O. Edfors, F. Tufvesson, and E. G. Larsson, “Massive MIMO in real propagation environments: Do all antennas contribute equally?” *IEEE Trans. Commun.*, vol. 63, no. 11, pp. 3917–3928, 2015.
- [33] C. Polk, “Optical Fresnel-zone gain of a rectangular aperture,” *IRE Trans. Antennas Propag.*, vol. 4, no. 1, pp. 65–69, 1956.
- [34] T. Kwon, Y.-G. Lim, B. Min, and C.-B. Chae, “RF lens-embedded massive MIMO systems: Fabrication issues and codebook design,” *IEEE Trans. Microw. Theory Tech.*, vol. 64, no. 7, pp. 2256–2271, July 2016.
- [35] Y.-J. Cho, G.-Y. Suk, B. Kim, D.-K. Kim, and C.-B. Chae, “RF lens embedded antenna array for mmWave MIMO: Design and performance,” *IEEE Commun. Mag.*, vol. 56, no. 6, pp. 42–48, July 2018.
- [36] S.-H. Park, S.-M. Kim, S. Kim, H.-I. Yoo, B. Kim, and C.-B. Chae, “A transparent antenna system for in-building networks,” in *Proc. IEEE Int. Conf. Commun. (ICC)*, 2022.
- [37] F. Pirz, “Design of a wideband, constant beamwidth, array microphone for use in the near field,” *Bell Syst. Tech. J.*, vol. 58, no. 8, pp. 1839–1850, Oct. 1979.
- [38] J. Ryan and R. Goubran, “Near-field beamforming for microphone arrays,” in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.*, vol. 1, Munich, Germany, Apr. 1997, pp. 363–366.
- [39] S. Nordholm, V. Rehbock, K. Tee, and S. Nordebo, “Chebyshev optimization for the design of broadband beamformers in the near field,” *IEEE Trans. Circuits Syst. II: Analog Digit. Signal Process.*, vol. 45, no. 1, pp. 141–143, Jan. 1998.
- [40] N. J. Myers and R. W. Heath, “InFocus: A spatial coding technique to mitigate misfocus in near-field LoS beamforming,” *IEEE Trans. Wireless Commun.*, vol. 21, no. 4, pp. 2193–2209, Apr. 2022.

- [41] H. Lu and Y. Zeng, "Near-field modeling and performance analysis for multi-user extremely large-scale MIMO communication," *IEEE Commun. Lett.*, vol. 26, no. 2, pp. 277–281, 2021.
- [42] B. Friedlander, "Localization of signals in the near-field of an antenna array," *IEEE Trans. Signal Process.*, vol. 67, no. 15, pp. 3885–3893, Aug. 2019.
- [43] S. Phang, M. T. Ivrlač, G. Gradoni, S. C. Creagh, G. Tanner, and J. A. Nossek, "Near-field MIMO communication links," *IEEE Trans. Circuits Syst. I: Regul. Pap.*, vol. 65, no. 9, pp. 3027–3036, 2018.
- [44] N. Deshpande, M. R. Castellanos, S. R. Khosravirad, J. Du, H. Viswanathan, and R. W. Heath, "A wideband generalization of the near-field region for extremely large phased-arrays," *IEEE Wireless Commun. Lett.*, vol. 12, no. 3, pp. 515–519, Mar. 2023.
- [45] G. Bacci, L. Sanguinetti, and E. Björnson, "Spherical wavefronts improve MU-MIMO spectral efficiency when using electrically large arrays," *IEEE Wireless Commun. Lett.*, vol. 12, no. 7, pp. 1219–1223, 2023.
- [46] F. Bohagen, P. Orten, and G. E. Oien, "Design of optimal high-rank line-of-sight MIMO channels," *IEEE Trans. Wireless Commun.*, vol. 6, no. 4, pp. 1420–1425, 2007.
- [47] E. Torkildson, U. Madhow, and M. Rodwell, "Indoor millimeter wave MIMO: Feasibility and performance," *IEEE Trans. Wireless Commun.*, vol. 10, no. 12, pp. 4150–4160, 2011.
- [48] H. Do, N. Lee, and A. Lozano, "Reconfigurable ULAs for line-of-sight MIMO transmission," *IEEE Trans. Wireless Commun.*, vol. 20, no. 5, pp. 2933–2947, 2021.
- [49] D. Tse and P. Viswanath, *Fundamentals of wireless communications*. Cambridge University Press, 2005.
- [50] G. Caire and S. Shamai, "On the achievable throughput of a multi-antenna Gaussian broadcast channel," *IEEE Trans. Inf. Theory*, vol. 49, no. 7, pp. 1691–1706, 2003.
- [51] A. Goldsmith, S. A. Jafar, N. Jindal, and S. Vishwanath, "Capacity limits of MIMO channels," *IEEE J. Sel. Areas Commun.*, vol. 21, no. 5, pp. 684–702, 2003.
- [52] P. Viswanath and D. N. C. Tse, "Sum capacity of the vector Gaussian broadcast channel and uplink-downlink duality," *IEEE Trans. Inf. Theory*, vol. 49, no. 8, pp. 1912–1921, 2003.
- [53] O. Bucci and G. Franceschetti, "On the spatial bandwidth of scattered fields," *IEEE Trans. Antennas Propag.*, vol. 35, no. 12, pp. 1445–1455, 1987.
- [54] O. Bucci, C. Gennarelli, and C. Savarese, "Representation of electromagnetic fields over arbitrary surfaces by a finite and nonredundant number of samples," *IEEE Trans. Antennas Propag.*, vol. 46, no. 3, pp. 351–359, 1998.
- [55] A. M. Sayeed, "Deconstructing multiantenna fading channels," *IEEE Trans. Signal Process.*, vol. 50, no. 10, pp. 2563–2579, 2002.
- [56] A. S. Y. Poon, R. W. Brodersen, and D. N. C. Tse, "Degrees of freedom in multiple-antenna channels: A signal space approach," *IEEE Trans. Inf. Theory*, vol. 51, no. 2, pp. 523–536, 2005.
- [57] R. A. Kennedy, P. Sadeghi, T. D. Abhayapala, and H. M. Jones, "Intrinsic limits of dimensionality and richness in random multipath fields," *IEEE Trans. Signal Processing*, vol. 55, no. 6, pp. 2542–2556, 2007.
- [58] L. W. Hanlen and T. D. Abhayapala, "Space-time-frequency degrees of freedom: Fundamental limits for spatial information," in *Proc. IEEE Int. Symp. Inf. Theory*, 2007, pp. 701–705.
- [59] M. Franceschetti, *Wave Theory of Information*. Cambridge University Press, 2017.
- [60] S. Hu, F. Rusek, and O. Edfors, "Beyond massive MIMO: The potential of data transmission with large intelligent surfaces," *IEEE Trans. Signal Process.*, vol. 66, no. 10, pp. 2746–2758, 2018.
- [61] A. Pizzo, T. L. Marzetta, and L. Sanguinetti, "Degrees of freedom of holographic MIMO channels," in *IEEE Int. Workshop on Sig. Proc. Advances in Wireless Commun. (SPAWC)*, 2020, pp. 1–5.
- [62] A. Pizzo, A. d. J. Torres, L. Sanguinetti, and T. L. Marzetta, "Nyquist sampling and degrees of freedom of electromagnetic fields," *IEEE Trans. Signal Process.*, vol. 70, pp. 3935–3947, 2022.
- [63] A. Pizzo and A. Lozano, "On Landau's eigenvalue theorem for line-of-sight MIMO channels," *IEEE Wireless Commun. Lett.*, vol. 11, no. 12, pp. 2565–2569, 2022.
- [64] W. C. Chew, *Waves and fields in inhomogeneous media*. Wiley-IEEE Press, 1995.
- [65] S. M. Kay, *Fundamentals of statistical signal processing: Estimation theory*. Prentice Hall, 1993.
- [66] J. Kotecha and A. Sayeed, "Transmit signal design for optimal estimation of correlated MIMO channels," *IEEE Trans. Signal Process.*, vol. 52, no. 2, pp. 546–557, 2004.
- [67] Ö. T. Demir, E. Björnson, and L. Sanguinetti, "Channel modeling and channel estimation for holographic massive MIMO with planar arrays," *IEEE Wireless Commun. Lett.*, vol. 11, no. 5, pp. 997–1001, 2022.
- [68] Ö. T. Demir, E. Björnson, and L. Sanguinetti, "Exploiting array geometry for reduced-subspace channel estimation in RIS-aided communications," in *Proc. IEEE Sens. Array Multichannel Signal Process. Workshop (SAM)*, 2022, pp. 455–459.
- [69] Z. Gao, L. Dai, S. Han, I. Chih-Lin, Z. Wang, and L. Hanzo, "Compressive sensing techniques for next-generation wireless communications," *IEEE Wirel. Commun.*, vol. 25, no. 3, pp. 144–153, 2018.
- [70] C. You, Y. Zhang, C. Wu, Y. Zeng, B. Zheng, L. Chen, L. Dai, and A. L. Swindlehurst, "Near-field beam management for extremely large-scale array communications," *arXiv preprint arXiv:2306.16206*, 2023.
- [71] Ö. T. Demir and E. Björnson, "A new polar-domain dictionary design for the near-field region of extremely large aperture arrays," in *IEEE CAMSAP*, 2023.
- [72] Z. Wu and L. Dai, "Multiple access for near-field communications: SDMA or LDMA?" *IEEE J. Sel. Areas Commun.*, 2023.
- [73] S. S. Yuan, Z. He, X. Chen, C. Huang, and E. Wei, "Electromagnetic effective degree of freedom in a MIMO system in free space," *IEEE Antennas Wirel. Propag. Lett.*, vol. 21, no. 3, pp. 446–450, 2021.
- [74] S. Sun and M. Tao, "Characteristics of channel eigenvalues and mutual coupling effects for holographic reconfigurable intelligent surfaces," *Sensors*, vol. 22, no. 14, p. 5297, 2022.
- [75] T. L. Marzetta, E. G. Larsson, and T. B. Hansen, "Massive MIMO and Beyond" in *Information Theoretic Perspectives on 5G Systems and Beyond*, I. Maric, O. Simeone, S. Shamai (editors). Cambridge University Press, 2019.
- [76] M. T. Ivrlac and J. A. Nossek, "Toward a circuit theory of communication," *IEEE Trans. Circuits and Systems*, vol. 57, pp. 1663–1683, 2010.
- [77] D. C. Youla, *Theory and Synthesis of Linear Passive Time-Invariant Networks*. Cambridge University Press, 2015.
- [78] T. B. Hansen and A. D. Yaghjian, *Plane-Wave Theory of Time-Domain Fields: Near-Field Scanning Applications*. Wiley-IEEE Press, 1999.
- [79] H. Weyl, "Ausbreitung elektromagnetischer wellen über einen ebenen leiter," *Ann. Physik*, vol. 60, pp. 481–500, 1919.
- [80] T. L. Marzetta, "Fundamental limitations on the capacity of wireless links that use polarimetric antenna arrays," in *Proc. IEEE Int. Symp. Inf. Theory*, 2002, p. 51.
- [81] Y. Li, H. Wang, and X.-G. Xia, "On quasi-orthogonal space-time block codes for dual-polarized MIMO channels," *IEEE Trans. Wireless Commun.*, vol. 11, no. 1, pp. 397–407, Jan. 2012.
- [82] T. Kim, B. Clerckx, D. J. Love, and S. J. Kim, "Limited feedback beamforming systems for dual-polarized MIMO channels," *IEEE Trans. Wireless Commun.*, vol. 9, no. 11, pp. 3425–3439, Nov. 2010.
- [83] N. Amitay and J. Salz, "Linear equalization theory in digital data transmission over dually polarized fading radio channels," *Bell Labs Tech. J.*, vol. 63, no. 10, pp. 2215–2259, 1984.
- [84] S. A. Schelkunoff, "A mathematical theory of linear arrays," *The Bell System Technical Journal*, vol. 22, 1943.
- [85] A. Uzkof, "An approach to the problem of optimum directive antenna design," *Comptes Rendus (Doklady) de l'Academie des Sci. de l'URSS*, vol. 53, pp. 35–38, 1946.
- [86] T. L. Marzetta, "Super-directive antenna arrays: Fundamentals and new perspectives," *Proc. Asilomar Conf. on Signals, Systems, and Computers*, 2019.
- [87] A. Pizzo, T. L. Marzetta, and L. Sanguinetti, "Spatially-stationary model for holographic MIMO small-scale fading," *IEEE J. Sel. Areas Commun.*, vol. 38, pp. 1964–1979, 2020.
- [88] R. H. Clarke, "A statistical theory of mobile-radio reception," *Bell System Technical Journal*, pp. 957–1000, 1968.
- [89] T. L. Marzetta and T. B. Hansen, "Rayleigh-Jeans-Clarke model for wireless noise in a resonant cavity: Scalar case," *Proc. IEEE Global Commun. Conf. (GLOBECOM)*, 2022.
- [90] P. M. Morse and H. Feshbach, *Methods of Theoretical Physics*. McGraw-Hill, 1953.
- [91] G. Joos, *Theoretical Physics*. Glasgow: I. M. Freeman, Blackie & Son Ltd, 1934.
- [92] A. Singh and T. L. Marzetta, "Shannon theory for wireless communication in a resonant chamber," *accepted to IEEE J. Sel. Areas Commun. Special Issue on Electromagnetic Info. Theory*, 2023.
- [93] A. Mezghani, M. Akrouf, M. R. Castellanos, S. Saab, B. Hochwald, R. W. Heath, and J. A. Nossek, "Reincorporating circuit theory into information theory," *IEEE BITS Inf. Theory Mag.*, pp. 1–17, Dec. 2023.

- [94] D. M. Kerns and E. S. Dayhoff, "Theory of diffraction in microwave interferometry," *Journal of Research of the National Bureau of Standards Section B Mathematics and Mathematical Physics*, p. 1, 1960.
- [95] H. C. Pocklington, "Electrical oscillations in wire," *Cambridge Philos. Soc. Proc.*, vol. 9, pp. 324–332, 1897.
- [96] W. Wasylikiwskij and W. Kahn, "Theory of mutual coupling among minimum-scattering antennas," *IEEE Trans. Antennas Propag.*, vol. 18, no. 2, pp. 204–216, 1970.
- [97] —, "Scattering properties and mutual coupling of antennas with prescribed radiation pattern," *IEEE Transactions on Antennas and Propagation*, vol. 18, no. 6, pp. 741–752, 1970.
- [98] A. Mezghani and R. W. Heath, Jr., "Massive MIMO precoding and spectral shaping with low resolution phase-only DACs and active constellation extension," *IEEE Trans. Wireless Commun.*, vol. 21, no. 7, pp. 5265–5278, 2022.
- [99] H. Nyquist, "Thermal agitation of electric charge in conductors," *Physical review*, vol. 32, no. 1, p. 110, 1928.
- [100] M. R. Castellanos and R. W. Heath, "Electromagnetic manifold characterization of antenna arrays," submitted to *IEEE Trans. Wireless Commun.*, Oct. 2023.
- [101] R. Bhagavatula, C. Oestges, and R. W. Heath, "A new double-directional channel model including antenna patterns, array orientation, and depolarization," *IEEE Trans. Veh. Technol.*, vol. 59, no. 5, pp. 2219–2231, Jun. 2010.
- [102] M. R. Castellanos and R. W. Heath, "Linear polarization optimization for wideband MIMO systems with reconfigurable arrays," *IEEE Trans. Wireless Commun.*, pp. 1–14, Nov. 2023.
- [103] B. Friedlander, "Polarization sensitivity of antenna arrays," *IEEE Trans. Signal Process.*, vol. 67, no. 1, pp. 234–244, Jan. 2019.
- [104] G. Efstathopoulos and A. Manikas, "Extended array manifolds: Functions of array manifolds," *IEEE Trans. Signal Process.*, vol. 59, no. 7, pp. 3272–3287, Jul. 2011.
- [105] B. Friedlander, "The extended manifold for antenna arrays," *IEEE Trans. Signal Process.*, vol. 68, pp. 493–502, Jan. 2020.
- [106] W. L. Stutzman and G. A. Thiele, *Antenna theory and design*. John Wiley & Sons, 2012.
- [107] A. Buffi, P. Nepa, and G. Manara, "Design criteria for near-field-focused planar arrays," *IEEE Antennas Propag. Mag.*, vol. 54, no. 1, pp. 40–50, 2012.
- [108] L. Wei, C. Huang, G. C. Alexandropoulos, W. E. I. Sha, Z. Zhang, M. Debbah, and C. Yuen, "Multi-user holographic MIMO surfaces: Channel modeling and spectral efficiency analysis," *IEEE J. Sel. Top. Signal Process.*, vol. 16, no. 5, pp. 1112–1124, 2022.
- [109] H. Lu and Y. Zeng, "Communicating with extremely large-scale array/surface: Unified modeling and performance analysis," *IEEE Trans. Wireless Commun.*, vol. 21, no. 6, pp. 4039–4053, 2021.
- [110] X. Zhang, A. F. Molisch, and S.-Y. Kung, "Variable-phase-shift-based RF-baseband codesign for MIMO antenna selection," *IEEE Trans. Signal Process.*, vol. 53, no. 11, pp. 4091–4103, 2005.
- [111] J. de Souza, A. Amiri, T. Abrao, E. de Carvalho, and P. Popovski, "Quasi-distributed antenna selection for spectral efficiency maximization in subarray switching XL-MIMO systems," *IEEE Trans. Veh. Technol.*, vol. 70, no. 7, pp. 6713–6725, 2021.
- [112] F. Sohrabi and W. Yu, "Hybrid digital and analog beamforming design for large-scale antenna arrays," *IEEE J. Sel. Top. Signal Process.*, vol. 10, no. 3, pp. 501–513, 2016.
- [113] S. S. Ioushua and Y. C. Eldar, "A family of hybrid analog–digital beamforming methods for massive MIMO systems," *IEEE Trans. Signal Process.*, vol. 67, no. 12, pp. 3243–3257, 2019.
- [114] K.-K. Wong, A. Shojaefard, K. F. Tong, and Y. Zhang, "Fluid antenna systems," *IEEE Trans. Wireless Commun.*, vol. 20, no. 3, pp. 1950–1962, 2020.
- [115] L. Sanguinetti, A. A. D. Amico, and M. Debbah, "Wavenumber-division multiplexing in line-of-sight holographic MIMO communications," *IEEE Trans. Wireless Commun.*, 2022.
- [116] S. Hu, F. Rusek, and O. Edfors, "Beyond massive MIMO: The potential of positioning with large intelligent surfaces," *IEEE Trans. Signal Process.*, vol. 66, no. 7, pp. 1761–1774, 2018.
- [117] J. Yang, Y. Zeng, S. Jin, C.-K. Wen, and P. Xu, "Communication and localization with extremely large lens antenna array," *IEEE Trans. Wireless Commun.*, vol. 20, no. 5, pp. 3031–3048, 2021.
- [118] Z. Zhang and L. Dai, "Pattern-division multiplexing for multi-user continuous-aperture MIMO," *IEEE J. Sel. Areas Commun.*, 2023.
- [119] Z. Wan, J. Zhu, and L. Dai, "Performance comparison between continuous aperture MIMO and discrete MIMO," in *Proc. IEEE Global Commun. Conf. (GLOBECOM)*, 2023.
- [120] K.-K. Wong, W. K. New, X. Hao, K.-F. Tong, and C.-B. Chae, "Fluid antenna system—part I: Preliminaries," *IEEE Commun. Letters*, vol. 27, no. 8, pp. 1919–1923, Aug. 2023.
- [121] K.-K. Wong, K.-F. Tong, and C.-B. Chae, "Fluid antenna system—part II: Research opportunities," *IEEE Commun. Letters*, vol. 27, no. 8, pp. 1924–1928, Aug. 2023.
- [122] K.-K. Wong, K.-F. Tong, and C.-B. Chae, "Fluid antenna system—part III: A new paradigm of distributed artificial scattering surfaces for massive connectivity," *IEEE Commun. Letters*, vol. 27, no. 8, pp. 1929–1933, Aug. 2023.
- [123] G. J. Hayes, J. Ju-Hee So, A. Qusba, M. D. Dickey, and G. Lazzi, "Flexible liquid metal alloy (EGaIn) microstrip patch antenna," *IEEE Trans. Antennas Propag.*, vol. 60, no. 5, pp. 2151–2156, 2012.
- [124] A. Dey, R. Guldiken, and G. Mumcu, "Microfluidically reconfigured wideband frequency-tunable liquid-metal monopole antenna," *IEEE Trans. Antennas Propag.*, vol. 64, no. 6, pp. 2572–2576, 2016.
- [125] A. Singh, I. Goode, and C. E. Saavedra, "A multistate frequency reconfigurable monopole antenna using fluidic channels," *IEEE Antennas Wirel. Propag. Lett.*, vol. 18, no. 5, pp. 856–860, 2019.
- [126] K.-K. Wong, K. F. Tong, Y. Shen, Y. Chen, and Y. Zhang, "Bruce lee-inspired fluid antenna system: Six research topics and the potentials for 6G," *Frontiers in Commun. Net.*, vol. 3, pp. 1–32, 2022.
- [127] W. K. New, K.-K. Wong, H. Xu, K. F. Tong, and C.-B. Chae, "Fluid antenna system: New insights on outage probability and diversity gain," *IEEE Trans. Wireless Commun.*, 2023.
- [128] J. M. Carlson, M. R. Castellanos, and R. W. Heath, Jr., "Dynamic metasurface antennas for energy-efficient MISO communications," in *Proc. IEEE Global Commun. Conf. (IEEE GLOBECOM)*, 2023.
- [129] J. D. Boerman and J. T. Bernhard, "Performance study of pattern reconfigurable antennas in MIMO communication systems," *IEEE Trans. Antennas Propag.*, vol. 56, no. 1, pp. 231–236, Jan. 2008.
- [130] D. Piazza, N. J. Kirsch, A. Forenza, R. W. Heath, and K. R. Dandekar, "Design and evaluation of a reconfigurable antenna array for MIMO systems," *IEEE Trans. Antennas Propag.*, vol. 56, no. 3, pp. 869–881, Mar. 2008.
- [131] A. M. Sayeed and V. Raghavan, "Maximizing MIMO capacity in sparse multipath with reconfigurable antenna arrays," *IEEE J. Sel. Areas Commun.*, vol. 1, no. 1, pp. 156–166, Jun. 2007.
- [132] M. Hasan, İ. Bahçeci, and B. A. Cetiner, "Downlink multi-user MIMO transmission for radiation pattern reconfigurable antenna systems," *IEEE Trans. Wireless Commun.*, vol. 17, no. 10, pp. 6448–6463, Oct. 2018.
- [133] H. Wang, Y. Shen, K. F. Tong, and K.-K. Wong, "Continuous electrowetting surface-wave fluid antenna for mobile communications," in *Proc. IEEE Region 10 Conference (TENCON)*, 2022, pp. 1–3.
- [134] P. K. Nisar *et al.*, "Liquid metal antennas: Materials, fabrication and applications," *Sensors*, vol. 20, no. 1, p. 177, 2019.
- [135] J. Wang *et al.*, "Metantenna: When metasurface meets antenna again," *IEEE Trans. Antennas Propag.*, vol. 68, no. 3, pp. 1332–1347, 2020.
- [136] T. J. Cui, S. Liu, and L. L. Li, "Information entropy of coding metasurface," *Light: Science & Applications*, vol. 5, no. 11, 2016.
- [137] M. Barbuto *et al.*, "Metasurfaces 3.0: a new paradigm for enabling smart electromagnetic environments," *IEEE Trans. Antennas Propag.*, vol. 70, no. 10, pp. 8883–8897, 2021.
- [138] R. L. Haupt and M. Lanagan, "Reconfigurable antennas," *IEEE Antennas Propag. Mag.*, vol. 55, no. 1, pp. 49–61, 2013.
- [139] D. R. Smith, O. Yurduseven, L. P. Mancera, P. Bowen, and N. B. Kundtz, "Analysis of a waveguide-fed metasurface antenna," *Phys. Rev. Appl.*, vol. 8, no. 5, p. 054048, 2017.
- [140] L. You, J. Xu, G. C. Alexandropoulos, J. Wang, W. Wang, and X. Gao, "Energy efficiency maximization of massive MIMO communications with dynamic metasurface antennas," *IEEE Trans. Wireless Commun.*, vol. 22, no. 1, pp. 393–407, 2022.
- [141] W. Huang *et al.*, "Joint microstrip selection and beamforming design for mmwave systems with dynamic metasurface antennas," in *Proc. IEEE Int. Conf. Acoust., Speech, Sig. Proc. (ICASSP)*, 2023, pp. 1–5.
- [142] H. Zhang *et al.*, "Beam focusing for near-field multiuser MIMO communications," *IEEE Trans. Wireless Commun.*, vol. 21, no. 9, pp. 7476–7490, 2022.
- [143] M. R. Castellanos, J. Carlson, and R. W. Heath, "Energy-efficient tri-hybrid precoding with dynamic metasurface antennas," in *Proc. Asilomar Conf. Signals, Syst., Comput.*, 2023.
- [144] J. Vieira *et al.*, "A flexible 100-antenna testbed for massive MIMO," in *Proc. IEEE Global Commun. Conf. (GLOBECOM)*, 2014.
- [145] H. Asplund *et al.*, *Advanced antenna systems for 5G network deployments: Bridging the gap between theory and practice*. Academic Press, 2020.

- [146] G. Callebaut, L. Liu, T. Eriksson, L. Van der Perre, O. Edfors, and C. Fager, "6G radio testbeds: Requirements, trends, and approaches," *IEEE Microwave Magazine*, 2024.
- [147] A. Pereira *et al.*, "Deployment strategies for large intelligent surfaces," *IEEE Access*, vol. 10, pp. 61 753–61 768, 2022.
- [148] J. Alegría, F. Rusek, and O. Edfors, "Trade-offs in decentralized multi-antenna architectures: The WAX decomposition," *IEEE Trans. Signal Proc.*, vol. 69, pp. 3627–3641, 2021.
- [149] J. Alegría and F. Rusek, "Trade-offs in decentralized multi-antenna architectures: Sparse combining modules for WAX decomposition," *IEEE Trans. Signal Proc. (to appear)*, vol. 71, 2023.
- [150] Z. H. Shaik, E. Björnson, and E. G. Larsson, "MMSE-optimal sequential processing for cell-free massive MIMO with radio stripes," *IEEE Trans. Commun.*, vol. 69, no. 11, pp. 7775–7789, 2021.