# Stacked Intelligent Metasurfaces for Efficient Holographic MIMO Communications in 6G

Jiancheng An, *Member, IEEE,* Chao Xu, *Senior Member, IEEE,* Derrick Wing Kwan Ng, *Fellow, IEEE,*
George C. Alexandropoulos, *Senior Member, IEEE,* Chongwen Huang, *Member, IEEE,* Chau Yuen, *Fellow, IEEE,*
and Lajos Hanzo, *Life Fellow, IEEE*

*Abstract*—The revolutionary technology of *Stacked Intelligent Metasurfaces (SIM)* has been recently shown to be capable of carrying out advanced signal processing directly in the native electromagnetic (EM) wave domain. An SIM is fabricated by a sophisticated amalgam of multiple stacked metasurface layers, which may outperform its single-layer metasurface counterparts, such as reconfigurable intelligent surfaces (RISd) and metasurface lenses. We harness this new SIM concept for implementing efficient holographic multiple-input multiple-output (HMIMO) communications that dot require excessive radio-frequency (RF) chains, which constitutes a substantial benefit compared to existing implementations. We first present an HMIMO communication system based on a pair of SIMs at the transmitter (TX) and receiver (RX), respectively. In sharp contrast to the conventional MIMO designs, the considered SIMs are capable of automatically accomplishing transmit precoding and receiver combining, as the EM waves propagate through them. As such, each information data stream can be directly radiated and recovered from the corresponding transmit and receive ports. Secondly, we formulate the problem of minimizing the error between the actual end-to-end SIMs'parametrized channel matrix and the target diagonal one, with the latter representing a flawless interference-free system of parallel subchannels. This is achieved by jointly optimizing the phase shifts associated with all the metasurface layers of both the TX-SIM and RX-SIM. We then design a gradient descent algorithm to solve the resultant non-convex problem. Furthermore, we theoretically analyze the HMIMO channel capacity bound and provide some useful fundamental insights. Extensive simulation results are provided for characterizing our SIM-based HMIMO system, quantifying its substantial performance benefits. Indicatively, it is demonstrated that a $150\%$ capacity improvement is feasible when compared with MIMO and RIS-aided communication systems.

*Index Terms*—Stacked intelligent metasurface (SIM), holographic MIMO (HMIMO), reconfigurable intelligent surface (RIS), 3D integrated metasurfaces, wave-based computing.

## I. INTRODUCTION

**W**ITH the completion of the 3GPP Release 17, it is high time for both industry and academia to begin conceptualizing the sixth-generation (6G) mobile networks [1]. Wireless network evolution has been primarily motivated by the pursuit of higher data rates and wider device connectivity. While this demand will continue to increase, the explosive proliferation of the Internet-of-Everything (IoE), ranging from extended reality to interconnected autonomous systems, is driving a fundamental paradigm shift [2]. It is envisaged that by 2030, the number of connected devices will reach 500 billion, according to Cisco's annual report [3]. To support these heterogeneous IoE services imposing extreme performance requirements, 6G wireless networks are expected to integrate communication, sensing, computing, and control capabilities, while drastically improving data rates, latency, and connectivity [4]. As such, 6G will undergo a revolutionary transformation by flexibly orchestrating both physical and virtual resources to support the envisioned heterogeneous IoE scenarios and by harnessing sophisticated disruptive techniques, including artificial intelligence (AI) [5] spanning all network layers [6], satellite communications [7], programmable reflective [8] and computing metasurfaces [9], and integrated sensing and communications devices and techniques [10], [11], to name a few.

## A. Emerging Metasurface-Based Technologies

We commence by reviewing a pair of 6G enabling techniques, namely, reconfigurable intelligent surfaces (RISs) and holographic multiple-input multiple-output (HMIMO) communications.

*1) Reconfigurable Intelligent Surfaces (RISs):* The emerging metasurface technology was shown to be able to shape smart reconfigurable environments [8], [12]–[16]. Specifically, a programmable metasurface is composed of a large number of low-cost passive reflecting elements, which are capable of manipulating the electromagnetic (EM) behavior of radio waves [17], allowing for proactive customization of the wireless propagation environment [18], [19]. By employing an RIS, [20] designed a scheme that substantially boosted the energy efficiency of the downlink in multiuser multiple-input single-output (MISO) communication systems, as compared to conventional relaying solutions. Following this pioneering work, a great deal of research has sprouted up by investigating the effects of realistic hardware imperfections [21] and RIS element responses [22], optimizing the RIS phase shifts for wideband operation [23] and based on the statistical channel state information (CSI) [24], [25], and enhancing the quality-of-service (QoS) [26]–[28] in RIS-assisted communication systems. However, the encouraging performance benefits of these RIS solutions rely on the availability of accurate CSI, which generally requires excessive pilot overhead for acquisition [16]. To address this issue, the authors in [18], [29] devised codebook-based frameworks for channel estimation in RIS-assisted MIMO systems. According to those frameworks, the estimation of the RIS-parametrized channel and the design of transmit beamforming were handled via conventional protocols, while simplifying the reflection coefficient optimization by selecting the best entry from otpimized RIS reflection beam codebooks [30], [31]. It has been shown that the designed codebook-based solutions are appealingly scalable exhibiting strong robustness against hardware imperfections [18], [30].

*2) Holographic MIMO (HMIMO):* Over the past decade, the massive MIMO technique has become one of the most crucial enablers for increasing wireless capacity [32]. Explicitly, massive MIMO has the potential of focusing energy into a smaller spatial region, thus attaining huge improvements in both spectral and energy efficiencies [33]. It also provides other benefits, such as it can be realized with inexpensive low-power components, exhibit reduced latency, lead to simplified protocol designs, and be robustness against jamming [32]. As the 6G research is ramping up, a natural question arises – *what will the next generation MIMO be like?* Recently, the innovative concept of HMIMO has emerged [34]–[37]. Specifically, by employing a large intelligent surface (LIS) constructed of an electromagnetically active material that integrates massive numbers of radiating and sensing elements [38], impressive improvements are expected [39]. Furthermore, [40] demonstrated that LIS-aided solutions are capable of improving the spatial multiplexing gain even in strong line-of-sight (LoS) propagation conditions. To characterize the fundamental capacity limits of HMIMO communications, the authors in [41] established a mathematically tractable channel model by considering the small-scale fading in the far-field as a spatially correlated random Gaussian field, while being consistent with the scalar Helmholtz equation. The same authors then developed a Fourier plane-wave series-based expansion of the HMIMO channel response for arbitrary scattering environments [42]. Following this channel model, a channel estimation scheme leveraging the specific array geometry for identifying a low-dimensional common subspace for arbitrary spatial correlation matrices was designed in [43]. Moreover, the authors in [44] studied the family of discrete amplitude-controlled holographic beamformers with the specific objective of satisfying a given sum-rate requirement, while multi-user HMIMO systems were investigated in [45], [46].

## B. Motivation

It has been recently proposed to cascade multiple metasurfaces[1] to realize *stacked intelligent metasurfaces (SIMs)*, which have the capability to implement signal processing in the EM wave regime [47], [48], [52]. In this paper, we propose the integration of SIMs with the transceivers to support HMIMO communications. Before proceeding, we first elaborate on our motivation by answering the question – *why do we need SIM?* – from the following three perspectives:

1) The existing research on HMIMO communications is still in its infancy and lacks practical implementations [37]. Since integrating an abundance of expensive active elements at the transceiver is an impractical option, recent research efforts focus on implementing HMIMO communications by employing programmable metasurfaces [34]–[36]. However, their performance remains limited by practical hardware constraints, such as the tunable amplitude/phase associated with each meta-atom of a single-layer metasurface. As a remedy, a multilayer metasurface architecture might be beneficial for improving both the spatial-domain gain and the design degrees of freedom, thus, flexibly forming diverse radio frequency (RF) waveforms compared to its single-layer counterparts.

2) Although integrating RISs into existing wireless networks has been numerically shown to improve both the spectral and energy efficiencies in various scenarios [20], [21], [53], there are still stumbling blocks in the way of practical deployments of RISs. On one hand, the two-hop multiplicative path loss coefficient severely impacts the resultant performance [54]. Indeed, several works investigate the RIS placement to reduce pathloss [21] or optimize a certain performance objective [55]. As a further impediment, the widespread RIS deployment significantly increases the burden of media access control (MAC) optimization, including both resource allocation and multiple access [6], [56], [57]. The joint optimization of coexisting distributed active and passive nodes in wireless networks also increases the computational burden and control signaling overhead [58]. Hence, we

---

[1]A range of other terminologies having a similar multilayer structure were also used in different research communities, such as programmable AI machine [47], stacked metasurface slab [48], 3D integrated metasurface device [49], cascaded metasurfaces [50], and multilayer metasurfaces [51].

embark on investigating HMIMO communications by intrinsically integrating programmable metasurfaces with the transceiver.

3) Over the past decade, we have witnessed the rise of deep learning (DL) techniques. Although DL has been widely utilized for improving the performance of wireless networks, the implementation of DL relies essentially on a central processing unit (CPU) or a graphical processing unit (GPU) [59]. Explicitly, DL is merely a computing architecture, whose processing speed is fundamentally constrained by the specific CPU/GPU utilized. To further improve computational efficiency and reduce power consumption, a novel *wave-based computing* paradigm has recently enjoyed much research attention [60]. Specifically, by constructing a diffractive neural network having a well-designed multilayer structure, the computational tasks can be performed on the profile of the EM wave by leveraging the amplitude/phase information [61]. By fully harnessing the benefits of this wave-based computing paradigm, it becomes possible to perform massively parallel signal processing in the native EM wave regime, where the forward propagation within the diffractive neural network can be realized at the speed of light.

### C. Contributions

Motivated by the aforementioned observations, in this paper, we present an SIM-enabled HMIMO communication system. The main contributions of this paper are summarized as follows:

1) We establish a novel HMIMO framework by harnessing an SIM at the transmitter (TX) and another one at the receiver (RX) for achieving substantial spatial gains. The proposed SIM-based transceivers can directly perform precoding/combining in the native EM wave regime with reduced the number of transmit/receive RF chains, which is attributed to the large metasurface aperture.

2) We formulate a channel fitting problem which focuses on approximating an end-to-end diagonal channel matrix by optimizing the phase shifts associated with the different metasurface layers. This allows each spatial stream to be radiated and recovered independently at the corresponding transmit and receive ports, thus effectively creating a set of interference-free parallel subchannels. By taking into account the constant-modulus constraint and the coupled variables in the objective function, we then propose an efficient gradient descent algorithm for iteratively solving the resultant non-convex problem.

3) We theoretically analyze the HMIMO channel capacity. Since it is non-trivial to derive the closed-form capacity expression, we provide both an upper and lower bound of the HMIMO capacity by assuming that all the spatial streams experience the best and worst sub-channel quality, respectively. Furthermore, we provide fundamental insights into the scaling law of the HMIMO channel capacity versus the number of data streams and meta-atoms.

4) Our extensive results demonstrate the benefits of the SIM-aided HMIMO framework conceived as well as

the accuracy of our analytical results. We also quantify the channel fitting performance as well as the channel capacity attained under various setups, shedding light on the optimal SIM design. Additionally, we verify the substantial performance improvements attained compared to the conventional MIMO schemes as well as to their RIS-aided counterparts.

### D. Organization

The rest of this paper is organized as follows. In Section II, we introduce the general SIM-based HMIMO system model, which encompasses the SIM structure together with a spatially correlated HMIMO channel model. Following this, we formulate a channel fitting problem in Section III-A and propose an efficient algorithm to address the resulting optimization problem in Section III-B. We analyze the HMIMO channel capacity and the computational complexity of the proposed algorithm in Section IV. Finally, our numerical results are provided in Section V before concluding the paper in Section VI.

### E. Notations

We adopt bold lowercase and uppercase letters to denote vectors and matrices, respectively; $(\cdot)^*$, $(\cdot)^T$, and $(\cdot)^H$ represent the conjugate, transpose, and Hermitian transpose, respectively; $|c|$, $\Re(c)$, and $\Im(c)$ refer to the magnitude, real part, and imaginary part, respectively, of a complex number $c$; $\|\cdot\|_{\mathrm{F}}$ is the Frobenius norm; $\mathbb{E}(\cdot)$ stands for the expectation operation; $\log_a(\cdot)$ is the logarithmic function with base $a$, while $\ln(\cdot)$ is the natural logarithm; $\mathrm{diag}(\mathbf{v})$ produces a diagonal matrix with the elements of $\mathbf{v}$ on the main diagonal; $\mathbf{S}^{1/2}$ denotes the square root of a square matrix $\mathbf{S}$; $\mathrm{vec}(\mathbf{M})$ denotes the vectorization of a matrix $\mathbf{M}$; $\mathbf{M}_{a:b,\,:}$, $\mathbf{M}_{:,\,c:d}$, $\mathbf{M}_{a:b,\,c:d}$ represent the matrices constructed by extracting $a$-to-$b$-th rows, $c$-to-$d$-th columns, as well as both $a$-to-$b$-th rows and $c$-to-$d$-th columns, respectively, of the matrix $\mathbf{M}$; $\mathrm{sinc}(x) = \sin(\pi x)/(\pi x)$ is the sinc function; $\mathbb{C}^{x \times y}$ represents the space of $x \times y$ complex-valued matrices; $\lceil x \rceil$ refers to the nearest integer greater than or equal to $x$; $\mathrm{mod}(x, n)$ returns the remainder after division of $x$ by $n$; $\partial f/\partial x$ means the partial derivative of a function $f$ with respect to (*w.r.t.*) the variable $x$; $\mathbf{0}$ and $\mathbf{1}$ denote all-zero and all-one vectors, respectively, with appropriate dimensions, while $\mathbf{I}_N \in \mathbb{C}^{N \times N}$ denotes the identity matrix; the distribution of a circularly symmetric complex Gaussian (CSCG) random vector with mean vector $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma} \succeq \mathbf{0}$ is denoted by $\sim \mathcal{CN}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, where $\sim$ stands for "distributed as".

## II. THE PROPOSED SIM-BASED HMIMO SYSTEM MODEL

In this section, we present the holistic system model of our SIM-based HMIMO. Specifically, we first introduce the proposed TX/RX SIM-based design and then elaborate on the spatially correlated HMIMO channel model based on the recent consolidated efforts in [43]. Finally, we discuss the optimal transmission regime of the HMIMO channel, given a limited number of information data streams.
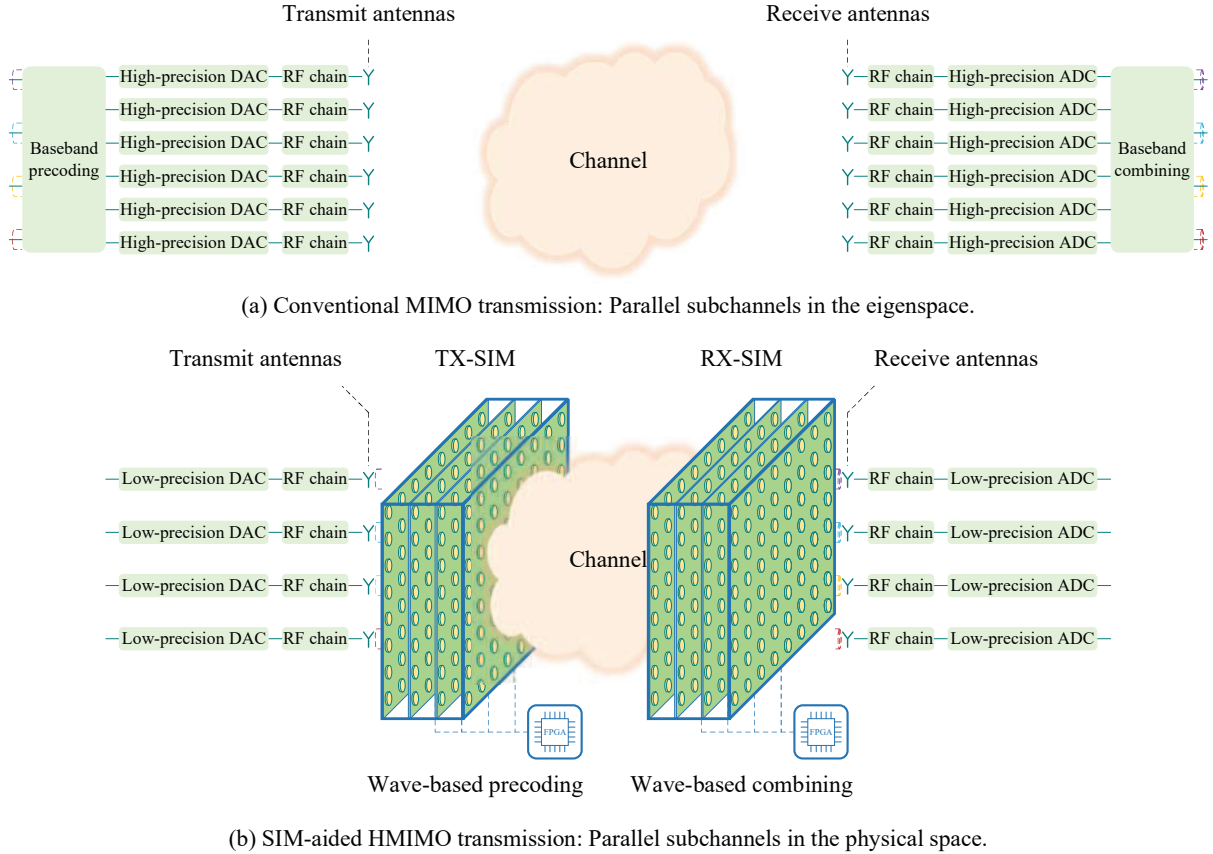
Transmit antennas · Receive antennas · Baseband precoding · High-precision DAC · RF chain · Channel · RF chain · High-precision ADC · Baseband combining

(a) Conventional MIMO transmission: Parallel subchannels in the eigenspace.

Transmit antennas · TX-SIM · RX-SIM · Receive antennas · Low-precision DAC · RF chain · Channel · RF chain · Low-precision ADC · Wave-based precoding · Wave-based combining

(b) SIM-aided HMIMO transmission: Parallel subchannels in the physical space.

Fig. 1. Transmission comparison of conventional MIMO and SIM-based HMIMO.

### A. Proposed SIM Design

Before proceeding, let us briefly review the conventional MIMO scheme illustrated in Fig. 1(a), where multiple data streams are first precoded and then fed to the corresponding transmit antennas. At the output of the wireless channel, receiver combining is adopted for recovering the different spatial streams. As a consequence, multiple parallel subchannels are constructed in the eigenspace domain [62], benefiting from the precoding and combining at the transmitter and receiver, respectively.

The proposed SIM-assisted HMIMO system is illustrated in Fig. 1(b). In sharp contrast to the conventional MIMO systems having only active antennas, an SIM is integrated with both the TX and RX for enhancing the QoS[2]. Specifically, a TX/RX-SIM is a closed vacuum container having several stacked metasurface layers [47]. Each metasurface is comprised of a large number of meta-atoms [63], which are connected to a smart controller, e.g., a field programmable gate array (FPGA) board. By appropriately tuning the drive level of the control circuit associated with each meta-atom, the system becomes capable of manipulating the EM behavior of the penetrating wave, and thus, producing a customized spatial waveform

shape at the output of the metasurface layer. Moreover, by compactly arranging large numbers of meta-atoms on the output metasurface of the TX-SIM as well as on the input metasurface of the RX-SIM, the desired information-bearing EM waves can be radiated from almost the entire surface into the ether and then collected in the same way. As such, both the TX-SIM and RX-SIM interact with the wireless channel in an almost continuous manner, thus, being capable to support low-latency *HMIMO communications*. Specifically, the TX-SIM undertakes the precoding task, casting the appropriate information-bearing EM wave into the ether, while the RX-SIM efficiently combines the impinging EM wave on the input surface for signal recovery. As a result, we are able to establish multiple parallel subchannels in the physical space, and the corresponding multiple data streams can be directly radiated and recovered from the associated TX and RX metasurfaces, respectively, without imposing any interference.

*Remark 1:* Here we elaborate on three core benefits of the proposed SIM-based HMIMO transmission paradigm, as compared to its conventional counterpart [62]. Firstly, the conventional MIMO design requires a large number of active components to achieve spatial gains, thus resulting in high hardware costs and energy consumption. In contrast, the proposed SIM-based HMIMO utilizes low-cost metasurfaces for gleaning spatial gains, which substantially reduces the number of active RF chains required. Secondly, the conventional MIMO transmission solution relies on high-precision

---

[2]Although, in this paper, we only consider the point-to-point HMIMO scenario, our SIM can also be utilized to perform the zero-forcing (ZF) precoding and combining for supporting multiuser HMIMO communications [45], [46]. The specific design is beyond the scope of this paper and reserved for our future research.
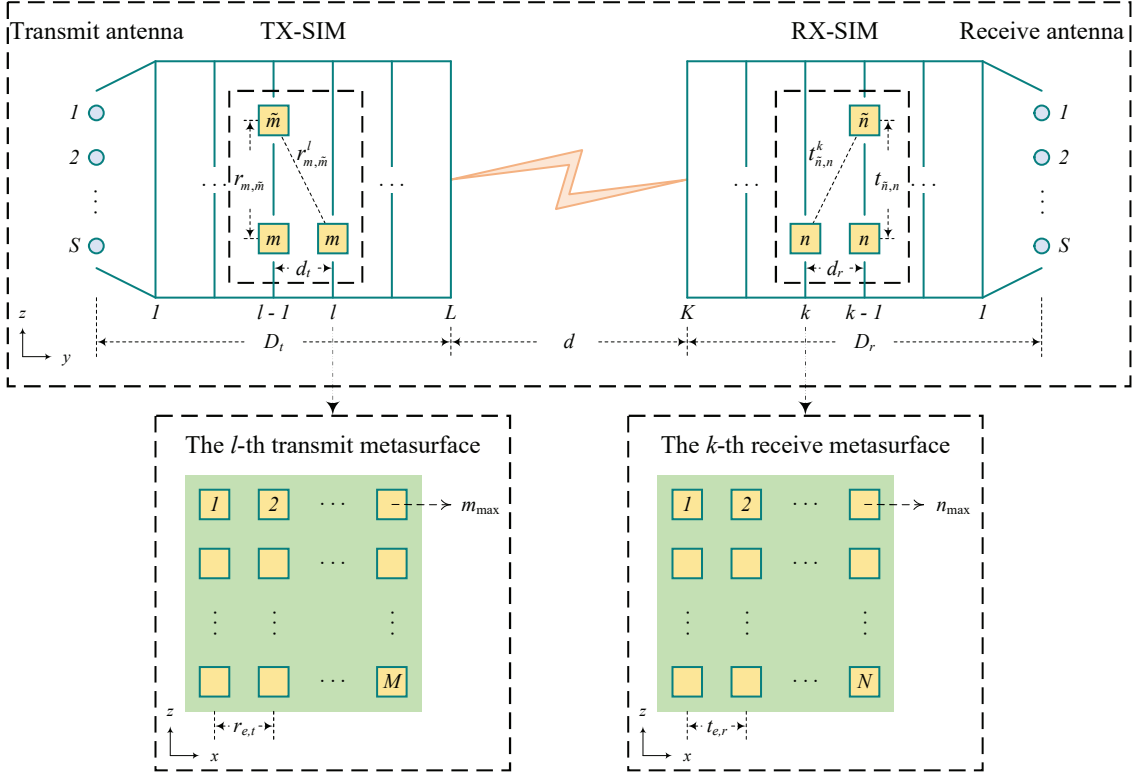
Fig. 2. Detailed schematic of the SIM-aided HMIMO system of Fig. 1(b).

power-thirsty digital-to-analog converters (DAC) and analog-to-digital converters (ADC). Instead, the TX/RX-SIM creates multiple parallel subchannels in the physical space, enabling each data stream to be individually processed by low-precision power-efficient DACs/ADCs. For example, 1-bit resolution may be used for binary phase shift keying (BPSK) without unduly compromising its communication performance. Thirdly, due to using precoding and combining in the wave domain, the power consumption of signal processing is significantly reduced. As such, the SIM substantially reduces the overall energy consumption compared to conventional digital transceiver designs [47]. Nevertheless, a quantitative evaluation of the energy efficiency of the proposed SIM relying on passive metasurfaces requires an accurate energy consumption model, as well as an accurate transmission model for characterizing the wave propagation between adjacent metasurfaces, both of which require further investigation.

A detailed schematic of the proposed SIM-based HMIMO system, which relies on wave-based precoding and combining, is provided in Fig. 2. Let $S$ and $\mathcal{S} = \{1, 2, \cdots, S\}$ denote the number of data streams and the corresponding set, respectively. Moreover, $L$ and $K$ denote the number of metasurface layers at the TX and RX, respectively, while their corresponding sets are represented by $\mathcal{L} = \{1, 2, \cdots, L\}$ and $\mathcal{K} = \{1, 2, \cdots, K\}$. For notational convenience, we assume that the number of meta-atoms on each metasurface layer of the TX-SIM is identical and so is for the RX-SIM of Fig. 2. Specifically, let $M$ and $N$ denote the number of meta-atoms on each metasurface layer associated with the

TX-SIM and the RX-SIM, respectively, satisfying $M \geq S$ and $N \geq S$, while representing the corresponding set as $\mathcal{M} = \{1, 2, \cdots, M\}$ and $\mathcal{N} = \{1, 2, \cdots, N\}$. Moreover, let $\phi_m^l = e^{j\theta_m^l}$ denote the transmission coefficient imposed by the $m$-th meta-atom on the $l$-th TX metasurface layer with $\theta_m^l$ representing the corresponding phase shift, which we assume that it can be continuously adjusted in the interval between $0$ and $2\pi$, i.e., $\theta_m^l \in [0, 2\pi)$, $m \in \mathcal{M}$, $l \in \mathcal{L}$. Thus, the transmission coefficient vector of the $l$-th TX metasurface layer and its corresponding matrix version are denoted by $\boldsymbol{\phi}^l = \left[\phi_1^l, \phi_2^l, \cdots, \phi_M^l\right]^T \in \mathbb{C}^{M \times 1}$ and $\boldsymbol{\Phi}^l = \text{diag}\left(\boldsymbol{\phi}^l\right) \in \mathbb{C}^{M \times M}$, respectively. Similarly, let $\psi_n^k = e^{j\xi_n^k}$ denote the transmission coefficient imposed by the $n$-th meta-atom on the $k$-th RX metasurface layer, where $\xi_n^k$ represents the corresponding phase shift satisfying $\xi_n^k \in [0, 2\pi)$, $n \in \mathcal{N}$, $k \in \mathcal{K}$. Then, the transmission coefficient vector of the $k$-th RX metasurface layer and its corresponding matrix version are respectively denoted by $\boldsymbol{\psi}^k = \left[\psi_1^k, \psi_2^k, \cdots, \psi_N^k\right]^T \in \mathbb{C}^{N \times 1}$ and $\boldsymbol{\Psi}^k = \text{diag}\left(\boldsymbol{\psi}^k\right) \in \mathbb{C}^{N \times N}$.

Furthermore, we assume that all the metasurface layers rely on an isomorphic lattice arrangement [47], while each metasurface is modeled as a uniform planar array. Specifically, the element spacing between the $\tilde{m}$-th meta-atom and the $m$-th one on the same TX metasurface and that between the $n$-th meta-atom and the $\tilde{n}$-th one on the same RX metasurface can be expressed as

$$r_{m,\tilde{m}} = r_{e,t}\sqrt{(m_z - \tilde{m}_z)^2 + (m_x - \tilde{m}_x)^2}, \qquad (1)$$

$$r_{m,s}^1 = \sqrt{\left[\left(m_z - \frac{m_{\max}+1}{2}\right)r_{e,t} - \left(s - \frac{S+1}{2}\right)\frac{\lambda}{2}\right]^2 + \left(m_x - \frac{m_{\max}+1}{2}\right)^2 r_{e,t}^2 + d_t^2}, \tag{7}$$

$$t_{s,n}^1 = \sqrt{\left[\left(n_z - \frac{n_{\max}+1}{2}\right)t_{e,r} - \left(s - \frac{S+1}{2}\right)\frac{\lambda}{2}\right]^2 + \left(n_x - \frac{n_{\max}+1}{2}\right)^2 t_{e,r}^2 + d_r^2}. \tag{8}$$

$$t_{\tilde{n},n} = t_{e,r}\sqrt{(\tilde{n}_z - n_z)^2 + (\tilde{n}_x - n_x)^2}, \tag{2}$$

respectively, where $r_{e,t}$ and $t_{e,r}$ denote the element spacing between adjacent meta-atoms on the same TX metasurface and that on the same RX metasurface, respectively (see Fig. 2). Additionally, $m_z$ and $m_x$ denote the indices of the $m$-th meta-atom along the $z$-axis and the $x$-axis, respectively, while $n_z$ and $n_x$ denote the indices of the $n$-th meta-atom along the $z$-axis and the $x$-axis, respectively, which are defined by

$$m_z = \lceil m/m_{\max} \rceil, \quad m_x = \mathrm{mod}\,(m-1, m_{\max}) + 1, \tag{3}$$
$$n_z = \lceil n/n_{\max} \rceil, \quad n_x = \mathrm{mod}\,(n-1, n_{\max}) + 1, \tag{4}$$

respectively, with $m_{\max}$ and $n_{\max}$ denoting the maximum number of meta-atoms on each row of the TX metasurface and that of the RX metasurface, respectively, as shown in Fig. 2. Throughout this paper, we consider square metasurface arrays associated with $M = m_{\max}^2$ and $N = n_{\max}^2$.

Let us now consider the transmission process between the adjacent metasurface layers. For the sake of simplicity, we assume that all the metasurface layers are parallel and have uniform spacing, as shown in Fig. 2. Specifically, let $d_t$ and $d_r$ denote the spacing between any two adjacent metasurface layers in the TX-SIM and that in the RX-SIM, respectively, while $D_t$ and $D_r$ represent the thickness of the TX-SIM and RX-SIM, respectively. Thus we have $d_t = D_t/L$ and $d_r = D_r/K$. As a result, the transmission distance from the $\tilde{m}$-th meta-atom on the $(l-1)$-st TX metasurface to the $m$-th one on the $l$-th TX metasurface and that from the $n$-th meta-atom of the $k$-th RX metasurface to the $\tilde{n}$-th one on the $(k-1)$-st RX metasurface are

$$r_{m,\tilde{m}}^l = \sqrt{r_{m,\tilde{m}}^2 + d_t^2}, \ l \in \mathcal{L}/\{1\}, \tag{5}$$
$$t_{\tilde{n},n}^k = \sqrt{t_{\tilde{n},n}^2 + d_r^2}, \ k \in \mathcal{K}/\{1\}, \tag{6}$$

respectively.

Next, we consider the transmission process from the transmit antenna array to the input metasurface of the TX-SIM and that from the output metasurface of the RX-SIM to the receive antenna array. The transmit and receive antennas are both arranged in a uniform linear array, with the element spacing of half-wavelength, i.e., $\lambda/2$, and the array centers aligned with those of all metasurfaces. By doing so, the transmission distance from the $s$-th source to the $m$-th meta-atom on the input metasurface of the TX-SIM and that from the $n$-th meta-atom on the output metasurface of the RX-SIM to the $s$-th destination is given by (7) and (8), respectively, as shown at the top of this page.

According to the Rayleigh-Sommerfeld diffraction theory [60], the transmission coefficient from the $\tilde{m}$-th meta-atom on the $(l-1)$-st TX metasurface layer to the $m$-th meta-atom on the $l$-th TX metasurface layer is expressed by

$$w_{m,\tilde{m}}^l = \frac{A_t \cos \chi_{m,\tilde{m}}^l}{r_{m,\tilde{m}}^l}\left(\frac{1}{2\pi r_{m,\tilde{m}}^l} - j\frac{1}{\lambda}\right)e^{j2\pi r_{m,\tilde{m}}^l/\lambda}, \ l \in \mathcal{L}, \tag{9}$$

where $r_{m,\tilde{m}}^l$ denotes the corresponding transmission distance, $A_t$ is the area of each meta-atom in the TX-SIM, while $\chi_{m,\tilde{m}}^l$ represents the angle between the propagation direction and the normal direction of the $(l-1)$-th TX metasurface layer. Let $\mathbf{W}^l \in \mathbb{C}^{M \times M}$, $l \in \mathcal{L}/\{1\}$ denote the transmission coefficient matrix between the $(l-1)$-st TX metasurface layer and the $l$-th TX metasurface layer. In particular, the transmission coefficient matrix from the transmit antenna array to the input metasurface layer of the TX-SIM is represented by $\mathbf{W}^1 \in \mathbb{C}^{M \times S}$. Thus, the effect of the TX-SIM in Fig. 2 is formulated by

$$\mathbf{P} = \mathbf{\Phi}^L \mathbf{W}^L \cdots \mathbf{\Phi}^2 \mathbf{W}^2 \mathbf{\Phi}^1 \mathbf{W}^1 \in \mathbb{C}^{M \times S}. \tag{10}$$

Moreover, the transmission coefficient from the $n$-th meta-atom on the $k$-th RX metasurface layer to the $\tilde{n}$-th meta-atom on the $(k-1)$-st RX metasurface layer is expressed by [60]

$$u_{\tilde{n},n}^k = \frac{A_r \cos \varsigma_{\tilde{n},n}^k}{t_{\tilde{n},n}^k}\left(\frac{1}{2\pi t_{\tilde{n},n}^k} - j\frac{1}{\lambda}\right)e^{j2\pi t_{\tilde{n},n}^k/\lambda}, \ k \in \mathcal{K}, \tag{11}$$

where $t_{\tilde{n},n}^k$ denotes the corresponding transmission distance, $A_r$ is the area of each meta-atom in the RX-SIM, while $\varsigma_{\tilde{n},n}^k$ represents the angle between the propagation direction and the normal direction of the $(k-1)$-th RX metasurface layer. Let $\mathbf{U}^k \in \mathbb{C}^{N \times N}$, $k \in \mathcal{K}/\{1\}$ represent the transmission coefficient matrix between the $k$-th RX metasurface layer to the $(k-1)$-st RX metasurface layer, while the transmission coefficient matrix from the output metasurface layer of the RX-SIM to the receive antenna array is denoted by $\mathbf{U}^1 \in \mathbb{C}^{S \times N}$. Hence, the effect of the RX-SIM in Fig. 2 is represented by

$$\mathbf{Q} = \mathbf{U}^1 \mathbf{\Psi}^1 \mathbf{U}^2 \mathbf{\Psi}^2 \cdots \mathbf{U}^K \mathbf{\Psi}^K \in \mathbb{C}^{S \times N}. \tag{12}$$

*Remark 2:* In order to maximize the energy efficiency and characterize the performance of SIM-aided HMIMO communications, we have applied the constant modulus constraint and assumed continuously-adjustable phase shifts for the transmission coefficients associated with each meta-atom [64]. While the tuning precision in practice is typically proportional to the hardware costs, the low-resolution meta-atoms may be used in practical SIM design. The specific SIM optimization and the resultant performance analysis taking into account these hardware constraints such as realistic discrete phase shifts

[47], adjustable magnitudes [13], as well as coupled phase and magnitude tuning mechanisms [8] deserve future exploration.

*Remark 3:* It should be noted that the transmission coefficients between adjacent metasurface layers may deviate from those specified in (9) and (11) due to the existence of practical hardware imperfections, irreversible fabrication shortcomings, as well as innate modeling errors [47]. Hence, it may be necessary to calibrate these transmission coefficients before the SIM's practical deployment, which can be carried out separately for each individual SIM. One effective method is to transmit a known excitation signal and measure the response at the receive panel, and then update the transmission coefficients by employing the classic error back-propagation algorithm [59]. As an initial exploration of the precoding and combining capability of SIM, the specific calibrate process is beyond the scope of this paper and reserved for our future research.

### B. Spatially-Correlated HMIMO Channel Model

Next, let us consider the HMIMO channel model between the TX-SIM and the RX-SIM, where the most prominent property is the spatial correlation among the tightly packed meta-atoms. Specifically, the spatially-correlated HMIMO channel spanning from the output metasurface of the TX-SIM to the input metasurface of the RX-SIM is written by [44]

$$\mathbf{G} = \mathbf{R}_{\text{Rx}}^{1/2} \tilde{\mathbf{G}} \mathbf{R}_{\text{Tx}}^{1/2} \in \mathbb{C}^{N \times M}, \quad (13)$$

where $\tilde{\mathbf{G}} \in \mathbb{C}^{N \times M}$ is the independent and identically distributed (i.i.d.) Rayleigh fading channel, i.e., $\tilde{\mathbf{G}} \sim \mathcal{CN}\left(\mathbf{0}, \rho^2 \mathbf{I}_N \otimes \mathbf{I}_M\right)$ with $\rho^2$ denoting the average path loss between the pair of wireless transceivers, while $\mathbf{R}_{\text{Tx}} \in \mathbb{C}^{M \times M}$ and $\mathbf{R}_{\text{Rx}} \in \mathbb{C}^{N \times N}$ represent the spatial correlation matrix at the TX-SIM and that at the RX-SIM, respectively. By considering far-field propagation in an isotropic scattering environment [41], [65], the spatial correlation matrix at the TX-SIM and that at the RX-SIM can be expressed by [43]

$$[\mathbf{R}_{\text{Tx}}]_{m,\tilde{m}} = \text{sinc}\left(2r_{m,\tilde{m}}/\lambda\right), \ \tilde{m} \in \mathcal{M}, \ m \in \mathcal{M}, \quad (14)$$

$$[\mathbf{R}_{\text{Rx}}]_{\tilde{n},n} = \text{sinc}\left(2t_{\tilde{n},n}/\lambda\right), \ n \in \mathcal{N}, \ \tilde{n} \in \mathcal{N}, \quad (15)$$

respectively.

Moreover, the path loss between the transmitter and the receiver is modeled by [66]

$$\text{PL}\left(d\right) = \text{PL}\left(d_0\right) + 10b \log_{10}\left(\frac{d}{d_0}\right) + X_\delta, \ d \geq d_0, \quad (16)$$

where $\text{PL}\left(d_0\right) = 20 \log_{10}\left(4\pi d_0/\lambda\right)$ dB is the free space path loss at the reference distance $d_0$, $b$ represents the path loss exponent, $X_\delta$ is a zero mean Gaussian random variable with a standard deviation $\delta$, characterizing the large-scale signal fluctuations of shadow fading.

*Remark 4:* Note that the spatial correlation matrix highly depends on the scattering environments surrounding both the transmitter and the receiver. Therefore, it is generally unrealistic to derive a universal spatially-correlated fading model that can be applied to all practical communication scenarios. Motivated readers are referred to [41]–[43] and references therein for gaining further insights concerning the channel models for HMIMO systems that are derived from the intrinsic EM propagation properties.

### C. SIM-Aided HMIMO Channel Capacity with Limited Number of Streams

In this subsection, we will consider both the HMIMO channel capacity as well as the optimal transmission given a limited number of data streams. Specifically, given the wireless channel $\mathbf{G}$ and the fixed number of data streams $S$, the optimal HMIMO transmission is achieved by applying the truncated singular value decomposition (SVD) policy [18]. The detailed procedures are outlined as follows:

*1:* First of all, we perform the SVD of $\mathbf{G}$ so that $\mathbf{G} = \mathbf{E}\boldsymbol{\Lambda}\mathbf{F}^H$, where we have $\boldsymbol{\Lambda} = \text{diag}\left(\lambda_1, \lambda_2, \cdots, \lambda_O\right)$ and $O = \min\left(M, N\right)$, while $\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_O$ denoting the singular values in non-increasing order.

*2:* Next, by spreading the data streams using a transmit precoder $\mathbf{P} = \mathbf{F}_{:,1:S} \in \mathbb{C}^{M \times S}$ and collecting the spatial signals using a receive combiner $\mathbf{Q} = \mathbf{E}_{:,1:S}^H \in \mathbb{C}^{S \times N}$, the resultant end-to-end channel becomes the following diagonal matrix

$$\mathbf{H} = \mathbf{Q}\mathbf{G}\mathbf{P} = \boldsymbol{\Lambda}_{1:S,1:S} \in \mathbb{C}^{S \times S}, \quad (17)$$

with $\boldsymbol{\Lambda}_{1:S,1:S}$ being the $S$-th order leading principal minor of $\boldsymbol{\Lambda}$.

*3:* Furthermore, in order to maximize the channel capacity, the optimal power allocation coefficients can be obtained by applying the well-known water-filling algorithm [28]. Specifically, the amount of power allocated to the $s$-th data stream is determined as

$$p_s = \max\left(0, \tau - \frac{\sigma^2}{\lambda_s^2}\right), \ s \in \mathcal{S}, \quad (18)$$

where $\tau$ is a threshold value satisfying the total transmit power constraint, i.e., $\sum_{s=1}^{S} p_s = P_t$ with $P_t$ denoting the total available power at the transmitter, which can be obtained by utilizing the bisection method, while $\sigma^2$ represents the average noise power at the receiver.

*4:* Therefore, the HMIMO channel capacity for a finite number of data streams is given by

$$C = \sum_{s=1}^{S} \log_2\left(1 + \frac{p_s \lambda_s^2}{\sigma^2}\right). \quad (19)$$

Next, let us get back to our SIM-aided HMIMO communication system shown in Fig. 2. In sharp contrast to conventional MIMO designs constructing multiple virtual streams in the eigenspace by employing the appropriate digital precoding and combining, we endeavor to naturally form multiple parallel physical subchannels between the transmit antennas and their corresponding receive antennas. Explicitly, the precoding and combining are implemented in the EM wave regime by optimizing the TX-SIM and RX-SIM as follows

$$\mathbf{F}_{:,1:S} \simeq \boldsymbol{\Phi}^L \mathbf{W}^L \cdots \boldsymbol{\Phi}^2 \mathbf{W}^2 \boldsymbol{\Phi}^1 \mathbf{W}^1, \quad (20)$$

$$\mathbf{E}_{:,1:S}^H \simeq \mathbf{U}^1 \boldsymbol{\Psi}^1 \mathbf{U}^2 \boldsymbol{\Psi}^2 \cdots \mathbf{U}^K \boldsymbol{\Psi}^K. \quad (21)$$

Thus, one might be able to construct an end-to-end diagonal channel $\mathbf{H}$, such that multiple data streams can be directly radiated and recovered, respectively, from the corresponding transmit and receive antennas. In a nutshell, SIM not only reaps spatial gains benefiting from the massive number of

meta-atoms on the metasurface layer but performs the precoding and combining at the speed of light, thanks to its direct wave-based computing capability attained by the multilayer structure.

## III. PROBLEM FORMULATION AND SOLUTION OF JOINT OPTIMIZING TX-SIM AND RX-SIM

### A. Problem Formulation

In this subsection, we formulate the problem of minimizing the error between the end-to-end channel of $\mathbf{H} = \mathbf{QGP}$ and the expected diagonal matrix $\mathbf{\Lambda}_{1:S,1:S}$ by optimizing the phase shifts of the TX-SIM and RX-SIM in Fig. 2. We adopt the Frobenius norm to characterize the fitting error of the desired channel fitting problem [67]. Specifically, the optimization problem is formulated as[3]

$$\underset{\phi_m^l, \psi_n^k, \alpha}{\text{minimize}} \ \Gamma = \|\alpha\mathbf{QGP} - \mathbf{\Lambda}_{1:S,1:S}\|_{\text{F}}^2 \tag{22a}$$

$$\text{subject to} \quad \mathbf{P} = \mathbf{\Phi}^L\mathbf{W}^L\cdots\mathbf{\Phi}^2\mathbf{W}^2\mathbf{\Phi}^1\mathbf{W}^1, \tag{22b}$$

$$\mathbf{Q} = \mathbf{U}^1\mathbf{\Psi}^1\mathbf{U}^2\mathbf{\Psi}^2\cdots\mathbf{U}^K\mathbf{\Psi}^K, \tag{22c}$$

$$\mathbf{\Phi}^l = \text{diag}\left(\left[\phi_1^l, \phi_2^l, \cdots, \phi_M^l\right]^T\right), \ l \in \mathcal{L}, \tag{22d}$$

$$\mathbf{\Psi}^k = \text{diag}\left(\left[\psi_1^k, \psi_2^k, \cdots, \psi_N^k\right]^T\right), \ k \in \mathcal{K}, \tag{22e}$$

$$\left|\phi_m^l\right| = 1, \ m \in \mathcal{M}, \ l \in \mathcal{L}, \tag{22f}$$

$$\left|\psi_n^k\right| = 1, \ n \in \mathcal{N}, \ k \in \mathcal{K}, \tag{22g}$$

$$\alpha \in \mathbb{C}, \tag{22h}$$

where $\alpha$ is a scaling factor compensated by SIM.

*Remark 5:* Our objective is to evaluate the signal processing capability of SIM having multiple stacked metasurface layers. As such, we have assumed that the TX-SIM and RX-SIM of Fig. 2 compensate for an adaptive gain $\alpha$, as seen in (22a). Although this assumption might seem somewhat simplified, it is reasonable due to the fact that conventional precoding and combining architectures relying on digital signal processing result in additional hardware costs and energy consumption. Thus, a fair performance comparison between these two paradigms requires special attention under the same resource consumption. Since the energy consumption of our proposed SIM is unexplored, it poses an open challenge in conducting a fair performance comparison.

Note that due to the non-convex constant modulus constraint on each transmission coefficient, i.e., (22f) and (22g), as well as the highly coupled variables in the objective function, i.e., (22a), it is non-trivial to obtain the optimal solution of Problem (22). As such, in Section III-B, we will provide an efficient algorithm for solving the channel fitting problem iteratively.

### B. The Proposed Gradient Descent Algorithm

In this section, an efficient gradient descent algorithm is proposed for solving the challenging Problem (22). To ensure

[3]Note that although the formulated problem seems to have a similar form to that in hybrid beamforming, e.g., [67], they are fundamentally different because the proposed SIM-assisted HMIMO fully removes the digital precoding and combining by utilizing the multilayer metasurface structure, while imposing the phase shifts in the wave regime.

compliance with the constant modulus constraints, i.e., (22f) and (22g), our gradient descent algorithm is implemented by deriving the partial derivative *w.r.t.* the phase shifts. As such, the constant modulus constraints can be guaranteed throughout the iteration process. The detailed steps of the iteration core of the proposed gradient descent algorithm are summarized as follows.

#### Step 1: Calculate the partial derivatives
First, the partial derivatives of the loss function $\Gamma$ *w.r.t.* the $m$-th phase shift $\theta_m^l$ of the $l$-th TX metasurface layer and that *w.r.t.* the $n$-th phase shift $\xi_n^k$ of the $k$-th RX metasurface layer are respectively given by

$$\frac{\partial\Gamma}{\partial\theta_m^l} = 2\sum_{s=1}^{S}\sum_{\tilde{s}=1}^{S}\Im\left[\left(\alpha\phi_m^l x_{m,s,\tilde{s}}^l\right)^*\left(\alpha h_{s,\tilde{s}} - \lambda_{s,\tilde{s}}\right)\right], \tag{23}$$

$$\frac{\partial\Gamma}{\partial\xi_n^k} = 2\sum_{s=1}^{S}\sum_{\tilde{s}=1}^{S}\Im\left[\left(\alpha\psi_n^k y_{n,s,\tilde{s}}^k\right)^*\left(\alpha h_{s,\tilde{s}} - \lambda_{s,\tilde{s}}\right)\right], \tag{24}$$

where $h_{s,\tilde{s}}$ and $\lambda_{s,\tilde{s}}$ denote the entries located at the $s$-th row and the $\tilde{s}$-th column of the matrix $\mathbf{H} = \mathbf{QGP}$ and that of the matrix $\mathbf{\Lambda}$, respectively; while $x_{m,s,\tilde{s}}^l$ and $y_{n,s,\tilde{s}}^k$ denote the cascaded channel spanning from the $\tilde{s}$-th transmit antenna to the $s$-th receive antenna via the $m$-th meta-atom of the $l$-th TX metasurface layer and that via the $n$-th meta-atom of the $k$-th RX metasurface layer, which are defined by

$$x_{m,s,\tilde{s}}^l = \mathbf{Q}_{s,:}\mathbf{G}\mathbf{\Phi}^L\mathbf{W}^L\cdots\mathbf{W}_{:,m}^{l+1}\mathbf{W}_{m,:}^l\cdots\mathbf{\Phi}^1\mathbf{W}_{:,\tilde{s}}^1, \tag{25}$$

$$y_{n,s,\tilde{s}}^k = \mathbf{U}_{s,:}^1\mathbf{\Psi}^1\cdots\mathbf{U}_{:,n}^k\mathbf{U}_{n,:}^{k+1}\cdots\mathbf{U}^K\mathbf{\Psi}^K\mathbf{G}\mathbf{P}_{:,\tilde{s}}, \tag{26}$$

respectively.

#### Step 2: Normalize the partial derivatives
In order to mitigate the potential gradient explosion and vanishing problems [68], we normalize the partial derivatives at each iteration as follows

$$\frac{\partial\Gamma}{\partial\theta_m^l} \leftarrow \frac{\pi}{\varrho_l}\cdot\frac{\partial\Gamma}{\partial\theta_m^l}, \ m \in \mathcal{M}, \ l \in \mathcal{L}, \tag{27}$$

$$\frac{\partial\Gamma}{\partial\xi_n^k} \leftarrow \frac{\pi}{\varepsilon_k}\cdot\frac{\partial\Gamma}{\partial\xi_n^k}, \ n \in \mathcal{N}, \ k \in \mathcal{K}, \tag{28}$$

where we have $\varrho_l = \max_{m\in\mathcal{M}}\left(\frac{\partial\Gamma}{\partial\theta_m^l}\right), l \in \mathcal{L}$ and $\varepsilon_k = \max_{n\in\mathcal{N}}\left(\frac{\partial\Gamma}{\partial\xi_n^k}\right), k \in \mathcal{K}$ denoting the maximum value of the partial derivative associated with the $l$-th TX metasurface layer and that with the $k$-th RX metasurface layer, respectively. Note that the normalization process also has the benefit of allowing us to readily pick an initial learning rate independent of the data value [68].

#### Step 3: Update the phase shifts
Then the phase shifts associated with the TX-SIM and RX-SIM in Fig. 2 can be updated by

$$\theta_m^l \leftarrow \theta_m^l - \eta\frac{\partial\Gamma}{\partial\theta_m^l}, \ m \in \mathcal{M}, \ l \in \mathcal{L}, \tag{29}$$

$$\xi_n^k \leftarrow \xi_n^k - \eta\frac{\partial\Gamma}{\partial\xi_n^k}, \ n \in \mathcal{N}, \ k \in \mathcal{K}, \tag{30}$$

respectively, where $\eta > 0$ denotes the learning rate that determines the step size at each iteration.

#### Step 4: Update the scaling factor

TABLE I
THE PROPOSED GRADIENT DESCENT ALGORITHM FOR SOLVING (22).

1: **INPUT:** $\mathbf{W}^l$, $l \in \mathcal{L}$, $\mathbf{G}$, $\mathbf{U}^k$, $k \in \mathcal{K}$, $\mathbf{\Lambda}_{1:S,1:S}$.
2: Randomly initializing the phase shifts, i.e., $\boldsymbol{\theta}^l$, $l \in \mathcal{L}$, $\boldsymbol{\xi}^k$, $k \in \mathcal{K}$;
3: Calculating the scaling factor $\alpha$ by applying (31);
4: **REPEAT**
5:   Calculating the partial derivatives of $\Gamma$ *w.r.t.* $\theta_m^l$ and that *w.r.t.* $\xi_n^k$ by applying (23) and (24), respectively;
6:   Normalizing the partial derivatives of $\Gamma$ *w.r.t.* $\theta_m^l$ and that *w.r.t.* $\xi_n^k$ by applying (27) and (28), respectively;
7:   Updating the phase shifts, i.e., $\theta_m^l$ and $\xi_n^k$, by applying (29) and (30), respectively;
8:   Updating the scaling factor $\alpha$ by applying (31);
9:   Diminishing the learning rate $\eta$ by applying (32);
10:   Calculating the objective function value $\Gamma$ by applying (22a);
11: **UNTIL** The decrement of $\Gamma$ is less than a preset threshold value or the number of iterations reaches the maximum;
12: **OUTPUT:** $\boldsymbol{\theta}^l$, $l \in \mathcal{L}$, $\boldsymbol{\xi}^k$, $k \in \mathcal{K}$.

Given a set of phase shifts associated with the TX-SIM and RX-SIM, the resultant end-to-end channel becomes $\mathbf{H} = \mathbf{QGP}$. Consequently, the optimal scaling factor can be readily obtained by applying the least-square technique as follows

$$\alpha = \left(\mathbf{h}^H \mathbf{h}\right)^{-1} \mathbf{h}^H \boldsymbol{\lambda}, \tag{31}$$

where we have $\boldsymbol{\lambda} = \text{vec}\left(\mathbf{\Lambda}_{1:S,1:S}\right) \in \mathbb{C}^{S^2 \times 1}$ and $\mathbf{h} = \text{vec}\left(\mathbf{H}\right) \in \mathbb{C}^{S^2 \times 1}$ denoting the vectorization of $\mathbf{\Lambda}_{1:S,1:S}$ and that of $\mathbf{H}$, respectively.

***Step 5: Update the learning rate***

For the sake of avoiding any overshooting effects [59], we adopt a negative exponentially decaying learning schedule for decreasing the learning rate, as the iteration proceeds. More specifically, the learning rate is updated by

$$\eta \leftarrow \eta\beta, \tag{32}$$

at each iteration, where $0 < \beta < 1$ is a hyperparameter controlling the decay rate.

After repeating (23) $\sim$ (32) several times, the loss function $\Gamma$ gradually approaches convergence. In order to prevent the gradient descent algorithm from getting trapped in a local optimum, we first generate multiple sets of phase shifts and then select the one that minimizes $\Gamma$ for initialization. For clarity, we summarize the detailed procedures of the proposed gradient descent in Table I.

## IV. PERFORMANCE ANALYSIS

### A. HMIMO Channel Capacity Analysis

In this subsection, we evaluate the channel capacity of our HMIMO system. We assume that the TX-SIM and RX-SIM shown in Fig. 2 have carried out perfect precoding and combining in the EM regime. However, due to the fact that (19) cannot be readily expressed in closed form, here we provide an upper and a lower bound for the channel capacity of our HMIMO system. Specifically, by assuming that all the data streams experience either the best or the worst subchannel,

respectively, the ergodic channel capacity is upper and lower bounded by

$$S \log_2 \left(1 + \frac{P_t \mathbb{E}\left(\lambda_S^2\right)}{S\sigma^2}\right) \leq \mathbb{E}\left(C\right) \leq S \log_2 \left(1 + \frac{P_t \mathbb{E}\left(\lambda_1^2\right)}{S\sigma^2}\right), \tag{33}$$

where $\mathbb{E}\left(\lambda_1^2\right)$ and $\mathbb{E}\left(\lambda_S^2\right)$ denote the statistical average of the 1-st and the $S$-th eigenvalues, respectively, which are obtained through numerical approximations.

In order to gain some fundamental insights into the HMIMO channel capacity, we next analyze its scaling law by considering some special cases. Specifically, the HMIMO channel capacity evaluated by considering a large number of data streams is characterized by *Proposition 1*.

*Proposition 1:* As $S \rightarrow \infty$, we have $P_t \log_2 e\mathbb{E}\left(\lambda_S^2\right)/\sigma^2 \leq \mathbb{E}\left(C\right) \leq P_t \log_2 e\mathbb{E}\left(\lambda_1^2\right)/\sigma^2$.

*Proof:* By taking the limit of lower bound and upper bound in (33) as $S \rightarrow \infty$ and applying the formula $\lim_{x \rightarrow \infty}\left(1 + 1/x\right)^x = e$, the proof is completed. ∎

*Proposition 1* demonstrates that blindly increasing the number of active components may not lead to substantial improvements in channel capacity. Specifically, the HMIMO channel capacity gradually saturates as the number of data streams increases due to the intrinsic limitations of *spatial multiplexing* in the HMIMO channel. Further increasing the number of active components may result in severe energy efficiency degradation. Therefore, the resultant HMIMO channel capacity critically depends on the statistical expectation of the eigenvalues, which results in a *selection gain* associated with the increased number of meta-atoms.

In order to characterize the fundamental scaling law of the HMIMO channel capacity versus the number of meta-atoms, we next consider the particular case of $S = 1$, $L = K = 1$ and assume i.i.d. Rayleigh fading for the sake of brevity. Specifically, the theoretical ergodic channel capacity is summarized in *Proposition 2*.

*Proposition 2:* As $M, N \rightarrow \infty$, we have $\mathbb{E}\left(C\right) \simeq \log_2\left(1 + \frac{\pi^2 P_t \rho^2}{4\sigma^2}M^2N^2\right)$.

*Proof:* As $M, N \rightarrow \infty$, the HMIMO channel capacity can be approximated by

$$\mathbb{E}\left(C\right) \simeq \log_2\left[1 + \frac{P_t}{\sigma^2}\mathbb{E}\left(|h|^2\right)\right], \tag{34}$$

with

$$\mathbb{E}\left(|h|^2\right) = \rho^2 \mathbb{E}\left(\left|\sum_{m=1}^{M}\phi_m h_{1,m}\right|^2 \left|\sum_{n=1}^{N}\psi_n h_{2,n}\right|^2\right), \tag{35}$$

where $h_{1,m} \sim \mathcal{CN}\left(0,1\right)$, $m \in \mathcal{M}$ and $h_{2,n} \sim \mathcal{CN}\left(0,1\right)$, $n \in \mathcal{N}$ denote the normalized channel coefficient of the link spanning from the source to the optimal scatterer via the $m$-th transmit meta-atom and that from the optimal scatterer to the destination via the $n$-th receive meta-atom, respectively. Note that the unnecessary superscripts and subscripts have been removed for the sake of brevity.

Furthermore, by applying the optimal phase shift configuration, i.e., $\phi_m = h^*_{1,m}/|h_{1,m}|$ and $\psi_n = h^*_{2,n}/|h_{2,n}|$, we have [21]

$$\mathbb{E}\left(\left|\sum_{m=1}^{M}\phi_m h_{1,m}\right|^2\right) = \mathbb{E}\left(\left|\sum_{m=1}^{M}|h_{1,m}|\right|^2\right) = \frac{\pi}{2}M^2, \quad (36)$$

$$\mathbb{E}\left(\left|\sum_{n=1}^{N}\psi_n h_{2,n}\right|^2\right) = \mathbb{E}\left(\left|\sum_{n=1}^{N}|h_{2,n}|\right|^2\right) = \frac{\pi}{2}N^2. \quad (37)$$

By substituting (36) and (37) into (34), the proof is completed. ∎

*Proposition 2* testifies to the quadratic scaling law of the channel gain versus the number of meta-atoms [21]. Note that in contrast to the RIS-aided system [21], both the TX-SIM and RX-SIM could attain spatial gains. In an ideal setup, one could obtain about 4 bps/Hz HMIMO channel capacity improvement upon every doubling of the number of meta-atoms at both the TX-SIM and RX-SIM. Again, we note that the rigorous capacity analysis of SIM-aided HMIMO systems having an arbitrary number of metasurface layers is a complex task due to the fact that a large number of matrix multiplications are involved during the forward propagation [47], [60]. To address this issue, effective matrix analysis tools might be employed for evaluating the fitting performance of the proposed SIM as well as the resultant channel capacity, which requires further investigation. Nevertheless, our numerical results of Section V demonstrate that harnessing a pair of SIMs having a moderate number of metasurface layers at both ends can fit the end-to-end channel with high accuracy. As such, motivated readers may refer to [42] for gaining deeper insights concerning the HMIMO channel capacity relying on conventional digital precoding and combining, which serves as an upper bound for our SIM-aided HMIMO system.

### B. Computational Complexity Analysis

Next, we analyze the computational complexity of the proposed gradient descent algorithm in terms of the number of real-valued multiplications. Specifically, the computational complexity of performing Step 1 includes that of performing the forward propagation, i.e., $\mathcal{O}_{1-1} = 4SM(ML - M + L) + 4SN(NK - N + K) + 4MSN + 2(M + N)S^2$, and that of calculating all partial derivatives, i.e., $\mathcal{O}_{1-2} = 2(ML + NK)S^2$. Additionally, the computational complexities of performing Steps 2 ∼ 5 are $\mathcal{O}_2 = 2(ML + NK)$, $\mathcal{O}_3 = (ML + NK)$, $\mathcal{O}_4 = 6S^2 + 2$, and $\mathcal{O}_5 = 1$, respectively. As a result, the total computational complexity of the proposed gradient descent algorithm is given by

$$\begin{aligned}
\mathcal{O} &= I\left[\mathcal{O}_{1-1} + \mathcal{O}_{1-2} + \mathcal{O}_2 + \mathcal{O}_3 + \mathcal{O}_4 + \mathcal{O}_5\right] \\
&= I\left[4SM(ML - M + L) + 4SN(NK - N + K)\right] \\
&\quad + I\left[4MSN + 2(M + N)S^2\right] + 2I(ML + NK)S^2 \\
&\quad + 3I(ML + NK) + I\left(6S^2 + 3\right) \\
&\simeq 4IS\left(M^2L + N^2K\right), \quad \text{for } M, N \gg S, \quad (38)
\end{aligned}$$

where $I$ denotes the number of iterations, which will be shown in Section V to be less than 20 under an empirical setup. Thus,

the proposed gradient descent algorithm is of polynomial-time complexity, when solving Problem (22).

## V. SIMULATION RESULTS

In this section, we provide numerical results for characterizing the performance of the proposed SIM-aided HMIMO system.

### A. Simulation Setups

As illustrated in Fig. 2, we consider an SIM-aided HMIMO system. In our simulations, the thicknesses of both the TX-SIM and RX-SIM are set to $D_t = D_r = 0.05$ m. Accordingly, the transmission coefficients between the adjacent metasurface layers in the TX-SIM and that in the RX-SIM are given by (9) and (11), respectively, while the HMIMO channel is generated by (13). Our SIM-aided HMIMO system operates at the frequency of $f_0 = 28$ GHz, which corresponds to the wavelength of $\lambda = 10.7$ mm. To account for the large-scale fading, we consider the reference distance of $d_0 = 1$ m and set $b = 3.5$ and $\delta = 9$ dB in our simulations [66]. Moreover, the distance between the transmitter and the receiver is set to $d = 250$ m.

Additionally, the total power available at the transmitter is set to $P_t = 20$ dBm, while the receiver sensitivity is set to $\sigma^2 = -110$ dBm. For the proposed gradient descent algorithm, the number of randomizations for initialization is set to 10. The maximum affordable number of iterations is set to 100, while the initial learning rate and decay parameter are set to $\eta_0 = 0.1$ and $\beta = 0.5$, respectively, unless otherwise specified. All the simulation results are obtained by averaging over 100 independent experiments. Moreover, we adopt a couple of different performance metrics. Specifically, for the sake of a fair comparison, we first quantify the normalized mean square error (NMSE) between the actual channel matrix and the target diagonal one defined as follows

$$\Delta = \mathbb{E}\left(\frac{\|\alpha\mathbf{QGP} - \mathbf{\Lambda}_{1:S,1:S}\|_{\mathrm{F}}^2}{\|\mathbf{\Lambda}_{1:S,1:S}\|_{\mathrm{F}}^2}\right), \quad (39)$$

while the other is the channel capacity of our SIM-assisted HMIMO system, which is defined by

$$C = \sum_{s=1}^{S}\log_2\left(1 + \frac{p_s|\alpha h_{s,s}|^2}{\sum_{\tilde{s}\neq s}^{S}p_{\tilde{s}}|\alpha h_{s,\tilde{s}}|^2 + \sigma^2}\right), \quad (40)$$

where $p_s$ denotes the amount of power allocated to the $s$-th stream obtained by (18). Note that in sharp contrast to conventional MIMO designs with extra digital combining, our SIM-aided HMIMO treats the residual signals from other data streams as interference, as seen in the denominator of (40).

### B. Performance versus System Parameters

Fig. 3 first evaluates the NMSE between the actual channel matrix and the target diagonal matrix for different numbers of metasurface layers, where we consider $S = 4$, $M = N = 100$, and $r_{e,t} = t_{e,r} = \lambda/2$. As such, $A_t$ and $A_r$ are determined accordingly. Observe from Fig. 3 that the channel fitting NMSE gradually decreases as the number of metasurface
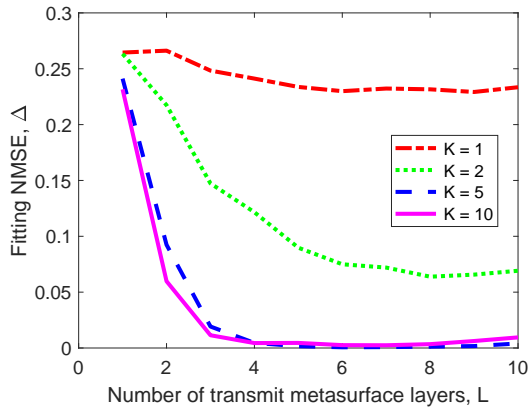
Fig. 3. The NMSE between the actual channel matrix and the target diagonal one versus the number of TX metasurface layers, where we have $S = 4$, $M = N = 100$, and $r_{e,t} = t_{e,r} = \lambda/2$.
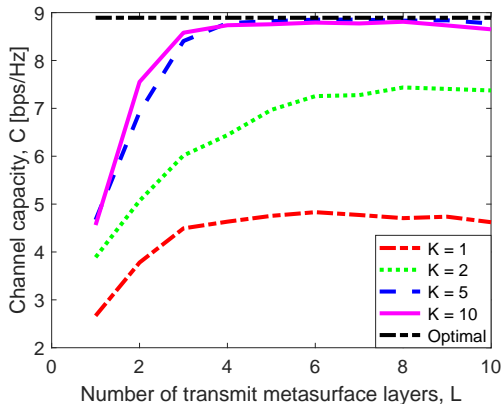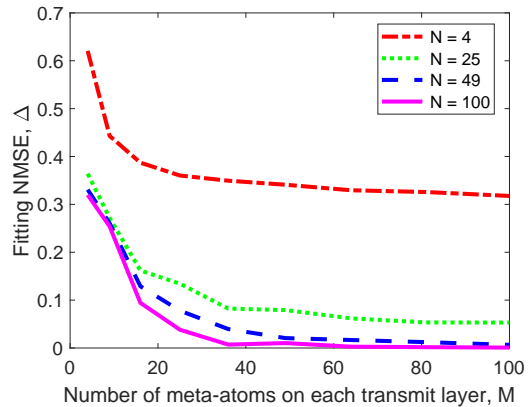


Fig. 5. The NMSE between the actual channel matrix and the target diagonal one versus the number of meta-atoms on each TX metasurface layer, where we have $S = 4$, $L = K = 7$, and $r_{e,t} = t_{e,r} = \lambda/2$.



Fig. 4. The channel capacity versus the number of TX metasurface layers, where we have $S = 4$, $M = N = 100$, and $r_{e,t} = t_{e,r} = \lambda/2$.
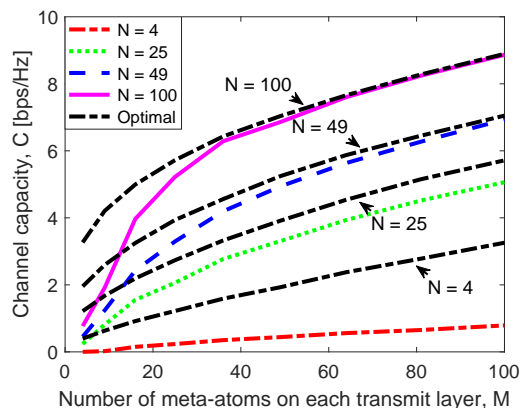


Fig. 6. The channel capacity versus the number of meta-atoms on each TX metasurface layer, where we have $S = 4$, $L = K = 7$, and $r_{e,t} = t_{e,r} = \lambda/2$.

layers increases and eventually bottoms out, when $L \geq 5$ or $K \geq 5$, thanks to the powerful inference capability of the multi-layer architecture advocated. Furthermore, Fig. 4 shows the corresponding channel capacity, where the channel capacity of the HMIMO system having a full-precision precoder and combiner is also plotted. It reveals that the SIM-aided HMIMO channel capacity gradually saturates as the number of metasurface layers increases. The SIM having an adequate number of metasurface layers might approach the capacity upper bound characterized by the full-precision precoding and combining. However, due to the fixed thickness of the SIM considered, i.e., $D_t$ and $D_r$, excessively dense metasurfaces may lead to a performance penalty, when the number of metasurface layers exceeds a certain threshold. For example, the SIM-aided HMIMO using $L = K = 10$ metasurface layers suffers both from some fitting performance erosion as well as from a capacity reduction compared to $L = K = 5$. *In a nutshell, both the channel fitting NMSE and the channel capacity approach their optimal values when using $L = 7$ metasurface layers, and further increasing the number of metasurface layers does not help to improve the fitting NMSE and channel capacity.*

In Fig. 5, we portray the channel fitting NMSE versus the number of meta-atoms, where the number of metasurface layers is set to $L = K = 7$, with all other system parameters remaining unchanged. Observe from Fig. 5 that the fitting NMSE decreases monotonically as the number of meta-atoms on each transmit or RX metasurface layer increases. As a benefit, the system becomes capable of establishing a perfectly diagonal end-to-end channel matrix for an infinite number of meta-atoms, i.e., $\Delta \to 0$ for $M \to \infty$ or $N \to \infty$. Furthermore, Fig. 6 plots the channel capacity versus the number of meta-atoms on each TX metasurface layer. Note that as a benefit of the substantial *selection gain* discussed in *Proposition 1*, one can always select the best $S$ subchannels for conveying information [62]. The channel capacity is improved as the number of meta-atoms increases, albeit the number of data streams is fixed. For example, the SIM-aided HMIMO system behaves competitively with its counterpart having full-precision precoding and combiner as $M, N \to 100$. *More specifically, for an adequate number of meta-atoms having tolerable fitting errors, say $N \geq 25$, Fig. 6 confirms our Proposition 2 that the channel capacity would increase by*
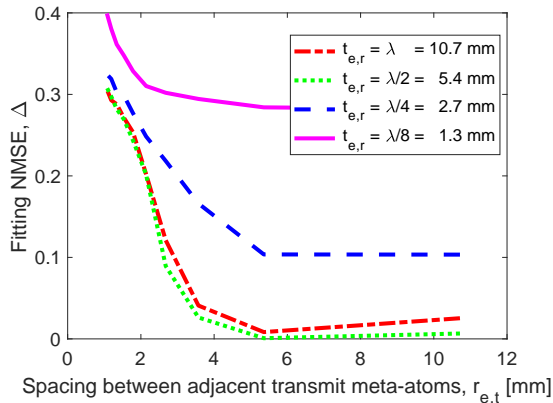
Fig. 7. The NMSE between the actual channel matrix and the target diagonal one versus the spacing between adjacent meta-atoms on each TX metasurface layer, where we have $S = 4$, $L = K = 7$, and $M = N = 100$.



Fig. 9. The NMSE between the actual channel matrix and the target diagonal one versus the number of data streams, where we have $L = K = 7$, and $r_{e,t} = t_{e,r} = \lambda/2$.
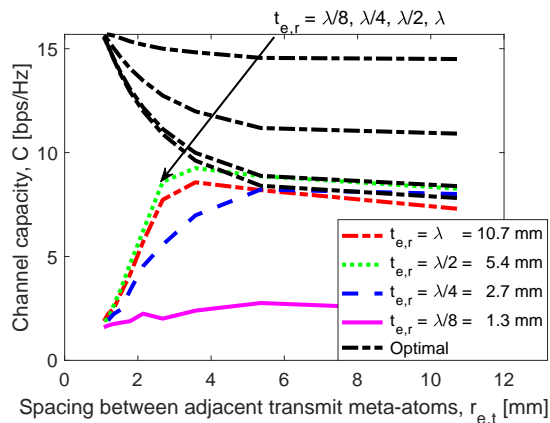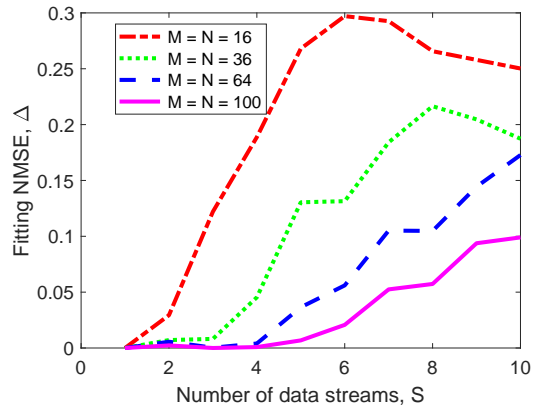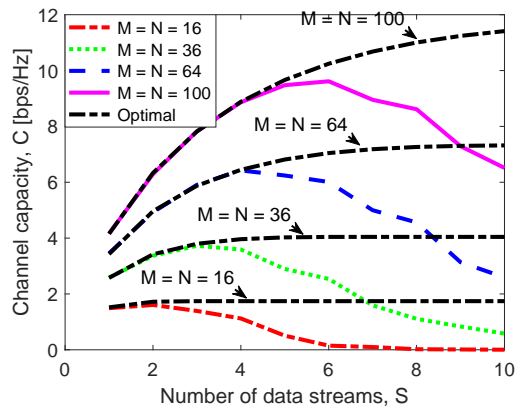


Fig. 8. The channel capacity versus the spacing between adjacent meta-atoms on each TX metasurface layer, where we have $S = 4$, $L = K = 7$, and $M = N = 100$.



Fig. 10. The channel capacity versus the number of data streams, where we have $L = K = 7$, and $r_{e,t} = t_{e,r} = \lambda/2$.

*about* 4 *bps/Hz when doubling the number of meta-atoms in both the TX-SIM and RX-SIM.*

Furthermore, Figs. 7 and 8 quantify the channel fitting NMSE and the channel capacity, respectively, versus the spacing between adjacent meta-atoms, where we set $L = K = 7$, $M = N = 100$, and increase the element spacing from $\lambda/10 = 1.1$ mm to $\lambda = 10.7$ mm. It is shown in Fig. 7 that the channel fitting NMSE achieves its minimum at $r_{e,t} = t_{e,r} = \lambda/2$. *Both a larger and a smaller element spacing would give rise to the similarity between the transmission coefficients as well as lead to undesired channel correlations, thus resulting in a poor channel fitting NMSE. Observe from Fig. 8 that RX-SIM having element spacing of about $t_{e,r} = \lambda/2$ attains the maximal capacity as well as the best fit with the full-precision counterpart.* Given the half-wavelength element spacing between adjacent meta-atoms on each TX metasurface, i.e., $r_{e,t} = \lambda/2$, both the setups of $t_{e,r} = \lambda$ and $t_{e,r} = \lambda/4$ suffer from a capacity loss of about 1 bps/Hz compared with that adopting $t_{e,r} = \lambda/2$. The inferior fitting NMSE by taking a small element spacing also widens the performance gap between the SIM-assisted HMIMO and its full-precision counterpart. In a nutshell, a

better channel fitting NMSE means that parallel subchannels are perfectly formed and suffer from less interference, thus achieving an improved channel capacity for our SIM-aided HMIMO system.

Next, Figs. 9 and 10 examine the channel fitting NMSE and channel capacity, respectively, versus the number of data streams, where we consider $L = K = 7$, and $r_{e,t} = t_{e,r} = \lambda/2$. It can be seen from Fig. 9 that we have $\Delta = 0$ for $S = 1$ under all the setups, and the fitting performance gradually degrades as the number of data streams increases due to the larger dimension of the channel matrix. By utilizing a pair of TX-SIM and RX-SIM having small metasurface profiles, such as $M = N = 16$, we achieve a fitting NMSE of $\Delta = 0.2$ for $S = 4$, which is reduced to $\Delta < 0.001$ upon increasing the number of meta-atoms on each metasurface layer to $M = N = 100$. Furthermore, in sharp contrast to the full-precision counterpart, the SIM-aided HMIMO channel capacity mainly relies on two factors: the number of data streams and the channel fitting NMSE. On one hand, the increasing number of data streams may offer a proportional *multiplexing gain* [62]. On the other hand, it becomes more challenging to acquire a low channel fitting NMSE for a growing number of data
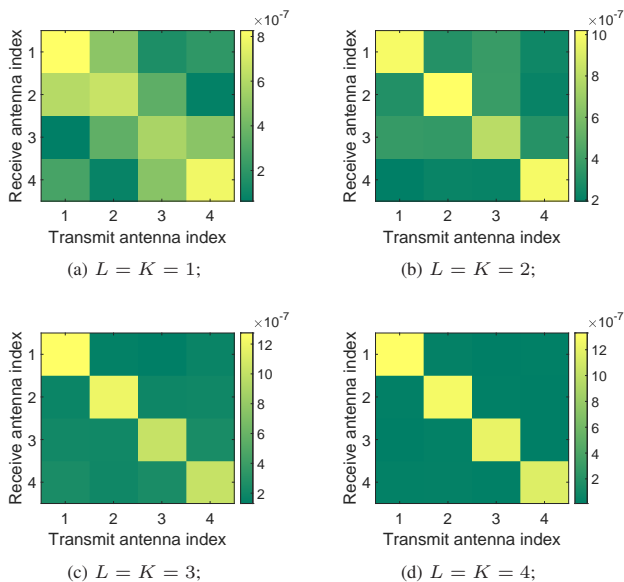
(a) $L = K = 1$;

(b) $L = K = 2$;

(c) $L = K = 3$;

(d) $L = K = 4$;

Fig. 11. The visualization of the end-to-end spatial channel matrix $\mathbf{H} = \mathbf{QGP}$.



(a) $\eta_0 = 0.1$;

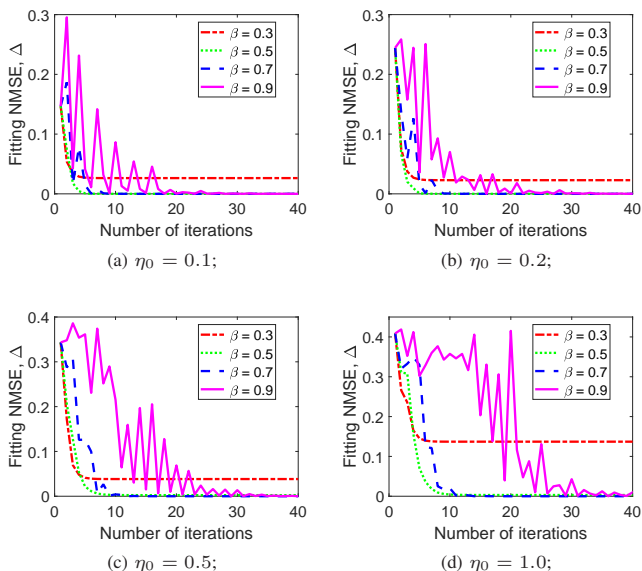(b) $\eta_0 = 0.2$;

(c) $\eta_0 = 0.5$;

(d) $\eta_0 = 1.0$;

Fig. 12. The convergence curves of the proposed gradient descent algorithm.

streams, thus leading to severe interference among different subchannels. *Due to this fundamental tradeoff between the multiplexing gain and the channel fitting NMSE, the channel capacity achieves its maximum for a certain number of data streams, e.g., $S = 2, 3, 4, 6$ for the four setups considered in Fig. 10.* In addition, note that the ideal channel capacity approaches saturation as the number of data streams increases, which is consistent with our *Proposition 1*.

### C. Validation of the Proposed Algorithm

For the sake of illustration, Fig. 11 visualizes the end-to-end channel matrix $\mathbf{H} = \mathbf{QGP}$ for different numbers of metasurface layers, where we consider $S = 4$, $M = N = 100$,
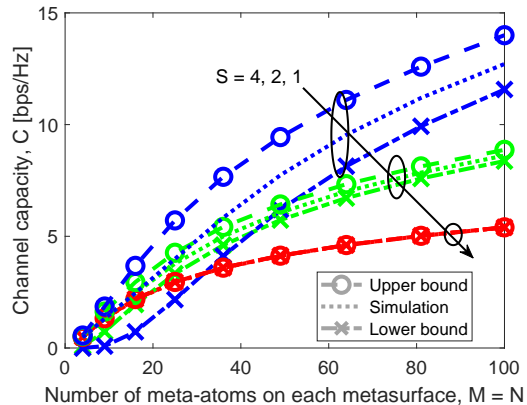


Fig. 13. Channel capacity comparison of the simulation and analytical results based on (33).

and $r_{e,t} = t_{e,r} = \lambda/2$. Observe from Fig. 11 that for a small number of metasurface layers, such as $L = K = 1$, the TX-SIM and RX-SIM struggle to form a diagonal channel matrix spanning from the source to the destination. Hence, each data stream suffers from the interference imposed by other streams, ultimately resulting in a reduced channel capacity (see Fig. 4). As the number of metasurface layers increases, the TX-SIM and RX-SIM attain a stronger inference capability and thus may form multiple parallel subchannels in the physical space. Fig. 11(d) shows that the TX-SIM and RX-SIM having four metasurface layers respectively succeed in forming an almost perfectly diagonal channel matrix, thus asymptotically achieving the maximal channel capacity.

Next, we examine the convergence performance of the proposed gradient descent algorithm by considering different values of the initial learning rate $\eta_0$ and the decay parameter $\beta$. As shown in Fig. 12(a), we begin by analyzing the case of $\eta_0 = 0.1$. Note that under all the setups, the fitting NMSE eventually decreases, thus facilitating convergence. However, for a small value of the decay parameter, e.g., $\beta = 0.3$, the channel fitting NMSE might converge to a local minimum. By contrast, a larger value of the decay parameter, e.g., $\beta = 0.9$, may result in overshooting effects. As a result, the fitting NMSE fluctuates violently during the initial iteration stage. Furthermore, when we increase the initial learning rate to $\eta_0 = 0.2, 0.5, 1.0$, the corresponding results are shown in Figs. 12(b), 12(c), and 12(d), respectively. It is demonstrated that an excessively high initial learning rate may require a long period to achieve convergence. For example, more than 40 iterations are required for the setup of $\eta_0 = 1.0$ and $\beta = 0.9$. Nonetheless, in all cases, the fitting NMSE can converge to the desired accuracy after a sufficient number of iterations.

### D. Performance Evaluation and Comparison to Existing Transmission Technologies

Furthermore, Fig. 13 verifies the accuracy of our theoretical analysis of the HMIMO channel capacity, where we consider $r_{e,t} = t_{e,r} = \lambda/4$ and increase the number of data streams from $S = 1$ to $S = 4$. Observe from Fig. 13 that the channel capacity increases with the number of meta-atoms on each
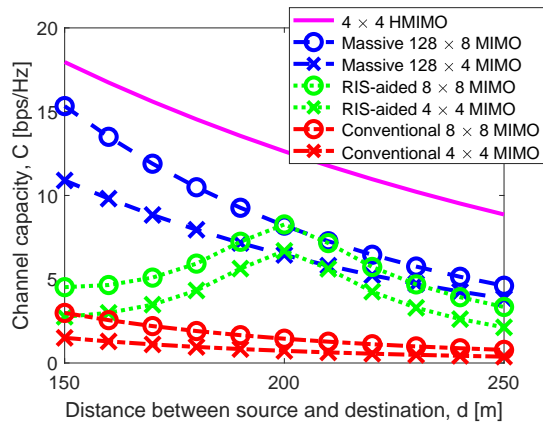
Fig. 14. Channel capacity comparison of our SIM-aided HMIMO ($S = 4$, $L = K = 7$, $M = N = 100$, $r_{e,t} = t_{e,r} = \lambda/2$) and other MIMO transmission schemes.



Fig. 16. The channel fitting NMSE comparison of the multilayer SIM and its single-layer counterpart that has the same total number of meta-atoms.
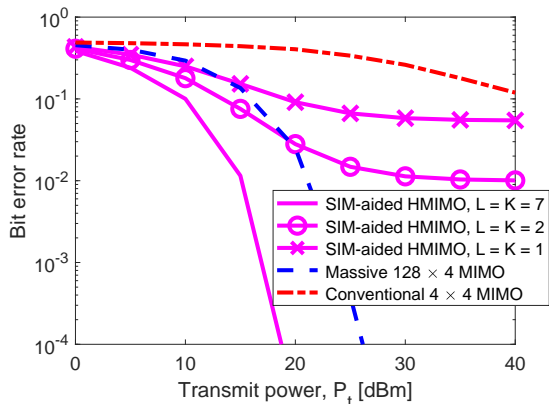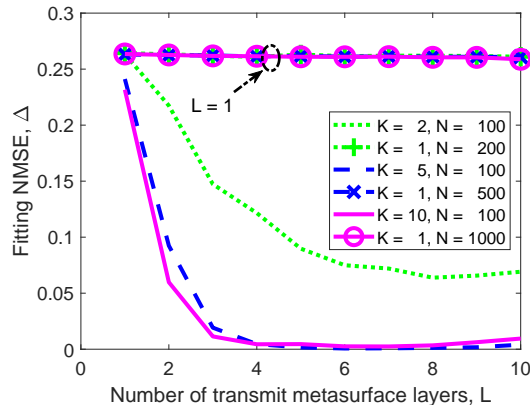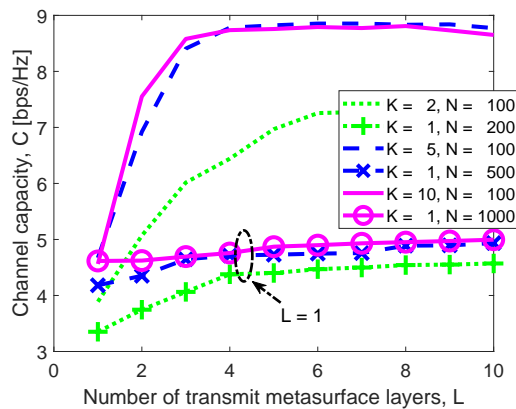


Fig. 15. Bit error rate comparison of our SIM-aided HMIMO ($S = 4$, $M = N = 100$, $r_{e,t} = t_{e,r} = \lambda/2$) and other MIMO transmission schemes, where the transmission rate is 4 bpcu.



Fig. 17. The channel capacity comparison of the multilayer SIM and its single-layer counterpart that has the same total number of meta-atoms.

metasurface as well as that of data streams, which is due to the substantial *selection gain* and *multiplexing gain*, respectively [62]. Specifically, when considering $S = 4$ data streams, the 8 bps/Hz capacity increase is observed by quadrupling the number of meta-atoms on each metasurface from $M = N = 25$ to $M = N = 100$, which is consistent with our previous analysis in *Proposition 2*. Furthermore, as expected, the actual channel capacity of our HMIMO communication system consistently lies between the upper and lower bounds derived. Specifically, both the upper and lower bounds are tight for a single data stream, i.e., $S = 1$. As the number of data streams increases, there is a widening gap between the analytical and simulation results due to the imperfect scaling operation in (33). Deriving the accurate capacity calls for future research.

In Fig. 14, we compare the channel capacity of our SIM-assisted HMIMO system to that of the massive MIMO scheme as well as to its RIS-aided MIMO counterpart. The detailed MIMO setups are shown in Fig. 14. Specifically, we adopt a pair of TX-SIM and RX-SIM ($S = 4$, $L = K = 7$, $M = N = 100$, $r_{e,t} = t_{e,r} = \lambda/2$) for performing the wave-based precoding and combining, while achieving the spatial gains. As for the RIS-aided MIMO scheme, a RIS having

1,000 elements is deployed at a site having a source-RIS distance of 200 m and a vertical spacing of 10 m *w.r.t.* the source-destination link to enhance the channel quality. Thus, we have a RIS-destination distance of $\sqrt{50^2 + 10^2} \approx 51$ m. All the channels are assumed to be Rayleigh fading along with the path loss model in (16). The path loss exponents are adjusted to $b = 2.2$ and $b = 2.7$ for the source-RIS and RIS-destination links, respectively [69]. Observe from Fig. 14 that as a benefit of the substantial spatial gain attained by the TX-SIM and RX-SIM, our HMIMO outperforms both its massive MIMO and RIS-aided counterparts under all the setups considered. Specifically, although RIS achieves significant capacity improvements over the conventional MIMO, it still suffers from a performance gap compared to the HMIMO due to the severe two-hop path loss. Even in the vicinity of RIS, the HMIMO attains a 150% capacity gain, which may increase to 200% at the cell edge, e.g., at $d = 250$ m. Additionally, the massive MIMO equipped with a huge number of active elements achieves impressive capacity improvements at the cost of an increasing number of active RF chains, which, however, still has at least a 3 bps/Hz capacity penalty compared to the HMIMO scheme.

Fig. 15 compares the error performance of our SIM-assisted

HMIMO system to conventional MIMO schemes. Specifically, we consider four data streams, each transmitting a BPSK symbol, which corresponds to a transmission rate of 4 bits per channel use (bpcu). The distance between the transmitter and receiver is set to 200 m, and the number of metasurface layers is set to $L = K = 1$, 2, 7, respectively. As shown in Fig. 15, deploying a sufficient number of metasurface layers, such as $L = K = 7$ in both the TX-SIM and RX-SIM, effectively mitigates the inter-stream interference, thanks to carrying out the precoding and combining in the wave domain. Furthermore, as a benefit of the spatial gain attained by the large transceiver surface aperture, the SIM-assisted HMIMO achieves a lower bit error rate (BER) than its large-scale MIMO counterpart. Observe from Fig. 15 that SIM-aided HMIMO has an 8 dB performance gain compared to massive MIMO in this setup. However, when reducing the number of metasurface layers to $L = K = 2$, the TX-SIM and RX-SIM modules failed to perfectly suppress the interference amongst the data streams. Consequently, the SIM-assisted HMIMO scheme suffers from performance erosion, which results in a residual BER as the transmit power $P_t$ increases. Note that we directly use the SIM phase shifts optimized in Section III-B to evaluate the BER performance and that directly optimizing the SIM for minimizing the BER may result in further performance improvements.

Finally, Fig. 16 compares the channel fitting NMSE of the multilayer SIM to its single-layer counterpart. Specifically, the total number of meta-atoms in the TX-SIM and RX-SIM are rearranged into a single-layer metasurface, respectively. In order to maintain the same transceiver surface area, the meta-atom spacing is set to $r_{e,t} = 5\lambda/\left\lceil\sqrt{ML}\right\rceil$ and $t_{e,r} = 5\lambda/\left\lceil\sqrt{NK}\right\rceil$. All other simulation parameters are consistent with those in Fig. 3. Observe from Fig. 16 that the single-layer SIM fails to accurately fit the expected end-to-end channel, even with an adequate number of meta-atoms. By contrast, the multilayer SIM structure achieves a superior channel fitting NMSE as the number of metasurface layers increases. Furthermore, Fig. 17 shows the channel capacity of these two transmission schemes. It is evident that the single-layer SIM provides only marginal capacity improvements as the number of meta-atoms increases, which is primarily due to the inability of the single-layer TX-SIM and RX-SIM to effectively suppress the inter-stream interference. Moreover, deploying a large number of meta-atoms in a limited space also leads to channel correlation. Observe from Fig. 17 that the proposed SIM-aided HMIMO system associated with $L = K = 10$ and $M = N = 100$ achieves almost twice the capacity improvement compared to its single-layer counterpart for the same total number of meta-atoms. In a nutshell, both the channel fitting NMSE and the corresponding channel capacity of the single-layer SIM suffer from significant performance penalties compared to its multi-layer SIM counterpart. Nevertheless, finding the optimal SIM design under a given total number of meta-atoms remains an open research question that requires further investigation.

## VI. Conclusions

In this paper, we proposed a novel SIM-based HMIMO communication paradigm, which attains substantial spatial gains while performing the precoding and combining functionalities directly in the native EM regime at the speed of light. We first formulated a channel fitting problem to approximate the MIMO-capacity-optimal diagonal channel matrix by optimizing the phase shifts of both the TX-SIM and RX-SIM. Then, we proposed an efficient gradient descent algorithm for iteratively solving that non-trivial fitting problem. Additionally, we derived a numerical approximation method for characterizing the HMIMO channel capacity and derived some fundamental capacity scaling laws. Finally, extensive simulations were provided for validating the benefits of the proposed SIM-based HMIMO system, demonstrating that substantial capacity improvements were attained upon increasing the number of the SIMs' meta-atoms.

In conclusion, our pivotal findings are as follows. Firstly, our experimental insights have shed light on the optimal SIM design. Specifically, we found that a 7-layer SIM having half-wavelength element spacing achieves an excellent channel fitting performance approaching the MIMO channel capacity. Secondly, both our theoretical analysis and simulation results have shown a quadratic channel gain when doubling the number of meta-atoms. Additionally, we have verified the performance advantages of the proposed HMIMO scheme over the existing benchmark schemes. Notably, a 150% capacity gain was attained over its conventional massive MIMO and RIS-assisted counterparts. As such, the multilayer SIM structure is capable of carrying out signal processing in the wave domain, which might lead to disruptive implementation-oriented advances. Moreover, an active SIM may be created by integrating small power amplifiers within some of the meta-atoms [47]. Upon adjusting the drive level of these power amplifiers, a non-linear module can be produced for further enhancing the inference capability of the SIM. Nonetheless, accurately evaluating the achievable performance gain of the active SIM requires further investigation.

## References

[1] 3GPP, "Release 17 Description; Summary of Rel-17 Work Items," Technical Report (TR) 21.917, 3rd Generation Partnership Project (3GPP), Jul. 2022.

[2] Samsung, "The next hyper-connected experience for all," 6G white paper, Samsung, Jul. 2020.

[3] Cisco, "The future of work," Annual Report, Cisco, Oct. 2021.

[4] S. Dang, O. Amin, B. Shihada, and M.-S. Alouini, "What should 6G be?," *Nat. Electron.*, vol. 3, pp. 20–29, Jan. 2020.

[5] K. B. Letaief, W. Chen, Y. Shi, J. Zhang, and Y.-J. A. Zhang, "The roadmap to 6G: AI empowered wireless networks," *IEEE Commun. Mag.*, vol. 57, pp. 84–90, Aug. 2019.

[6] G. C. Alexandropoulos, K. Stylianopoulos, C. Huang, C. Yuen, M. Bennis, and M. Debbah, "Pervasive machine learning for smart radio environments enabled by reconfigurable intelligent surfaces," *Proc. IEEE*, vol. 110, pp. 1494–1525, Sep. 2022.

[7] O. Kodheli, E. Lagunas, N. Maturo, S. K. Sharma, B. Shankar, J. F. M. Montoya, J. C. M. Duncan, D. Spano, S. Chatzinotas, S. Kisseleff, J. Querol, L. Lei, T. X. Vu, and G. Goussetis, "Satellite communications in the new space era: A survey and future challenges," *IEEE Commun. Surveys Tuts.*, vol. 23, pp. 70–109, 1st Quart. 2021.

[8] J. An, C. Xu, Q. Wu, D. W. K. Ng, M. Di Renzo, C. Yuen, and L. Hanzo, "Codebook-based solutions for reconfigurable intelligent surfaces and their open challenges," *IEEE Wireless Commun.*, pp. 1–8, 2022, Early Access.

[9] B. Yang, X. Cao, C. Huang, G. C. Alexandropoulos, L. Dai, C. Yuen, M. Debbah, and H. V. Poor, "Next generation reconfigurable metasurfaces: When wave propagation control meets computing," *IEEE Wireless Commun.*, to appear, 2023.

[10] J. An, H. Li, D. W. K. Ng, and C. Yuen, "Fundamental detection probability vs. achievable rate tradeoff in integrated sensing and communication systems," *arXiv preprint arXiv:2305.02847*, 2023.

[11] C. Xu, L. Xiang, J. An, C. Dong, S. Sugiura, R. G. Maunder, L.-L. Yang, and L. Hanzo, "OTFS-aided RIS-assisted SAGIN systems outperform their OFDM counterparts in doubly selective high-doppler scenarios," *IEEE Int. Things J.*, vol. 10, pp. 682–703, Jan. 2023.

[12] E. Calvanese Strinati, G. C. Alexandropoulos, H. Wymeersch, B. Denis, V. Sciancalepore, R. D'Errico, A. Clemente, D.-T. Phan-Huy, E. De Carvalho, and P. Popovski, "Reconfigurable, intelligent, and sustainable wireless environments for 6G smart connectivity," *IEEE Commun. Mag.*, vol. 59, pp. 99–105, Oct. 2021.

[13] Q. Wu and R. Zhang, "Towards smart and reconfigurable environment: Intelligent reflecting surface aided wireless network," *IEEE Commun. Mag.*, vol. 58, pp. 106–112, Jan. 2020.

[14] X. Yu, V. Jamali, D. Xu, D. W. K. Ng, and R. Schober, "Smart and reconfigurable wireless communications: From IRS modeling to algorithm design," *IEEE Wireless Commun.*, vol. 28, pp. 118–125, Dec. 2021.

[15] H. Zhang, S. Zeng, B. Di, Y. Tan, M. Di Renzo, M. Debbah, Z. Han, H. V. Poor, and L. Song, "Intelligent omni-surfaces for full-dimensional wireless communications: Principles, technology, and implementation," *IEEE Commun. Mag.*, vol. 60, pp. 39–45, Feb. 2022.

[16] M. Jian, G. C. Alexandropoulos, E. Basar, C. Huang, R. Liu, Y. Liu, and C. Yuen, "Reconfigurable intelligent surfaces for wireless communications: Overview of hardware designs, channel models, and estimation techniques," *Intelligent, Converged Netw.*, vol. 3, pp. 1–32, Mar. 2022.

[17] G. C. Alexandropoulos, G. Lerosey, M. Debbah, and M. Fink, "Reconfigurable intelligent surfaces and metamaterials: The potential of wave propagation control for 6G wireless communications," *IEEE ComSoc TCCN Newslett.*, vol. 6, pp. 25–37, Jun. 2020.

[18] J. An, C. Xu, L. Gan, and L. Hanzo, "Low-complexity channel estimation and passive beamforming for RIS-assisted MIMO systems relying on discrete phase shifts," *IEEE Trans. Commun.*, vol. 70, pp. 1245–1260, Feb. 2022.

[19] G. C. Alexandropoulos, N. Shlezinger, and P. del Hougne, "Reconfigurable intelligent surfaces for rich scattering wireless communications: Recent experiments, challenges, and opportunities," *IEEE Commun. Mag.*, vol. 59, pp. 28–34, Jun. 2021.

[20] C. Huang, A. Zappone, G. C. Alexandropoulos, M. Debbah, and C. Yuen, "Reconfigurable intelligent surfaces for energy efficiency in wireless communication," *IEEE Trans. Wireless Commun.*, vol. 18, pp. 4157–4170, Aug. 2019.

[21] Q. Wu, S. Zhang, B. Zheng, C. You, and R. Zhang, "Intelligent reflecting surface-aided wireless communications: A tutorial," *IEEE Trans. Commun.*, vol. 69, pp. 3313–3351, May 2021.

[22] W. Xu, J. An, Y. Xu, C. Huang, L. Gan, and C. Yuen, "Time-varying channel prediction for RIS-assisted MU-MISO networks via deep learning," *IEEE Trans. Cognitive Commun. Netw.*, vol. 8, pp. 1802–1815, Dec. 2022.

[23] C. Xu, J. An, T. Bai, S. Sugiura, R. G. Maunder, Z. Wang, L.-L. Yang, and L. Hanzo, "Channel estimation for reconfigurable intelligent surface assisted high-mobility wireless systems," *IEEE Trans. Veh. Technol.*, vol. 72, pp. 718–734, Jan. 2023.

[24] W. Xu, J. An, C. Huang, L. Gan, and C. Yuen, "Deep reinforcement learning based on location-aware imitation environment for RIS-aided mmwave MIMO systems," *IEEE Wireless Commun. Lett.*, vol. 11, pp. 1493–1497, Jul. 2022.

[25] Y. Han, W. Tang, S. Jin, C.-K. Wen, and X. Ma, "Large intelligent surface-assisted wireless communication exploiting statistical CSI," *IEEE Trans. Veh. Technol.*, vol. 68, pp. 8238–8242, Aug. 2019.

[26] C. Xu, J. An, T. Bai, L. Xiang, S. Sugiura, R. G. Maunder, L.-L. Yang, and L. Hanzo, "Reconfigurable intelligent surface assisted multi-carrier wireless systems for doubly selective high-mobility Ricean channels," *IEEE Trans. Veh. Technol.*, vol. 71, pp. 4023–4041, Apr. 2022.

[27] C. Pan, H. Ren, K. Wang, W. Xu, M. Elkashlan, A. Nallanathan, and L. Hanzo, "Multicell MIMO communications relying on intelligent reflecting surfaces," *IEEE Trans. Wireless Commun.*, vol. 19, pp. 5218–5233, Aug. 2020.

[28] J. An, Q. Wu, and C. Yuen, "Scalable channel estimation and reflection optimization for reconfigurable intelligent surface-enhanced OFDM systems," *IEEE Wireless Commun. Lett.*, vol. 11, pp. 796–800, Apr. 2022.

[29] V. Jamali, G. C. Alexandropoulos, R. Schober, and H. V. Poor, "Low-to-zero-overhead IRS reconfiguration: Decoupling illumination and channel estimation," *IEEE Commun. Lett.*, vol. 26, pp. 932–936, Apr. 2022.

[30] J. An and L. Gan, "The low-complexity design and optimal training overhead for IRS-assisted MISO systems," *IEEE Wireless Commun. Lett.*, vol. 10, pp. 1820–1824, Aug. 2021.

[31] X. Jia, J. An, H. Liu, H. Liao, L. Gan, and C. Yuen, "Environment-aware codebook for reconfigurable intelligent surface-aided MISO communications," *IEEE Wireless Commun. Lett.*, pp. 1–1, 2023, Early Access.

[32] E. G. Larsson, O. Edfors, F. Tufvesson, and T. L. Marzetta, "Massive MIMO for next generation wireless systems," *IEEE Commun. Mag.*, vol. 52, pp. 186–195, Feb. 2014.

[33] T. L. Marzetta, "Noncooperative cellular wireless with unlimited numbers of base station antennas," *IEEE Trans. Wireless Commun.*, vol. 9, pp. 3590–3600, Nov. 2010.

[34] C. Huang, S. Hu, G. C. Alexandropoulos, A. Zappone, C. Yuen, R. Zhang, M. D. Renzo, and M. Debbah, "Holographic MIMO surfaces for 6G wireless networks: Opportunities, challenges, and trends," *IEEE Wireless Commun.*, vol. 27, pp. 118–125, Oct. 2020.

[35] 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, pp. 2746–2758, May 2018.

[36] N. Shlezinger, G. C. Alexandropoulos, M. F. Imani, Y. C. Eldar, and D. R. Smith, "Dynamic metasurface antennas for 6G extreme massive MIMO communications," *IEEE Wireless Commun.*, vol. 28, pp. 106–113, Apr. 2021.

[37] T. Gong, P. Gavriilidis, R. Ji, C. Huang, G. C. Alexandropoulos, L. Wei, M. Debbah, H. V. Poor, and C. Yuen, "Holographic MIMO communications: Theoretical foundations, enabling technologies, and future directions," *arXiv preprint arXiv:2212.01257*, Apr. 2023.

[38] C. Xu, J. An, T. Bai, S. Sugiura, R. G. Maunder, L.-L. Yang, M. Di Renzo, and L. Hanzo, "Antenna selection for reconfigurable intelligent surfaces: A transceiver-agnostic passive beamforming configuration," *IEEE Trans. Wireless Commun.*, pp. 1–1, 2023, Early Access.

[39] Z. Wan, Z. Gao, F. Gao, M. Di Renzo, and M.-S. Alouini, "Terahertz massive MIMO with holographic reconfigurable intelligent surfaces," *IEEE Trans. Commun.*, vol. 69, pp. 4732–4750, Jul. 2021.

[40] D. Dardari, "Communicating with large intelligent surfaces: Fundamental limits and models," *IEEE J. Sel. Areas Commun.*, vol. 38, pp. 2526–2537, Nov. 2020.

[41] 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, Sep. 2020.

[42] A. Pizzo, L. Sanguinetti, and T. L. Marzetta, "Fourier plane-wave series expansion for holographic MIMO communications," *IEEE Trans. Wireless Commun.*, pp. 1–16, 2022, Early Access.

[43] O. 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, pp. 997–1001, May 2022.

[44] X. Hu, R. Deng, B. Di, H. Zhang, and L. Song, "Holographic beamforming for ultra massive MIMO with limited radiation amplitudes: How many quantized bits do we need?," *IEEE Commun. Lett.*, pp. 1–5, 2022, Early Access.

[45] R. Deng, B. Di, H. Zhang, and L. Song, "HDMA: Holographic-pattern division multiple access," *IEEE J. Sel. Areas Commun.*, vol. 40, pp. 1317–1332, Apr. 2022.

[46] L. Wei, C. Huang, G. C. Alexandropoulos, W. E. Sha, Z. Zhang, M. Debbah, and C. Yuen, "Multi-user holographic MIMO surfaces: Channel modeling and spectral efficiency analysis," *IEEE J. Sel. Topics Signal Process.*, vol. 16, pp. 1112–1124, Aug. 2022.

[47] C. Liu, Q. Ma, Z. J. Luo, Q. R. Hong, Q. Xiao, H. C. Zhang, L. Miao, W. M. Yu, Q. Cheng, L. Li, *et al.*, "A programmable diffractive deep neural network based on a digital-coding metasurface array," *Nat. Electron.*, vol. 5, pp. 113–122, Feb. 2022.

[48] N. Chamanara, Y. Vahabzadeh, and C. Caloz, "Stacked metasurface slab," in *Proc. Int. Congress Artificial Materials, Novel Wave Phenomena*, pp. 70–72, IEEE, Aug. 2018.

[49] Y. Hu, X. Luo, Y. Chen, Q. Liu, X. Li, Y. Wang, N. Liu, and H. Duan, "3D-integrated metasurfaces for full-colour holography," *Light: Science & Applications*, vol. 8, no. 1, pp. 1–9, 2019.

[50] C. Pfeiffer and A. Grbic, "Cascaded metasurfaces for complete phase and polarization control," *Applied Physics Lett.*, vol. 102, p. 231116, Jun. 2013.

[51] Y. Zhou, I. I. Kravchenko, H. Wang, J. R. Nolen, G. Gu, and J. Valentine, "Multilayer noninteracting dielectric metasurfaces for multiwavelength metaoptics," *Nano Lett.*, vol. 18, pp. 7529–7537, Nov. 2018.

[52] J. An, M. Di Renzo, M. Debbah, and C. Yuen, "Stacked intelligent metasurfaces for multiuser beamforming in the wave domain," in *Proc. IEEE Int. Conf. Commun. (ICC)*, pp. 1–6, Rome, Italy, May 2023.

[53] J. An, C. Xu, L. Wang, Y. Liu, L. Gan, and L. Hanzo, "Joint training of the superimposed direct and reflected links in reconfigurable intelligent surface assisted multiuser communications," *IEEE Trans. Green Commun. Netw.*, vol. 6, pp. 739–754, Jun. 2022.

[54] W. Tang, M. Z. Chen, X. Chen, J. Y. Dai, Y. Han, M. Di Renzo, Y. Zeng, S. Jin, Q. Cheng, and T. J. Cui, "Wireless communications with reconfigurable intelligent surface: Path loss modeling and experimental measurement," *IEEE Trans. Wireless Commun.*, vol. 20, pp. 421–439, Jan. 2021.

[55] A. L. Moustakas, G. C. Alexandropoulos, and M. Debbah, "Reconfigurable intelligent surfaces and capacity optimization: A large system analysis," *IEEE Trans. Wireless Commun.*, to appear, 2023.

[56] X. Cao, B. Yang, C. Huang, C. Yuen, M. Di Renzo, Z. Han, D. Niyato, H. V. Poor, and L. Hanzo, "AI-assisted MAC for reconfigurable intelligent-surface-aided wireless networks: Challenges and opportunities," *IEEE Commun. Mag.*, vol. 59, pp. 21–27, Jun. 2021.

[57] X. Cao, B. Yang, C. Huang, G. C. Alexandropoulos, C. Yuen, Z. Han, and H. V. Poor, "Massive access of static and mobile users via reconfigurable intelligent surfaces: Protocol design and performance analysis," *IEEE J. Sel. Areas Commun.*, vol. 40, pp. 1253–1269, Apr. 2022.

[58] J. An, L. Wang, C. Xu, L. Gan, and L. Hanzo, "Optimal pilot power based channel estimation improves the throughput of intelligent reflective surface assisted systems," *IEEE Trans. Veh. Technol.*, vol. 69, pp. 16202–16206, Dec. 2020.

[59] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *Nat.*, vol. 521, pp. 436–444, May 2015.

[60] X. Lin, Y. Rivenson, N. T. Yardimci, M. Veli, Y. Luo, M. Jarrahi, and A. Ozcan, "All-optical machine learning using diffractive deep neural networks," *Sci.*, vol. 361, pp. 1004–1008, Jul. 2018.

[61] X. Yao, K. Klyukin, W. Lu, M. Onen, S. Ryu, D. Kim, N. Emond, I. Waluyo, A. Hunt, J. A. Del Alamo, *et al.*, "Protonic solid-state electrochemical synapse for physical neural networks," *Nat. Commun.*, vol. 11, pp. 1–10, Jun. 2020.

[62] D. Tse and P. Viswanath, *Fundamentals of wireless communication*. Cambridge university press, 2005.

[63] T. J. Cui, M. Q. Qi, X. Wan, J. Zhao, and Q. Cheng, "Coding metamaterials, digital metamaterials and programmable metamaterials," *Light, Sci. Appl.*, vol. 3, pp. 1–9, Oct. 2014.

[64] E. Arbabi, A. Arbabi, S. M. Kamali, Y. Horie, M. Faraji-Dana, and A. Faraon, "MEMS-tunable dielectric metasurface lens," *Nat. Commun.*, vol. 9, pp. 1–9, Feb. 2018.

[65] L. Dai, B. Wang, M. Wang, X. Yang, J. Tan, S. Bi, S. Xu, F. Yang, Z. Chen, M. Di Renzo, C.-B. Chae, and L. Hanzo, "Reconfigurable intelligent surface-based wireless communications: Antenna design, prototyping, and experimental results," *IEEE Access*, vol. 8, pp. 45913–45923, Mar. 2020.

[66] T. S. Rappaport, G. R. MacCartney, M. K. Samimi, and S. Sun, "Wideband millimeter-wave propagation measurements and channel models for future wireless communication system design," *IEEE Trans. Commun.*, vol. 63, pp. 3029–3056, Sep. 2015.

[67] R. W. Heath, N. González-Prelcic, S. Rangan, W. Roh, and A. M. Sayeed, "An overview of signal processing techniques for millimeter wave MIMO systems," *IEEE J. Sel. Topics Signal Process.*, vol. 10, pp. 436–453, Apr. 2016.

[68] S. Basodi, C. Ji, H. Zhang, and Y. Pan, "Gradient amplification: An efficient way to train deep neural networks," *Big Data Mining and Analytics*, vol. 3, pp. 196–207, Sep. 2020.

[69] S. Zhang and R. Zhang, "Capacity characterization for intelligent reflecting surface aided MIMO communication," *IEEE J. Sel. Areas Commun.*, vol. 38, pp. 1823–1838, Aug. 2020.