

Observation Matrix Design for Densifying MIMO Channel Estimation via 2D Ice Filling

Zijian Zhang , *Graduate Student Member, IEEE*, and Mingyao Cui , *Graduate Student Member, IEEE*

Abstract—In recent years, densifying multiple-input multiple-output (MIMO) has attracted much attention from the communication community. Thanks to the subwavelength antenna spacing, the strong correlations among densifying antennas provide sufficient prior knowledge about channel state information (CSI). This inspires the careful design of observation matrices (e.g., transmit precoders and receive combiners), that exploits the CSI prior knowledge, to boost channel estimation performance. Aligned with this vision, this work proposes to jointly design the combiners and precoders by maximizing the mutual information between the received pilots and densifying MIMO channels. A two-dimensional ice-filling (2DIF) algorithm is proposed to efficiently accomplish this objective. The algorithm is motivated by the fact that the eigenspace of MIMO channel covariance can be decoupled into two sub-eigenspaces, which are associated with the correlations of transmitter antennas and receiver antennas, respectively. By properly setting the precoder and the combiner as the eigenvectors from these two sub-eigenspaces, the 2DIF promises to generate near-optimal observation matrices. Moreover, we further extend the 2DIF method to the popular hybrid combining systems, where a two-stage 2DIF (TS-2DIF) algorithm is developed to handle the analog combining circuits realized by phase shifters. Simulation results demonstrate that, compared to the state-of-the-art schemes, the proposed 2DIF and TS-2DIF methods can achieve superior channel estimation accuracy.

Index Terms—Channel estimation, densifying MIMO, dense array systems (DAS), observation matrix design.

I. INTRODUCTION

IN recent years, densifying multiple-input multiple-output (MIMO) has attracted considerable attention from the wireless communication community [2], [3], [4], [5], [6], [7]. Different from the conventional MIMO whose antennas are

usually spaced of half wavelength $\lambda/2$, the antenna spacing of densifying MIMO is much smaller, such as $\lambda/6$ [8], $\lambda/8$ [9], $\lambda/10$ [10], or even $\lambda/23$ [11]. By densely arranging massive subwavelength-spaced antennas in a compact space, densifying MIMO promises to realize the ultimate control of the radiated/received electromagnetic waves on limited apertures. To this end, many dense-antenna transceiver architectures have emerged, such as holographic MIMO (H-MIMO) [12], holographic reconfigurable surfaces (RHSs) [3], continuous-aperture MIMO (CAP-MIMO) [4], superdirective antenna arrays [5], reconfigurable intelligent surfaces (RISs) [13], and fluid antenna systems (FASs) [7]. Particularly, in high-frequency (e.g., millimeter-wave or terahertz) communications, densifying MIMO systems are common due to their shorter wavelengths, reduced grating lobes, and enhanced beamforming precision [2]. Utilizing the extensive channel observations facilitated by a multitude of antennas, densifying MIMO is anticipated to achieve significant array gains and multiplexing-diversity gains [14], [15], [16], [17]. Furthermore, densifying MIMO can mitigate the effects of grating lobes and offer enhanced performance for large oblique angles of incidence [18]. Some studies have also highlighted their capabilities to realize super-directivity [5], [19] or super-bandwidth [20] in wireless transmissions.

Enabled by their phase shifters and radio frequency (RF) chains, the transmission performance of MIMO is determined by the constructive precoders/combiners at transceivers [21]. To implement effective precoding/combining, an indispensable technology for MIMO systems is the acquisition of channel state information (CSI) [22], [23]. To date, numerous technologies have been proposed to estimate the channels of classical MIMO systems. For example, when the available pilot length exceeds the number of antennas, some classical estimators [24], such as the least square (LS) estimator and the minimum mean square error (MMSE) estimator, can be used to recover MIMO channels in a non-parametric way. By exploiting the channel sparsity in the angular domain, compressed sensing (CS)-based channel estimators can enhance the estimation accuracy and reduce the pilot overhead [25], [26]. Relevant techniques include the orthogonal matching pursuit (OMP)-based estimator [27] and the approximate message passing (AMP)-based estimator [25], [28]. Additionally, some deep learning (DL) approaches, which involve training neural networks based on channel datasets, are utilized to realize data-driven channel estimation in MIMO systems [29], [30], [31].

Received 14 July 2025; revised 20 December 2025; accepted 19 January 2026. Date of publication 23 January 2026; date of current version 27 February 2026. This work was supported by the National Natural Science Foundation of China under Grant 624B2123. An earlier version of this paper was presented in part at the IEEE ICC'25, Montreal, Canada, 2025 [DOI:10.1109/ICC52391.2025.11161466]. Compared with the conference version [1], this journal version provides theoretical analysis, generalized channel modeling, and a new algorithm for hybrid MIMO architectures. The associate editor coordinating the review of this article and approving it for publication was Rodrigo C. de Lamare. (*Corresponding author: Mingyao Cui.*)

Zijian Zhang is with the Department of Electronic Engineering, Tsinghua University, Beijing 100084, China (e-mail: zhangzij15@tsinghua.org.cn).

Mingyao Cui is with the Department of Electrical and Electronic Engineering, The University of Hong Kong, Hong Kong SAR 999077, China (e-mail: mycui@eee.hku.hk).

Digital Object Identifier 10.1109/TSP.2026.3657342

Although many channel estimators in the literature can be adopted in densifying MIMO systems, they often exhibit a non-negligible performance gap compared to the optimal estimator [24]. This is because most existing estimators overlook the strong correlations among densifying MIMO antennas. Specifically, since the antenna spacing of densifying MIMO is very small, the channels associated with close-by antennas are spatially similar [17]. Besides, the circuit mutual coupling induces signal interactions between adjacent antennas, which will influence the channel correlation in densifying MIMO systems [8], [9], [10], [11], [32]. These facts lead to the highly structured covariance matrices of densifying MIMO channels. Existing works have revealed that, such an informative covariance matrix can provide appreciable prior knowledge for the specific design of observation matrices (e.g., combiners and precoders) in channel estimation, thus significantly improving the accuracy of CSI acquisition [33], [34], [35].

To exploit the strong channel correlations for improved CSI acquisition, our prior work [36] proposes an ice filling (IF) based observation matrix design in dense array systems, which is inspired by the idea of Gaussian Process Regression (GPR). By maximizing the mutual information (MI) between the received pilot and wireless channel, which is characterized by the channel covariance matrix, the IF algorithm can sequentially produce the observation vectors of receivers in a pilot-by-pilot manner. Through optimizing pilot allocation to the channel covariance eigenvectors, this method works like filling ice blocks onto different orthogonal channels. Then, the designed observation matrix is shown to have near-optimal channel estimation performance in DASs, which achieves much higher estimation accuracy compared to the state-of-the-art schemes [36]. Nevertheless, despite its ability to exploit the channel covariance, IF method is only feasible to design the *vector-form receive combiner* in single-input multiple-output (SIMO) system with a single-antenna transmitter and a single-RF-chain receiver [36]. For a general densifying MIMO system with multiple antennas and multiple RF chains at both transceivers, *the receive and transmit channel covariance pair*, as well as the *matrix-form receive combiner and transmit precoder pair* are coupled together. As the IF scheme fails to tackle these couplings, it is far from optimal for densifying MIMO. To the best of our knowledge, the full exploitation of densifying MIMO's channel covariance for designing observation matrices is still an unaddressed challenge.

To fill in this gap, this work generalizes the IF scheme to a two-dimensional ice filling (2DIF) scheme, whose core idea is to design precoders and combiners by decoupling the coupled channel covariance matrices in their eigenspace. Our key contributions and findings are summarized as follows.

- **Generalized framework of observation matrix design:** Inspired by the IF algorithm, we apply the technique of GPR into densifying MIMO channel estimation. Our key idea is to maximize the MI between the received pilots and the MIMO channel by jointly optimizing the receive combiners and transmit precoders. The formulated observation matrix design problem is shown to be a generalization of that discussed in IF [36] because of the

consideration of practical MIMO systems with multi-RF-chain receivers. To be specific, the receiver-side channel covariance considered by IF is generalized to a MIMO channel covariance, which relies on the correlations at both sides of the transceivers. Moreover, the observation matrix is no longer a vector-form combiner, but the Kronecker product of the matrix-form combiner and precoder. These properties fundamentally distinguish our design from the IF scheme.

- **2DIF based observation matrix design:** To overcome the design challenges imposed by the coupling of matrix-form combiners and precoders in observation matrices, a 2DIF based observation matrix design is developed. The proposed design employs a greedy method to jointly produce the combiners and precoders in a block-by-block way. Concretely, we first prove that the eigenspace of the channel covariance can be decoupled into two sub-eigenspaces, which are associated with the correlations of transmitter antennas and receiver antennas, respectively. Then, utilizing the eigenspace invariance, we show that the near-optimal observation matrix can be obtained by properly setting the precoder and the combiner as the eigenvectors from these two sub-eigenspaces, which can be realized by a linear search algorithm. Besides, we also provide an intuitive and insightful explanation for 2DIF to clarify its physical significance. Similar to the water-filling precoding which maximizes the MIMO capacity, the implementation of 2DIF can be viewed as a two-dimensional ice-filling process.
- **TS-2DIF based observation matrix design:** The proposed 2DIF method requires that the amplitude of each receiver antenna can be controlled independently. However, for many hybrid MIMO structures, only the phase shifts of analog combiners can be reconfigured, thus the proposed 2DIF cannot be directly adopted in these scenarios. To address this issue, we propose the two-stage 2DIF (TS-2DIF) method. Different from the exiting works on hybrid MIMO beamforming that precoders and combiners can be designed independently, in our considered observation matrix design, the beamformer and the combiner are coupled due to a Kronecker-product structure. To overcome this challenge, the proposed TS-2DIF alternately optimizes the analog combiner, the digital combiner, and the precoder at the transceivers to approach the ideal observation matrix.

The rest of this paper is organized as follows. In Section II, the system model and problem formulation are introduced. In Section III, the proposed 2DIF based observation matrix design for channel estimation is illustrated. In Section IV, the proposed TS-2DIF based observation matrix design is provided. In Section V, the computational complexities of the proposed methods are analyzed, and the kernel selection is discussed. In Section VI, simulations are carried out to verify the effectiveness of the proposed schemes. In Section VII, conclusions are drawn.

Notation: $[\cdot]^T$, $[\cdot]^H$, $[\cdot]^*$, and $[\cdot]^{-1}$ denote the transpose, conjugate-transpose, conjugate, and inverse operations, respectively; $\|\cdot\|$ denotes the l_2 -norm operation; $\|\cdot\|_F$ denotes the

Frobenius-norm operation; $z(i)$ denotes the i -th entry of vector \mathbf{z} ; $\mathbf{Z}(i, j)$, $\mathbf{Z}(j, :)$ and $\mathbf{Z}(:, j)$ denote the (i, j) -th entry, the j -th row, and the j -th column of matrix \mathbf{Z} , respectively; $\text{Tr}(\cdot)$ denotes the trace of its argument; $\text{diag}(\cdot)$ and $\text{blkdiag}(\cdot)$ are the diagonal and the block-diagonal operations, respectively; $\mathbb{E}(\cdot)$ is the expectation operator; $\Re\{\cdot\}$ denotes the real part of the argument; $\ln(\cdot)$ denotes the natural logarithm of its argument; $\mathcal{CN}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ denotes the complex Gaussian distribution with mean $\boldsymbol{\mu}$ and covariance $\boldsymbol{\Sigma}$; $\mathcal{U}(a, b)$ denotes the uniform distribution between a and b ; \mathbf{I}_L is an $L \times L$ identity matrix; $\mathbf{1}_L$ is an all-one vector or matrix with dimension L ; and $\mathbf{0}_L$ is a zero vector or matrix with dimension L .

II. SYSTEM MODEL AND PROBLEM FORMULATION

A. Transceiver Model

This paper considers the uplink channel estimation of a densifying MIMO system, consisting of an N_R -antenna base station (BS) equipped with N_{RF} RF chains and an N_T -antenna user. The antennas at transceivers are densely arranged with sub-wavelength antenna spacing d . We define $\mathbf{H} \in \mathbb{C}^{N_R \times N_T}$ as the wireless channel and Q as the number of transmit pilots within a coherence-time frame. The received signal $\mathbf{y}_q \in \mathbb{C}^{N_{\text{RF}}}$ at the BS in timeslot q is modeled as

$$\begin{aligned} \mathbf{y}_q &= \mathbf{W}_q^H \mathbf{H} \mathbf{v}_q s_q + \mathbf{W}_q^H \mathbf{z}_q \\ &= (\mathbf{v}_q^T \otimes \mathbf{W}_q^H) \mathbf{h} s_q + \mathbf{W}_q^H \mathbf{z}_q, \end{aligned} \quad (1)$$

where $\mathbf{h} \equiv \text{vec}(\mathbf{H})$ is defined as the vectorized channel matrix, s_q the pilot symbol, and $\mathbf{z}_q \sim \mathcal{CN}(\mathbf{0}_M, \sigma^2 \mathbf{I}_M)$ is the additive white Gaussian noise (AWGN). Vector $\mathbf{v}_q \in \mathbb{C}^{N_T}$ denotes the precoder at the user. For the transmitter, the user equipment typically employs a fully digital precoder with a moderate number of antennas, N_T . Thereby, the coefficient of \mathbf{v}_q can be freely configured as long as the power constraint is satisfied: $\|\mathbf{v}_q\|^2 = P$, with P the per-pilot transmit power. For the receiver, $\mathbf{W}_q := \mathbf{A}_q \mathbf{D}_q \in \mathbb{C}^{N_R \times N_{\text{RF}}}$ is the hybrid combiner at the BS, with $\mathbf{A}_q \in \mathbb{C}^{N_R \times N_{\text{RF}}}$ and $\mathbf{D}_q \in \mathbb{C}^{N_{\text{RF}} \times N_{\text{RF}}}$ being the analog and digital combiners, respectively. As presented in Fig. 1, we focus on two typical implementations of the analog combiners: the amplitude-and-phase controllable combiner and the phase-only controllable combiner. The former deploys one phase shifter and one low-noise-amplifier (LNA) between each antenna-to-RF chain link. In this architecture, both the amplitude and phase of the elements of \mathbf{A}_q are adjustable, and thus the coefficients of \mathbf{W}_q can be freely controlled. The latter, however, employs one phase shifter to connect each antenna-to-RF chain link, while the signals aggregated on each RF chain are jointly processed by a global LNA. In this context, only the phases of \mathbf{A}_q are adjustable, which imposes a structural constraint on the feasible set of $\mathbf{W}_q = \mathbf{A}_q \mathbf{D}_q$.

Without loss of generality, we assume $s_q = 1, \forall q \in \{1, \dots, Q\}$. Considering the total Q timeslots for pilot transmission, we arrive at

$$\mathbf{y} = \mathbf{X}^H \mathbf{h} + \mathbf{z}, \quad (2)$$

where $\mathbf{y} := [\mathbf{y}_1^T, \dots, \mathbf{y}_Q^T]^T$, $\mathbf{z} := [\mathbf{z}_1^H \mathbf{W}_1, \dots, \mathbf{z}_Q^H \mathbf{W}_Q]^H$, $\mathbf{X} = [\mathbf{X}_1, \dots, \mathbf{X}_Q]$, and $\mathbf{X}_q := \mathbf{v}_q^* \otimes \mathbf{W}_q$ is defined as the

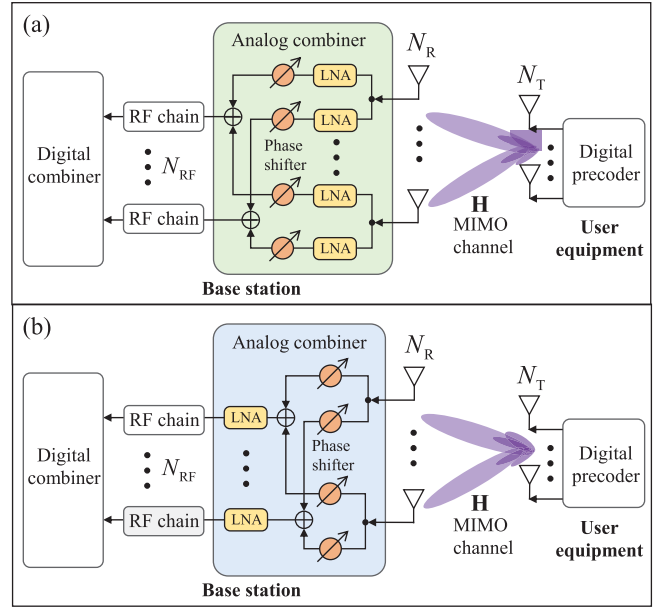


Fig. 1. An illustration of hybrid analog and digital MIMO architectures, where the analog combiner of BS can adopt either: (a) the amplitude-and-phase controllable structure; or (b) the phase-only controllable structure.

observation matrix for each pilot. This paper aims at accurately estimating \mathbf{h} from \mathbf{y} by jointly designing combiners $\{\mathbf{W}_q\}_{q=1}^Q$ and precoders $\{\mathbf{v}_q\}_{q=1}^Q$.

B. Channel Model

We consider the general correlated Rayleigh-fading channel model, which simultaneously takes the spatial correlation, spherical wavefront, and electromagnetic mutual coupling into account. Following this well-known model [37], [38], [39], the channel \mathbf{H} between the BS and the user is expressed as:

$$\mathbf{H} = \mathbf{C}_{\text{rx}}^{1/2} \mathbf{R}_{\text{rx}}^{1/2} \mathbf{H}_{\text{iid}} \mathbf{R}_{\text{tx}}^{1/2} \mathbf{C}_{\text{tx}}^{1/2}, \quad (3)$$

where $\mathbf{R}_{\text{rx}} \in \mathbb{C}^{N_R \times N_R}$ and $\mathbf{R}_{\text{tx}} \in \mathbb{C}^{N_T \times N_T}$ are the spatial correlation matrices at the BS and the user, respectively; $\mathbf{C}_{\text{rx}} \in \mathbb{C}^{N_R \times N_R}$ and $\mathbf{C}_{\text{tx}} \in \mathbb{C}^{N_T \times N_T}$ characterize the electromagnetic mutual coupling among the BS antennas and the user antennas, respectively; and \mathbf{H}_{iid} is a complex Gaussian matrix in which each elements are i.i.d with zero mean and unit variance. This channel model is widely recognized by the studies on correlated MIMO channels, also adopted by commercial channel simulators, such as Matlab Antenna toolbox [40].

Example 1 (Example of settings): Note that, the settings of \mathbf{R}_{rx} , \mathbf{R}_{tx} , \mathbf{C}_{rx} , and \mathbf{C}_{tx} depend on the specific transmission scenarios in practice. As a typical example, the (m, n) -th entry of \mathbf{R}_{rx} is usually modeled as [12], [15], [41]:

$$\mathbf{R}_{\text{rx}}(m, n) = \int \int_{-\pi/2}^{+\pi/2} f_{\text{rx}}(\varphi, \theta) e^{j\mathbf{k}(\varphi, \theta)^T (\mathbf{r}_m - \mathbf{r}_n)} d\theta d\varphi, \quad (4)$$

where φ and θ denote the azimuth angle and elevation angle of incident signals, respectively; $f_{\text{rx}}(\varphi, \theta)$ is the scenario-dependent spatial-scattering function,

which describes the angular multipath distribution and the directivity gain of receiver antennas in practice; $\mathbf{k}(\varphi, \theta) := \frac{2\pi}{\lambda} (\cos \theta \cos \varphi, \cos \theta \sin \varphi, \sin \theta)^T$ is the wave vector with λ being the wavelength; and $\mathbf{r}_m \in \mathbb{R}^3$ is the location of the m -th receiver antenna. To characterize the mutual coupling, \mathbf{C}_{rx} is modeled as [37], [42]:

$$\mathbf{C}_{\text{rx}} = (\mathbf{Z}_{\text{rx}} + R_{\text{rx}}\mathbf{I})^{-1}, \quad (5)$$

where \mathbf{Z}_{rx} is the mutual impedance matrix and $R_{\text{rx}} > 0$ is the dissipation resistance of antennas. Given the specific array parameters, \mathbf{Z}_{rx} and R_{rx} can be calculated by the well-known induced electromagnetic fields (EMF) method [43]. Similar settings can be adopted to generate \mathbf{R}_{tx} and \mathbf{C}_{tx} .

C. Problem Formulation

As antennas of densifying MIMO are packed within a small area, the channels across close-by antennas are strongly correlated. Define the covariance of vectorized channel $\mathbf{h} = \text{vec}(\mathbf{H})$ as $\mathbf{E}(\mathbf{h}\mathbf{h}^H) = \mathbf{\Sigma}_{\mathbf{h}} \in \mathbb{C}^{N_{\text{R}}N_{\text{T}} \times N_{\text{R}}N_{\text{T}}}$, which is also called the *kernel* of channel. The high-correlation property of wireless channel indicates that the kernel $\mathbf{\Sigma}_{\mathbf{h}}$ can provide prior knowledge to achieve accurate channel estimation [44], [45], [46], [47]. To this end, we follow the idea of GPR to design the estimator and the observation matrix. According to the model in (3), the channel is sampled from the Gaussian process $\mathcal{CN}(\mathbf{0}_{N_{\text{R}}N_{\text{T}}}, \mathbf{\Sigma}_{\mathbf{h}})$. The joint probability distribution of \mathbf{h} and \mathbf{y} then satisfies

$$\begin{bmatrix} \mathbf{h} \\ \mathbf{y} \end{bmatrix} \sim \mathcal{CN} \left(\begin{bmatrix} \mathbf{0}_{N_{\text{R}}N_{\text{T}}} \\ \mathbf{0}_{N_{\text{R}}FQ} \end{bmatrix}, \begin{bmatrix} \mathbf{\Sigma}_{\mathbf{h}} & \mathbf{\Sigma}_{\mathbf{h}}\mathbf{X} \\ \mathbf{X}^H\mathbf{\Sigma}_{\mathbf{h}} & \mathbf{X}^H\mathbf{\Sigma}_{\mathbf{h}}\mathbf{X} + \mathbf{\Xi} \end{bmatrix} \right), \quad (6)$$

where $\mathbf{\Xi} = \sigma^2 \text{blkdiag}(\mathbf{W}_1^H\mathbf{W}_1, \dots, \mathbf{W}_Q^H\mathbf{W}_Q)$ represents the covariance matrix of the noise \mathbf{z} . Then, the posterior mean and the posterior covariance of \mathbf{h} are expressed as

$$\boldsymbol{\mu}_{\mathbf{h}|\mathbf{y}} = \mathbf{\Sigma}_{\mathbf{h}}\mathbf{X}(\mathbf{X}^H\mathbf{\Sigma}_{\mathbf{h}}\mathbf{X} + \mathbf{\Xi})^{-1}\mathbf{y}, \quad (7)$$

$$\mathbf{\Sigma}_{\mathbf{h}|\mathbf{y}} = \mathbf{\Sigma}_{\mathbf{h}} - \mathbf{\Sigma}_{\mathbf{h}}\mathbf{X}(\mathbf{X}^H\mathbf{\Sigma}_{\mathbf{h}}\mathbf{X} + \mathbf{\Xi})^{-1}\mathbf{X}^H\mathbf{\Sigma}_{\mathbf{h}}. \quad (8)$$

The posterior mean $\boldsymbol{\mu}_{\mathbf{h}|\mathbf{y}}$ is exactly the adopted channel estimator, which is equivalent to the linear MMSE (LMMSE) estimator. Besides, the posterior covariance $\mathbf{\Sigma}_{\mathbf{h}|\mathbf{y}}$ characterizes the estimation error, which is largely dependent on the observation matrix \mathbf{X} . This dependence indicates that well-designed combiners and precoders, $\{\mathbf{W}_q\}_{q=1}^Q$ and $\{\mathbf{v}_q\}_{q=1}^Q$, can substantially reduce the channel estimation error. Motivated by this fact, GPR attempts to produce the observation matrices to gain as much information of \mathbf{h} as possible from the received signal \mathbf{y} . Henceforth, our objective is to find $\{\mathbf{W}_q\}_{q=1}^Q$ and $\{\mathbf{v}_q\}_{q=1}^Q$ that maximize the MI¹ between \mathbf{y} and \mathbf{h} , which is formulated as:

$$\max_{\mathbf{W}, \|\mathbf{v}_q\|^2=P} I(\mathbf{y}; \mathbf{h}) = \log \det (\mathbf{I}_{N_{\text{RF}}Q} + \mathbf{\Xi}^{-1}\mathbf{X}^H\mathbf{\Sigma}_{\mathbf{h}}\mathbf{X}). \quad (9)$$

¹According to the information-theoretic properties [48], the maximization of MI is asymptotically the minimization of Cramer-Rao bound.

where \mathcal{W} stands for the feasible set of hybrid combiners depending on the receiver hardware. In the subsequent sections, we will first elaborate on the observation matrix design while considering the amplitude-and-phase controllable analog combiners in Section III. Then, our design will be extended to the case of phase-only controllable analog combiners in Section IV.

III. PROPOSED TWO-DIMENSIONAL ICE FILLING (2DIF) BASED OBSERVATION MATRIX DESIGN

In this section, we consider the ideal case when the amplitudes and phases of all precoder and combiner coefficients are adjustable, as shown in Fig. 1(a). In this context, $\{\mathbf{W}_q\}_{q=1}^Q$ can be freely configured and $\{\mathbf{v}_q\}_{q=1}^Q$ should satisfy the transmit power constraints $\|\mathbf{v}_q\|^2 = P$ for all $q \in \{1, \dots, Q\}$.

A. Precoder/Combiner Design Using Greedy Method

Observing problem (9), one can find that the MI $I(\mathbf{y}; \mathbf{h})$ is non-concave with respect to (w.r.t) the overall observation matrix \mathbf{X} . Besides, due to the coupled term $\mathbf{v}_q^T \otimes \mathbf{W}_q^H$ in \mathbf{X} and the colored-noise covariance matrix $\mathbf{\Xi}$ introduced by combiners $\{\mathbf{W}_q\}_{q=1}^Q$, the global optimal solution to (9) is hard to obtain. To address this issue, we adopt a greedy method to generate $\{\mathbf{W}_q\}_{q=1}^Q$ and $\{\mathbf{v}_q\}_{q=1}^Q$ in a pilot-by-pilot manner. Specifically, we define $\bar{\mathbf{X}}_t = [\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_t]$ as the overall observation matrix for timeslots $1 \sim t$, where $t \leq Q$. Let $\bar{\mathbf{y}}_t = \bar{\mathbf{X}}_t^H \mathbf{h} + \bar{\mathbf{z}}_t$ denote the corresponding received signal, wherein $\bar{\mathbf{y}}_t = [y_1^T, y_2^T, \dots, y_t^T]^T$ and $\bar{\mathbf{z}}_t := [z_1^H \mathbf{W}_1, \dots, z_t^H \mathbf{W}_t]^H$. Given the current observation matrices $\{\mathbf{W}_q\}_{q=1}^t$ and vectors $\{\mathbf{v}_q\}_{q=1}^t$ in the first t timeslots, our greedy strategy aims to find the combiner \mathbf{W}_{t+1} and the precoder \mathbf{v}_{t+1} in the next timeslot, which maximize the MI increment from timeslot t to $t+1$:

$$\max_{\mathbf{W}_{t+1}, \mathbf{v}_{t+1}} \Delta I_{t+1} := I(\bar{\mathbf{y}}_{t+1}; \mathbf{h}) - I(\bar{\mathbf{y}}_t; \mathbf{h}). \quad (10)$$

For clarity, we summarize the proposed design strategy in **Algorithm 1**, and the sequential designs of $\{\mathbf{W}_q\}_{q=1}^Q$ and $\{\mathbf{v}_q\}_{q=1}^Q$ are illustrated as follows.

1) *When $t = 1$* : For ease of understanding, we begin with handling the first timeslot, i.e., $t = 1$. In this context, problem (10) can be rewritten as

$$\max_{\|\mathbf{v}_1\|^2=P, \mathbf{W}_1} I(\mathbf{y}_1; \mathbf{h}) \quad (11)$$

where the mutual information $I(\mathbf{y}_1; \mathbf{h})$ is given by

$$I(\mathbf{y}_1; \mathbf{h}) = \log \det \left(\mathbf{I}_{N_{\text{RF}}} + \frac{1}{\sigma^2} (\mathbf{W}_1^H \mathbf{W}_1)^{-1} \times (\mathbf{v}_1^T \otimes \mathbf{W}_1^H) \mathbf{\Sigma}_{\mathbf{h}} (\mathbf{v}_1^* \otimes \mathbf{W}_1) \right). \quad (12)$$

Since the reformulated problem (11) is still intricate, we seek to simplify it using the following lemmas.

Lemma 1: For the correlated Rayleigh-fading model in (3), the covariance of the vectored channel $\mathbf{\Sigma}_{\mathbf{h}}$ can be rewritten as the form of a Kronecker product of two kernels, i.e.,

$$\mathbf{\Sigma}_{\mathbf{h}} = \mathbf{\Sigma}_{\text{T}} \otimes \mathbf{\Sigma}_{\text{R}}, \quad (13)$$

Algorithm 1 2DIF Based Combiner and Precoder Design**Input:** Number of pilots Q , kernel $\Sigma_{\mathbf{h}}$.**Output:** Designed precoders $\{\mathbf{v}_q^{\text{opt}}\}_{q=1}^Q$ and combiners $\{\mathbf{W}_q^{\text{opt}}\}_{q=1}^Q$.

- 1: Rewrite kernel as $\Sigma_{\mathbf{h}} = \Sigma_{\mathbf{T}} \otimes \Sigma_{\mathbf{R}}$
- 2: Find the eigenvectors $[\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_{N_{\mathbf{T}}}]$ and the corresponding eigenvalues $[\alpha_1, \alpha_2, \dots, \alpha_{N_{\mathbf{T}}}]$ of $\Sigma_{\mathbf{T}}$
- 3: Find the eigenvectors $[\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_{N_{\mathbf{R}}}]$ and the corresponding eigenvalues $[\beta_1, \beta_2, \dots, \beta_{N_{\mathbf{R}}}]$ of $\Sigma_{\mathbf{R}}$
- 4: Initialize: $[\lambda_{1,1}^0, \lambda_{1,2}^0, \dots, \lambda_{N_{\mathbf{T}}, N_{\mathbf{R}}}^0] = [\alpha_1 \beta_1, \alpha_1 \beta_2, \dots, \alpha_{N_{\mathbf{T}}} \beta_{N_{\mathbf{R}}}]$
- 5: **for** $t = 0, \dots, Q - 1$ **do**
- 6: Find the optimal $n_{\mathbf{T}}^{\text{opt}}$ and $\{n_{\mathbf{R},k}^{\text{opt}}\}_{k=1}^{N_{\mathbf{RF}}}$ via **Algorithm 2**
- 7: Eigenvector-assignment: $\mathbf{v}_{t+1}^{\text{opt}} = \sqrt{P} \mathbf{a}_{n_{\mathbf{T}}^{\text{opt}}}^*$ and $\mathbf{W}_{t+1}^{\text{opt}} = [\mathbf{b}_{\mathbf{R},1}^{\text{opt}}, \dots, \mathbf{b}_{\mathbf{R}, N_{\mathbf{RF}}}^{\text{opt}}]$
- 8: Eigenvalue-update for all $n_{\mathbf{T}} \in \{1, \dots, N_{\mathbf{T}}\}$ and $n_{\mathbf{R}} \in \{1, \dots, N_{\mathbf{R}}\}$ via
- 9:
$$\lambda_{n_{\mathbf{T}}, n_{\mathbf{R}}}^{t+1} = \begin{cases} \frac{\lambda_{n_{\mathbf{T}}, n_{\mathbf{R}}}^t \sigma^2}{P \lambda_{n_{\mathbf{T}}, n_{\mathbf{R}}}^t + \sigma^2}, & n_{\mathbf{T}} = n_{\mathbf{T}}^{\text{opt}} \ \& \ n_{\mathbf{R}} \in \{n_{\mathbf{R},k}^{\text{opt}}\}_{k=1}^{N_{\mathbf{RF}}} \\ \lambda_{n_{\mathbf{T}}, n_{\mathbf{R}}}^t, & \text{else.} \end{cases}$$
- 10: **end for**
- 11: **return** Designed precoders $\{\mathbf{v}_q^{\text{opt}}\}_{q=1}^Q$ and combiners $\{\mathbf{W}_q^{\text{opt}}\}_{q=1}^Q$ for channel estimation.

where $\Sigma_{\mathbf{T}} := (\mathbf{C}_{\text{tx}}^{1/2})^T \mathbf{R}_{\text{tx}}^* (\mathbf{C}_{\text{tx}}^{1/2})^* \in \mathbb{C}^{N_{\mathbf{T}} \times N_{\mathbf{T}}}$ and $\Sigma_{\mathbf{R}} := \mathbf{C}_{\text{rx}}^{1/2} \mathbf{R}_{\text{rx}} (\mathbf{C}_{\text{rx}}^{1/2})^H \in \mathbb{C}^{N_{\mathbf{R}} \times N_{\mathbf{R}}}$ characterize the channel correlation among transmit antennas and that among receive antennas, respectively.

Proof: (See Appendix A). \blacksquare

Lemma 2: Introducing the orthogonality constraints $\mathbf{W}_q^H \mathbf{W}_q = \mathbf{I}_{N_{\mathbf{RF}}}$ for all $q \in \{1, \dots, Q\}$ does not influence the optimal value of MI $I(\mathbf{y}; \mathbf{h})$ in (9).

Proof: (See Appendix B). \blacksquare

Utilizing **Lemma 1** and **Lemma 2**, problem (11) can be equivalently rewritten as

$$\begin{aligned} \max_{\mathbf{v}_1, \mathbf{W}_1} I(\mathbf{y}_1; \mathbf{h}) &= \log \det \left(\mathbf{I}_{N_{\mathbf{RF}}} + \frac{\mathbf{v}_1^T \Sigma_{\mathbf{T}} \mathbf{v}_1^*}{\sigma^2} \mathbf{W}_1^H \Sigma_{\mathbf{R}} \mathbf{W}_1 \right) \\ \text{s.t. } \|\mathbf{v}_1\|^2 &= P, \\ \mathbf{W}_1^H \mathbf{W}_1 &= \mathbf{I}_{N_{\mathbf{RF}}}, \end{aligned} \quad (14)$$

where the reorganized objective function $I(\mathbf{y}_1; \mathbf{h})$ is obtained by substituting $\Sigma_{\mathbf{h}} = \Sigma_{\mathbf{T}} \otimes \Sigma_{\mathbf{R}}$ and $\mathbf{W}_1^H \mathbf{W}_1 = \mathbf{I}_{N_{\mathbf{RF}}}$ into (11). Observing (14), one can prove with ease that the optimal \mathbf{v}_1^* is the eigenvector of $\Sigma_{\mathbf{T}}$ associated with its largest eigenvalue, and the optimal \mathbf{W}_1 is composed of the eigenvectors of $\Sigma_{\mathbf{R}}$ associated with its top $N_{\mathbf{RF}}$ eigenvalues. Define the eigenvalue decompositions (EVDs) $\Sigma_{\mathbf{T}} = \mathbf{U}_{\mathbf{T}} \Lambda_{\mathbf{T}} \mathbf{U}_{\mathbf{T}}^H$ and $\Sigma_{\mathbf{R}} = \mathbf{U}_{\mathbf{R}} \Lambda_{\mathbf{R}} \mathbf{U}_{\mathbf{R}}^H$, wherein the eigenvalues in $\Lambda_{\mathbf{T}} = \text{diag}(\alpha_1, \dots, \alpha_{N_{\mathbf{T}}})$ and $\Lambda_{\mathbf{R}} = \text{diag}(\beta_1, \dots, \beta_{N_{\mathbf{R}}})$ are arranged in descending order. In this way, the maximum MI in (14) can be derived as

$$\max_{\mathbf{v}_1, \mathbf{W}_1} I(\mathbf{y}_1; \mathbf{h}) = \sum_{n=1}^{N_{\mathbf{RF}}} \log_2 \left(1 + \frac{P \alpha_1 \beta_n}{\sigma^2} \right), \quad (15)$$

and its achievable solution are expressed as

$$\mathbf{v}_1^{\text{opt}} = \sqrt{P} \mathbf{U}_{\mathbf{T}}^*(:, 1) \text{ and } \mathbf{W}_1^{\text{opt}} = \mathbf{U}_{\mathbf{R}}(:, [1, \dots, N_{\mathbf{RF}}]). \quad (16)$$

Equation (16) can be viewed as the optimal initialization settings for our proposed observation matrix design.

2) *From t to $t+1$:* Given the optimized precoder \mathbf{v}_1 and combiner \mathbf{W}_1 at timeslot $t=1$, we then consider designing $\{\mathbf{v}_q\}_{q=2}^Q$ and $\{\mathbf{W}_q\}_{q=2}^Q$ in a sequential manner. By invoking the principle of recursion, we only need to address the design of precoder \mathbf{v}_{t+1} and combiner \mathbf{W}_{t+1} with given $\{\mathbf{W}_q\}_{q=1}^t$ and $\{\mathbf{v}_q\}_{q=1}^t$. As proved in Appendix C, the MI increment, ΔI_{t+1} , can be equivalently rewritten as

$$\Delta I_{t+1} = \log_2 \det \left(\mathbf{I}_{N_{\mathbf{RF}}} + \frac{1}{\sigma^2} \mathbf{X}_{t+1}^H \Sigma_t \mathbf{X}_{t+1} \right), \quad (17)$$

wherein $\mathbf{X}_{t+1} := \mathbf{v}_{t+1}^* \otimes \mathbf{W}_{t+1}$ refers to the Kronecker-constrained observation matrix, and

$$\Sigma_t = \Sigma_{\mathbf{h}} - \Sigma_{\mathbf{h}} \bar{\mathbf{X}}_t (\bar{\mathbf{X}}_t^H \Sigma_{\mathbf{h}} \bar{\mathbf{X}}_t + \sigma^2 \mathbf{I}_{N_{\mathbf{RF}}})^{-1} \bar{\mathbf{X}}_t^H \Sigma_{\mathbf{h}} \quad (18)$$

is the posterior kernel of channel \mathbf{h} given the observation $\bar{\mathbf{y}}_t$. In particular, we have $\Sigma_0 := \Sigma_{\mathbf{h}}$. Utilizing **Lemma 1** and **Lemma 2**, the optimal precoder \mathbf{v}_{t+1} and combiner \mathbf{W}_{t+1} at the $(t+1)$ -th timeslot can be obtained by solving

$$\begin{aligned} \max_{\mathbf{v}_{t+1}, \mathbf{W}_{t+1}} \log \det \left(\mathbf{I}_{N_{\mathbf{RF}}} + \frac{1}{\sigma^2} (\mathbf{v}_{t+1}^T \otimes \mathbf{W}_{t+1}^H) \Sigma_t (\mathbf{v}_{t+1}^* \otimes \mathbf{W}_{t+1}) \right) \\ \text{s.t. } \|\mathbf{v}_{t+1}\|^2 = P, \\ \mathbf{W}_{t+1}^H \mathbf{W}_{t+1} = \mathbf{I}_{N_{\mathbf{RF}}}. \end{aligned} \quad (19)$$

Note that, in problem (14), the kernel $\Sigma_{\mathbf{h}}$ is decoupled into $\Sigma_{\mathbf{T}} \otimes \Sigma_{\mathbf{R}}$ such that \mathbf{v}_1 and \mathbf{W}_1 can be obtained by selecting the appropriate eigenvectors of $\Sigma_{\mathbf{T}}$ and $\Sigma_{\mathbf{R}}$ as in (16). We attempt the similar idea to solve for \mathbf{v}_{t+1} and \mathbf{W}_{t+1} . To this end, we first define the EVD: $\Sigma_t = \mathbf{U}_t \Lambda_t \mathbf{U}_t^H$ where $\mathbf{U}_t \in \mathbb{C}^{N_{\mathbf{T}} N_{\mathbf{R}} \times N_{\mathbf{T}} N_{\mathbf{R}}}$. Notice that the constraints $\|\mathbf{v}_{t+1}\|^2 = P$ and $\mathbf{W}_{t+1}^H \mathbf{W}_{t+1} = \mathbf{I}_{N_{\mathbf{RF}}}$ make the overall matrix observation \mathbf{X}_{t+1} orthogonal, i.e., $\mathbf{X}_{t+1}^H \mathbf{X}_{t+1} = P \mathbf{I}_{N_{\mathbf{RF}}}$. If we temporarily omit the Kronecker constraint $\mathbf{X}_{t+1} = \mathbf{v}_{t+1}^* \otimes \mathbf{W}_{t+1}$ and try to solve (19) by considering the orthogonal constraint $\mathbf{X}_{t+1}^H \mathbf{X}_{t+1} = P \mathbf{I}_{N_{\mathbf{RF}}}$ only, it becomes evident that the global optimal solution to (19) is $\mathbf{X}_{t+1} = \sqrt{P} \mathbf{U}_t(:, [1, \dots, N_{\mathbf{RF}}])$. Motivated by this discovery, a natural question arises: *when the Kronecker constraint holds, is it possible to set \mathbf{X}_{t+1} as the principal eigenvectors of Σ_t by properly designing \mathbf{v}_{t+1} and \mathbf{W}_{t+1} ?* Addressing this question is crucial for generating near-optimal observation matrices. We would like to investigate it by analyzing the impacts and feasibility of setting \mathbf{X}_{t+1} as the principal eigenvectors of Σ_t .

i) Influence of setting \mathbf{X}_{t+1} as the principal eigenvectors of Σ_t . Before evaluating the feasibility of setting $\mathbf{X}_{t+1} = \sqrt{P} \mathbf{U}_t(:, [1, \dots, N_{\mathbf{RF}}])$, we first need to exploit its influence on the evolution rule of the posterior kernel Σ_{t+1} . The following lemma characterizes the relationship between Σ_{t+1} and Σ_t .

Lemma 3: Let $\lambda_n(\cdot)$ denote the n -th largest eigenvalue of the matrix in its argument, e.g., $\lambda_n(\Sigma_t) = \Lambda_t(n, n)$ for

$\Sigma_t = \mathbf{U}_t \Lambda_t \mathbf{U}_t^H$. If $\mathbf{X}_{t+1} = \sqrt{P} \mathbf{U}_t(:, [1, \dots, N_{\text{RF}}])$, then the EVD of Σ_{t+1} can be derived from $\Sigma_t = \mathbf{U}_t \Lambda_t \mathbf{U}_t^H$ by

$$\Sigma_{t+1} = \mathbf{U}_t \text{diag} \left(\frac{\lambda_1(\Sigma_t) \sigma^2}{P \lambda_1(\Sigma_t) + \sigma^2}, \frac{\lambda_2(\Sigma_t) \sigma^2}{P \lambda_2(\Sigma_t) + \sigma^2}, \dots, \frac{\lambda_{N_{\text{RF}}}(\Sigma_t) \sigma^2}{P \lambda_{N_{\text{RF}}}(\Sigma_t) + \sigma^2}, \lambda_{N_{\text{RF}}+1}(\Sigma_t), \lambda_{N_{\text{RF}}+2}(\Sigma_t), \dots, \lambda_{N_{\text{R}} N_{\text{T}}}(\Sigma_t) \right) \mathbf{U}_t^H. \quad (20)$$

Proof: (See Appendix D). ■

Lemma 3 reveals that, when $\mathbf{X}_{t+1} = \sqrt{P} \mathbf{U}_t(:, [1, \dots, N_{\text{RF}}])$, the posterior covariance matrix Σ_{t+1} shares the same eigenvector space, \mathbf{U}_t , as Σ_t . The only difference on their EVDs is that the N_{RF} -largest eigenvalues of Σ_t , i.e., $\{\lambda_n(\Sigma_t)\}_{n=1}^{N_{\text{RF}}}$, are replaced by $\left\{ \frac{\lambda_n(\Sigma_t) \sigma^2}{P \lambda_n(\Sigma_t) + \sigma^2} \right\}_{n=1}^{N_{\text{RF}}}$ in Σ_{t+1} . Considering the generality of t , we can conclude that, the eigenvectors of channel covariance $\Sigma_0 := \Sigma_{\text{h}}$ are preserved by all subsequent posterior kernels $\Sigma_1, \Sigma_2, \dots$, and Σ_{Q-1} . In other words, the only difference among $\mathbf{U}_0, \mathbf{U}_1, \dots, \mathbf{U}_{Q-1}$ is their column arrangement orders.

ii) *Feasibility of setting \mathbf{X}_{t+1} as the principal eigenvectors of Σ_t .* Given the eigenvector-preserving property in Lemma 3, the feasibility of setting \mathbf{X}_{t+1} as the principal eigenvectors of Σ_t lies in the feasibility of setting \mathbf{X}_{t+1} as the eigenvectors of $\Sigma_0 = \Sigma_{\text{h}}$. To assess this feasibility, we examine the structure of the eigenspace of Σ_{h} below.

Corollary 1: The EVD of the prior covariance Σ_{h} can be written as

$$\begin{aligned} \Sigma_{\text{h}} &= \mathbf{U}_0 \Lambda_0 \mathbf{U}_0^H \\ &= \underbrace{(\mathbf{U}_{\text{T}} \otimes \mathbf{U}_{\text{R}})}_{\text{Orthogonal matrix}} \underbrace{(\Lambda_{\text{T}} \otimes \Lambda_{\text{R}})}_{\text{Eigenvalue matrix}} (\mathbf{U}_{\text{T}}^H \otimes \mathbf{U}_{\text{R}}^H) \\ &= \sum_{n_{\text{T}}=1}^{N_{\text{T}}} \sum_{n_{\text{R}}=1}^{N_{\text{R}}} \alpha_{n_{\text{T}}} \beta_{n_{\text{R}}} (\mathbf{a}_{n_{\text{T}}} \otimes \mathbf{b}_{n_{\text{R}}}) (\mathbf{a}_{n_{\text{T}}}^H \otimes \mathbf{b}_{n_{\text{R}}}^H), \quad (21) \end{aligned}$$

where $\mathbf{a}_{n_{\text{T}}} := \mathbf{U}_{\text{T}}(:, n_{\text{T}})$ denotes the n_{T} -th eigenvector of Σ_{T} and $\mathbf{b}_{n_{\text{R}}} := \mathbf{U}_{\text{R}}(:, n_{\text{R}})$ denotes the n_{R} -th eigenvector of Σ_{R} . In particular, $\{\alpha_{n_{\text{T}}} \beta_{n_{\text{R}}}\}_{n_{\text{T}}=1, n_{\text{R}}=1}^{N_{\text{T}}, N_{\text{R}}}$ and $\{\mathbf{a}_{n_{\text{T}}} \otimes \mathbf{b}_{n_{\text{R}}}\}_{n_{\text{T}}=1, n_{\text{R}}=1}^{N_{\text{T}}, N_{\text{R}}}$ are the eigenvalues and the corresponding eigenvectors of Σ_{h} , respectively.

Proof: (See Appendix E). ■

Recalling that $\mathbf{X}_{t+1} := \mathbf{v}_{t+1}^* \otimes \mathbf{W}_{t+1}$, we find that the analytical form of \mathbf{X}_{t+1} perfectly matches that of the eigenvectors of Σ_{h} , i.e., $\{\mathbf{a}_{n_{\text{T}}} \otimes \mathbf{b}_{n_{\text{R}}}\}_{n_{\text{T}}=1, n_{\text{R}}=1}^{N_{\text{T}}, N_{\text{R}}}$. This encouraging fact inspires us that, the desired optimal \mathbf{X}_{t+1} to solve (19) can be achieved by setting \mathbf{v}_{t+1} and \mathbf{W}_{t+1} to the appropriate eigenvectors from $\{\sqrt{P} \mathbf{a}_{n_{\text{T}}}^*\}_{n_{\text{T}}=1}^{N_{\text{T}}}$ and $\{\mathbf{b}_{n_{\text{R}}}\}_{n_{\text{R}}=1}^{N_{\text{R}}}$, respectively.

We provide an example to show the implementation process. Firstly, to achieve $\mathbf{X}_1 = \sqrt{P} \mathbf{U}_0(:, [1, \dots, N_{\text{RF}}])$ at timeslot $t = 1$, we can set $\mathbf{v}_1 = \sqrt{P} \mathbf{a}_1^* = \sqrt{P} \mathbf{U}_{\text{T}}^*(:, 1)$ and $\mathbf{W}_1 = [\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_{N_{\text{RF}}}] = \mathbf{U}_{\text{R}}(:, [1, \dots, N_{\text{RF}}])$, which coincides with the optimal solution in (16). Note that, according to Lemma 3, this process does not change the eigenvectors of

Σ_1 .² In this context, the desired $\mathbf{X}_2 = \sqrt{P} \mathbf{U}_1(:, [1, \dots, N_{\text{RF}}])$ at timeslot $t = 2$ can also be achieved by carefully selecting \mathbf{v}_2 and \mathbf{W}_2 from $\{\sqrt{P} \mathbf{a}_{n_{\text{T}}}^*\}_{n_{\text{T}}=1}^{N_{\text{T}}}$ and $\{\mathbf{b}_{n_{\text{R}}}\}_{n_{\text{R}}=1}^{N_{\text{R}}}$ respectively, which does not influence the eigenvectors of Σ_2 . Analogously, all our desired $\{\mathbf{X}_q\}_{q=1}^Q$ can be obtained by this successive process. As a result, the problem is transformed into: *How to select appropriate eigenvectors such that the objective function in problem (19) is maximized?*

B. Eigenvector Selection

In this subsection, the problem of eigenvector selection is investigated. According to Lemma 3 and Corollary 1, all posterior kernels $\{\Sigma_t\}_{t=0}^{Q-1}$ can be rewritten as the form of

$$\begin{aligned} \Sigma_t &= \mathbf{U}_0 \text{diag}(\lambda_{1,1}^t, \lambda_{1,2}^t, \dots, \lambda_{N_{\text{T}}, N_{\text{R}}}^t) \mathbf{U}_0^H \\ &= \sum_{n_{\text{T}}=1}^{N_{\text{T}}} \sum_{n_{\text{R}}=1}^{N_{\text{R}}} \lambda_{n_{\text{T}}, n_{\text{R}}}^t (\mathbf{a}_{n_{\text{T}}} \otimes \mathbf{b}_{n_{\text{R}}}) (\mathbf{a}_{n_{\text{T}}}^H \otimes \mathbf{b}_{n_{\text{R}}}^H), \quad (22) \end{aligned}$$

where $\lambda_{n_{\text{T}}, n_{\text{R}}}^t$ is the $(n_{\text{T}}, n_{\text{R}})$ -th eigenvalue of Σ_t associated with the eigenvector $\mathbf{a}_{n_{\text{T}}} \otimes \mathbf{b}_{n_{\text{R}}}$. In particular, $\lambda_{n_{\text{T}}, n_{\text{R}}}^0 = \alpha_{n_{\text{T}}} \beta_{n_{\text{R}}}$, and its update from t to $t+1$ is expressed by

$$\lambda_{n_{\text{T}}, n_{\text{R}}}^{t+1} = \begin{cases} \frac{\lambda_{n_{\text{T}}, n_{\text{R}}}^t \sigma^2}{P \lambda_{n_{\text{T}}, n_{\text{R}}}^t + \sigma^2}, & \mathbf{a}_{n_{\text{T}}}^* \otimes \mathbf{b}_{n_{\text{R}}} \text{ is selected,} \\ \lambda_{n_{\text{T}}, n_{\text{R}}}^t, & \text{else.} \end{cases} \quad (23)$$

Based on the above derivation, we prove the following lemma to further simplify the original problem (19):

Lemma 4: When $\mathbf{v}_{t+1} \in \{\sqrt{P} \mathbf{a}_{n_{\text{T}}}^*\}_{n_{\text{T}}=1}^{N_{\text{T}}}$ and the columns

of \mathbf{W}_{t+1} are belonging to $\{\mathbf{b}_{n_{\text{R}}}\}_{n_{\text{R}}=1}^{N_{\text{R}}}$, the original problem (19) can be transformed into an eigenvector selection problem, written as

$$\begin{aligned} \max_{n_{\text{T}}, \{n_{\text{R},k}\}_{k=1}^{N_{\text{RF}}}} & \sum_{k=1}^{N_{\text{RF}}} \log_2 \left(1 + \frac{P \lambda_{n_{\text{T}}, n_{\text{R},k}}^t}{\sigma^2} \right) \\ \text{s.t.} & n_{\text{T}} \in \{1, \dots, N_{\text{T}}\}, \\ & n_{\text{R},k} \in \{1, \dots, N_{\text{R}}\}, \forall k, \\ & n_{\text{R},k} \neq n_{\text{R},k'}, \forall k \neq k'. \quad (24) \end{aligned}$$

Proof: (See Appendix F). ■

Problem (24) aims to find one eigenvector index n_{T} on the transmitter side and N_{RF} different eigenvector indices $\{n_{\text{R},k}\}_{k=1}^{N_{\text{RF}}}$ on the receiver side, such that the selected eigenvalues $\{\lambda_{n_{\text{T}}, n_{\text{R},k}}^t\}_{k=1}^{N_{\text{RF}}}$ can maximize the objective (24). A linear search algorithm is proposed to solve (24) optimally, as summarized in Algorithm 2. The key idea is to fix an n_{T} and then find N_{RF} -largest values from $\{\lambda_{n_{\text{T}}, n_{\text{R}}}\}_{n_{\text{R}}=1}^{N_{\text{R}}}$ to calculate the objective $\sum_{k=1}^{N_{\text{RF}}} \log_2(1 + P \lambda_{n_{\text{T}}, n_{\text{R},k}}^t / \sigma^2)$. After traversing all $n_{\text{T}} \in \{1, \dots, N_{\text{T}}\}$, the optimal $n_{\text{T}}^{\text{opt}}$ and $\{n_{\text{R},k}^{\text{opt}}\}_{k=1}^{N_{\text{RF}}}$ can be obtained from the indices of the maximum objective. Finally, the desired precoder and combiner are thereby expressed as

$$\mathbf{v}_{t+1}^{\text{opt}} = \sqrt{P} \mathbf{a}_{n_{\text{T}}^{\text{opt}}}^* \text{ and } \mathbf{W}_{t+1}^{\text{opt}} = [\mathbf{b}_{n_{\text{R},1}^{\text{opt}}}, \dots, \mathbf{b}_{n_{\text{R},N_{\text{RF}}}^{\text{opt}}}], \quad (25)$$

²Thus, \mathbf{U}_1 and $\mathbf{U}_0 := \mathbf{U}_{\text{T}} \otimes \mathbf{U}_{\text{R}}$ share the same columns, but their column arrangement orders may be different due to the eigenvalue updates.

Algorithm 2 Linear Search for Eigenvalue Selection

Input: Eigenvalues $\{\lambda_{n_T, n_R}^t\}_{n_T=1, n_R=1}^{N_T, N_R}$ in timeslot t .
Output: Optimal eigenvalue indices n_T^{opt} and $\{n_{R,k}^{\text{opt}}\}_{k=1}^{N_{RF}}$ that maximize $\sum_{k=1}^{N_{RF}} \log_2(1 + \lambda_{n_T, n_{R,k}}^t / \sigma^2)$.

- 1: Initialize indexes: $n_T^{\text{opt}} = 1$ and $[n_{R,1}^{\text{opt}}, \dots, n_{R,N_{RF}}^{\text{opt}}] = [1, \dots, N_{RF}]$.
- 2: Initialize the maximum objective: $\zeta_{\max} = \sum_{k=1}^{N_{RF}} \log_2(1 + P\lambda_{n_T^{\text{opt}}, n_{R,k}^{\text{opt}}}^t / \sigma^2)$
- 3: **for** $n_T = 1, \dots, N_T$ **do**
- 4: Find the N_{RF} -largest values from $\{\lambda_{n_T, n_R}^t\}_{n_R=1}^{N_R}$, and then denote their second indexes as $\{n_{R,k}\}_{k=1}^{N_{RF}}$.
- 5: **if** $\sum_{k=1}^{N_{RF}} \log_2(1 + P\lambda_{n_T, n_{R,k}}^t / \sigma^2) > \zeta_{\max}$ **then**
- 6: Update the optimal indexes by $n_T^{\text{opt}} = n_T$ and $[n_{R,1}^{\text{opt}}, \dots, n_{R,N_{RF}}^{\text{opt}}] = [n_{R,1}, \dots, n_{R,N_{RF}}]$
- 7: Update the maximum objective: $\zeta_{\max} = \sum_{k=1}^{N_{RF}} \log_2(1 + P\lambda_{n_T^{\text{opt}}, n_{R,k}^{\text{opt}}}^t / \sigma^2)$
- 8: **end if**
- 9: **end for**
- 10: **return** Optimal n_T^{opt} and $\{n_{R,k}^{\text{opt}}\}_{k=1}^{N_{RF}}$.

respectively, which generate a feasible observation matrix $\mathbf{X}_{t+1}^{\text{opt}} = \mathbf{v}_{t+1}^{\text{opt}*} \otimes \mathbf{W}_{t+1}^{\text{opt}}$ at timeslot $t+1$. One can verify without difficulty that the initialized precoder and combiner obtained in (16) are a special case of (25) when $t=0$.

To summarize, the eigenvalue updating rule in (23), as well as the eigenvector selection method stated in **Algorithm 2** and (25), allow us to calculate all observation matrices.

C. Insightful Interpretation to 2DIF

In this subsection, we provide insightful explanations to the proposed 2DIF algorithm to clarify its physical significance. At first, the relationship between the well-known water-filling method and the proposed 2DIF method is discussed. Then, the superiority of the proposed 2DIF method over the existing IF method [36] is illustrated.

1) *Ideal Water-Filling:* To better understand the proposed observation matrix design, we first interpret problem (9) from the view of *water-filling*. Specifically, the orthogonal property of \mathbf{W}_q proved in **Lemma 2** allows us to replace the noise covariance matrix Ξ in (9) by $\sigma^2 \mathbf{I}_{N_{RF}Q}$ without affecting the optimal $I(\mathbf{y}; \mathbf{h})$. Then, by further relaxing the constraints $\mathbf{W}_q^T \mathbf{W}_q = \mathbf{I}_{N_{RF}}$ and $\|\mathbf{v}_q\|^2 = P$, we focus only on the total power constraint imposed on the overall observation matrix \mathbf{X} , i.e., $\text{Tr}(\mathbf{X}\mathbf{X}^H) = PN_{RF}Q$. In this case, the optimal value of problem (9) is shown to have an upper bound:

$$\begin{aligned} \max_{\mathbf{X}} \quad & \log_2 \det \left(\mathbf{I}_{N_{RF}Q} + \frac{1}{\sigma^2} \mathbf{X}^H \Sigma_{\mathbf{h}} \mathbf{X} \right) \\ \text{s.t.} \quad & \text{Tr}(\mathbf{X}\mathbf{X}^H) = PN_{RF}Q. \end{aligned} \quad (26)$$

Notably, this upper bound is equivalent to the channel capacity of a point-to-point MIMO system equipped with $N_T N_R$

transmit antennas and $N_{RF}Q$ receive antennas. Thereafter, the overall observation matrix can be optimally solved as³

$$\mathbf{X}^{\text{ideal}} = \mathbf{U}_0(:, [1, \dots, N_{RF}Q]) \mathbf{P}, \quad (27)$$

where $\mathbf{P} = \text{diag}(\sqrt{p_1}, \dots, \sqrt{p_{N_{RF}Q}})$ is the power allocation matrix. The power allocated to the n -th eigenvector is determined by the water-filling principle, i.e., $p_n = \left(\beta - \frac{\sigma^2}{\lambda_n(\Sigma_{\mathbf{h}})}\right)^+$, where the water-level β is adjusted to satisfy the total power constraint $\text{Tr}(\mathbf{X}^{\text{ideal}}(\mathbf{X}^{\text{ideal}})^H) = \sum_{n=1}^{N_{RF}Q} p_n = PN_{RF}Q$.

Although the ideal observation matrix $\mathbf{X}^{\text{ideal}}$, that maximizes the upper bound (26), might not be implementable in practice (as $\mathbf{X}^{\text{ideal}}$ may violate the constraints $\mathbf{W}_q^T \mathbf{W}_q = \mathbf{I}_{N_{RF}}$ and $\|\mathbf{v}_q\|^2 = P$), it can give us two pivotal intuitions. First, the observation matrix should align with the eigenspace $\{\mathbf{a}_{n_T} \otimes \mathbf{b}_{n_R}\}_{n_T=1, n_R=1}^{N_T, N_R}$ of the full MIMO channel covariance, $\Sigma_{\mathbf{h}} = \Sigma_T \otimes \Sigma_R$. Second, we need to fill in more power (water) to the eigenvectors having larger eigenvalues $\{\lambda_n(\Sigma_{\mathbf{h}})\}_{n=1}^{N_{RF}Q}$ (or equivalently lower base levels $\left\{\frac{\sigma^2}{\lambda_n(\Sigma_{\mathbf{h}})}\right\}_{n=1}^{N_{RF}Q}$).

2) *2DIF Versus Water-Filling:* The proposed 2DIF algorithm materializes the above two intuitions under the practical constraints $\mathbf{W}_q^T \mathbf{W}_q = \mathbf{I}_{N_{RF}}$ and $\|\mathbf{v}_q\|^2 = P$ via the eigenvector selection process in (24). The first intuition is automatically achieved by assigning \mathbf{v}_{t+1} with a eigenvector from $\{\sqrt{P}\mathbf{a}_{n_T}^*\}_{n_T=1}^{N_T}$ and assigning the columns of \mathbf{W}_{t+1} with different eigenvectors from $\{\mathbf{b}_{n_R}\}_{n_R=1}^{N_R}$, as proved in **Lemma 4**. The second intuition is approximately accomplished via selecting eigenvectors that have lower base levels, $\left\{\frac{\sigma^2}{\lambda_n(\Sigma_{\mathbf{h}})}\right\}$, by more times. This is attributed to the fact that the maximization of (24) tends to select an eigenvalue combination $\{\lambda_{n_T, n_R, k}^t\}_{k=1}^{N_{RF}}$ that has the lowest $\left\{\frac{\sigma^2}{\lambda_{n_T, n_R, k}^t}\right\}_{k=1}^{N_{RF}}$ on average. To see this approximation more clearly, we rewrite the updating rule of the selected eigenvalue in (23) as

$$\underbrace{\frac{\sigma^2}{\lambda_{n_T, n_R}^{t+1}}}_{\text{Updated ice-level}} = \underbrace{\frac{\sigma^2}{\lambda_{n_T, n_R}^t}}_{\text{Current ice-level}} + \underbrace{P}_{\text{Height of an ice block}}. \quad (28)$$

Equation (28) reveals that every time the eigenvector $\mathbf{a}_{n_T} \otimes \mathbf{b}_{n_R}$ is selected, the value of $\frac{\sigma^2}{\lambda_{n_T, n_R}^t}$ increases by P . Similar to the water-filling process, the process described by (28) can be vividly interpreted as allocating an ice block having P -unit power to the (n_T, n_R) -th orthogonal channel, where $\frac{\sigma^2}{\lambda_{n_T, n_R}^t}$ is viewed as the ice level in the t -th timeslot. To summarize, due to the consideration of N_{RF} RF chains and $N_T \times N_R$ MIMO systems, in each timeslot, the 2DIF first selects N_{RF} orthogonal channels, which have the deepest ice-levels $\left\{\frac{\sigma^2}{\lambda_{n_T, n_R, k}^t}\right\}_{k=1}^{N_{RF}}$ on average, from the total $N_T \times N_R$ channels. Then, the 2DIF will fill N_{RF} ice blocks (i.e., N_{RF} pilots) of height P onto them. As illustrated in Fig. 2(a), after Q timeslots, the final ice-levels of all channels can have a similar height with the water-level, β , determined by the water-filling principle. In this case, the second intuition is approximately achieved.

³For ease of discussion, we assume that $N_{RF}Q$ is smaller than the rank of channel covariance, $\Sigma_{\mathbf{h}}$.

IV. PROPOSED TWO-STAGE 2DIF (TS-2DIF) BASED OBSERVATION MATRIX DESIGN

Turn now to the receiver architecture with phase-only controllable analog combiner presented in Fig. 1(b). As the coefficients of matrices $\{\mathbf{W}_q = \mathbf{A}_q \mathbf{D}_q\}_{q=1}^Q$ can no longer be freely manipulated, the proposed 2DIF algorithm might not be implementable in these scenarios. To address this problem, a TS-2DIF algorithm is proposed in this section.

A. Overview of TS-2DIF Observation Matrix Design

Recall that the phase-only controllable structure in Fig. 1(b) requires to express the hybrid combiner as $\mathbf{W}_q = \mathbf{A}_q \mathbf{D}_q$. Each element of the analog combiner \mathbf{A}_q is restricted by the modulus constraint $|\mathbf{A}_q(n_R, n_{RF})| = \frac{1}{\sqrt{N_R}}$, for $n_R \in \{1, \dots, N_R\}$ and $n_{RF} \in \{1, \dots, N_{RF}\}$. This poses a structural constraint on the feasible set of \mathbf{W}_q , destroying its eigenvector structure, thereby making the 2DIF algorithm inapplicable. Therefore, it becomes necessary to redesign a new set of observation matrices $\{\mathbf{v}_q\}_{q=1}^Q$ and $\{\mathbf{W}_q = \mathbf{A}_q \mathbf{D}_q\}_{q=1}^Q$ tailored for the phase-only controllable architecture.

For this purpose, we propose a TS-2DIF algorithm as summarized in **Algorithm 3**, which consists of two stages. In the first stage, **Algorithm 1** is carried out to obtain the ideal observation matrix $\mathbf{X}_q^{\text{IF}} = (\mathbf{v}_q^{\text{IF}})^* \otimes \mathbf{W}_q^{\text{IF}}$ for $\forall q$, where the superscript "IF" is used to indicate the observation matrices generated by the 2DIF algorithm. Subsequently, the second stage aims to make the newly designed observation matrix $\mathbf{X}_q = \mathbf{v}_q^* \otimes (\mathbf{A}_q \mathbf{D}_q)$ sufficiently close to the ideal observation matrix \mathbf{X}_q^{IF} . To this end, the joint optimization of $\{\mathbf{v}_q\}_{q=1}^Q$ and $\{\mathbf{A}_q \mathbf{D}_q\}_{q=1}^Q$ are formulated as:

$$\min_{\mathbf{v}_q, \mathbf{A}_q, \mathbf{D}_q} \|\mathbf{X}_q^{\text{IF}} - \mathbf{v}_q^* \otimes (\mathbf{A}_q \mathbf{D}_q)\|_F^2, \quad (29)$$

$$\text{s.t. } \|\mathbf{v}_q\|^2 = P, \quad (29a)$$

$$|\mathbf{A}_q| = \frac{1}{\sqrt{N_R}} \mathbf{1}_{N_R \times N_{RF}}, \quad (29b)$$

where $\mathbf{1}_{N_R \times N_{RF}}$ is an N_R -by- N_{RF} all-one matrix. By solving problem (29), the newly designed observation matrices $\{\mathbf{X}_q\}_{q=1}^Q$ are expected to achieve a comparable channel estimation performance with $\{\mathbf{X}_q^{\text{IF}}\}_{q=1}^Q$.

B. Joint Optimization of Precoders and Hybrid Combiners

It is intricate to directly solve problem (29) owing to the non-convex modulus constraint in (29b) as well as the coupled relationship of \mathbf{v}_q , \mathbf{A}_q , and \mathbf{D}_q in the objective (29). To overcome these challenges, we exploit the alternating minimization method to iteratively update \mathbf{A}_q and \mathbf{D}_q , and \mathbf{v}_q until an convergence condition triggers. The detailed optimization procedures are elaborated one by one as follows.

1) *Fix \mathbf{A}_q and \mathbf{v}_q , and Optimize \mathbf{D}_q* : For ease of discussion, we denote $\mathbf{v}_q = [v_q(1), v_q(2), \dots, v_q(N_T)]^T$ and define $\mathbf{X}_{q,n}^{\text{IF}} \in \mathbb{C}^{N_R \times N_{RF}}$ as the n -th block component of \mathbf{X}_q^{IF} such that $\mathbf{X}_q^{\text{IF}} = [(\mathbf{X}_{q,1}^{\text{IF}})^T, (\mathbf{X}_{q,2}^{\text{IF}})^T, \dots, (\mathbf{X}_{q,N_T}^{\text{IF}})^T]^T$. Then, when

Algorithm 3 TS-2DIF Based Combiner and Precoder Design

Input: Number of pilots Q , kernel $\Sigma_{\mathbf{h}}$.

Output: Designed precoders $\{\mathbf{v}_q^{\text{opt}}\}_{q=1}^Q$ and hybrid combiners $\{\mathbf{A}_q^{\text{opt}}\}_{q=1}^Q$ and $\{\mathbf{D}_q^{\text{opt}}\}_{q=1}^Q$.

Stage 1 (Optimal observation matrix design)

1: Obtain the ideal precoders $\{\mathbf{v}_q^{\text{IF}}\}_{q=1}^Q$ and overall combiners $\{\mathbf{W}_q^{\text{IF}}\}_{q=1}^Q$ from **Algorithm 1**

2: Get the ideal observation matrix $\mathbf{X}_q^{\text{IF}} = (\mathbf{v}_q^{\text{IF}})^* \otimes \mathbf{W}_q^{\text{IF}}$ for all $q \in \{1, \dots, Q\}$

Stage 2 (Joint hybrid combiner and precoder design)

3: **for** $q = 1, \dots, Q$ **do**

4: **while** no convergence of $\|\mathbf{X}_q^{\text{IF}} - \mathbf{v}_q^* \otimes (\mathbf{A}_q \mathbf{D}_q)\|_F^2$ **do**

5: Update the digital combiner \mathbf{D}_q by (31)

6: Update the analog combiner \mathbf{A}_q by (36)

7: Update the precoder \mathbf{v}_q by (38)

8: **end while**

9: **end for**

10: **return** Designed precoders $\{\mathbf{v}_q\}_{q=1}^Q$ and hybrid combiners $\{\mathbf{A}_q\}_{q=1}^Q$ and $\{\mathbf{D}_q\}_{q=1}^Q$ for channel estimation.

keeping the combiner \mathbf{A}_q and the precoder \mathbf{v}_q fixed, the sub-problem for optimizing \mathbf{D}_q is expressed as

$$\min_{\mathbf{D}_q} \sum_{n=1}^{N_T} \|\mathbf{X}_{q,n}^{\text{IF}} - v_q^*(n) \mathbf{A}_q \mathbf{D}_q\|_F^2. \quad (30)$$

Problem (30) is a standard quadratic programming (QP), which can be optimally solved by making the gradient of objective function to zero, i.e.,

$$\begin{aligned} \mathbf{D}_q^{\text{opt}} &= \left(\sum_{n=1}^{N_T} |v_q(n)|^2 \mathbf{A}_q^H \mathbf{A}_q \right)^{-1} \left(\sum_{n=1}^{N_T} v_q(n) \mathbf{A}_q^H \mathbf{X}_{q,n}^{\text{IF}} \right) \\ &\stackrel{(a)}{=} \sum_{n=1}^{N_T} \frac{v_q(n)}{P} (\mathbf{A}_q^H \mathbf{A}_q)^{-1} \mathbf{A}_q^H \mathbf{X}_{q,n}^{\text{IF}}, \end{aligned} \quad (31)$$

where (a) holds because $\|\mathbf{v}_q\|^2 = \sum_{n=1}^{N_T} |v_q(n)|^2 = P$.

2) *Fix \mathbf{D}_q and \mathbf{v}_q , and Optimize \mathbf{A}_q* : We then fix the digital combiner \mathbf{D}_q and precoder \mathbf{v}_q , and seeks an analog combiner that optimizes the following sub-problem:

$$\min_{\mathbf{A}_q} \sum_{n=1}^{N_T} \|\mathbf{X}_{q,n}^{\text{IF}} - \mathbf{A}_q \mathbf{D}_q v_q^*(n)\|_F^2, \quad (32)$$

$$\text{s.t. } |\mathbf{A}_q| = \frac{1}{\sqrt{N_R}} \mathbf{1}_{N_R \times N_{RF}}. \quad (32a)$$

Directly optimizing problem (32) is challenging owing to the constant modulus constraint (32a) and the product of \mathbf{A}_q and \mathbf{D}_q . To address this issue, we notice that the objective function has the an upper bound due to the Cauchy-Schwarz inequality,

$$\begin{aligned} &\sum_{n=1}^{N_T} \|\mathbf{X}_{q,n}^{\text{IF}} - \mathbf{A}_q \mathbf{D}_q v_q^*(n)\|_F^2 \\ &\leq \sum_{n=1}^{N_T} \|\mathbf{X}_{q,n}^{\text{IF}} \mathbf{D}_q^{-1} - \mathbf{A}_q v_q^*(n)\|_F^2 \|\mathbf{D}_q\|_F^2. \end{aligned} \quad (33)$$

In (33), the analog combiner \mathbf{A}_q has escaped from the product form with \mathbf{D}_q , which can significantly simplify the optimization problem. Taking this into account, we replace the original objective function with $\sum_{n=1}^{N_T} \|\mathbf{X}_{q,n}^{\text{IF}} \mathbf{D}_q^{-1} - \mathbf{A}_q v_q^*(n)\|_F^2$. Then, by defining $\mathbf{J}_{q,n} = \mathbf{X}_{q,n}^{\text{IF}} \mathbf{D}_q^{-1}$, the new objective function can be further simplified as

$$\begin{aligned} & \sum_{n=1}^{N_T} \|\mathbf{J}_{q,n} - \mathbf{A}_q v_q^*(n)\|_F^2 \\ &= C_1 + \sum_{n=1}^{N_T} (\|\mathbf{A}_q\|_F^2 \|v_q(n)\|^2 - 2\text{Tr}\{\Re\{v_q^*(n) \mathbf{J}_{q,n}^H \mathbf{A}_q\}\}) \\ &= C_1 + N_{\text{RF}} P - 2\text{Tr}\{\Re\{\mathbf{J}^H \mathbf{A}_q\}\}, \end{aligned} \quad (34)$$

where $C_1 = \sum_{n=1}^{N_T} \|\mathbf{J}_{q,n}\|_F^2$ and $\mathbf{J} := \sum_{n=1}^{N_T} v_q(n) \mathbf{J}_{q,n}$. By combining (33) and (34), the new optimization problem is formulated as

$$|\mathbf{A}_q| = \max_{\substack{\mathbf{A}_q \\ \|\mathbf{A}_q\|_F = \sqrt{N_{\text{RF}} P}}} \text{Tr}\{\Re\{\mathbf{J}^H \mathbf{A}_q\}\}, \quad (35)$$

Evidently, the optimal solution of (35) is given by

$$\mathbf{A}_q^{\text{opt}} = \frac{1}{\sqrt{N_{\text{RF}}}} \exp(j\angle \mathbf{J}), \quad (36)$$

which completes the update of \mathbf{A}_q in step 6 of **Algorithm 3**.

3) *Fix \mathbf{A}_q and \mathbf{D}_q , and Optimize \mathbf{v}_q* : With given combiner matrices \mathbf{A}_q and \mathbf{D}_q , the objective function in (29) for optimizing \mathbf{v}_q can be rewritten as

$$\begin{aligned} & \sum_{n=1}^{N_T} \|\mathbf{X}_{q,n}^{\text{IF}} - \mathbf{A}_q \mathbf{D}_q v_q^*(n)\|_F^2 \\ &= \|\mathbf{X}_q^{\text{IF}}\|_F^2 + \|\mathbf{A}_q \mathbf{D}_q\|_F^2 - 2\Re\{\mathbf{c}_q^H \mathbf{v}_q\}, \end{aligned} \quad (37)$$

where $\mathbf{c}_q = [\text{Tr}\{(\mathbf{X}_{q,1}^{\text{IF}})^H \mathbf{A}_q \mathbf{D}_q\}, \dots, \text{Tr}\{(\mathbf{X}_{q,N_T}^{\text{IF}})^H \mathbf{A}_q \mathbf{D}_q\}]^T$. By further considering the power constraint in (29a), the optimal $\mathbf{v}_q^{\text{opt}}$ is given as⁴

$$\mathbf{v}_q^{\text{opt}} = \sqrt{P} \mathbf{c}_q / \|\mathbf{c}_q\|. \quad (38)$$

This completes the update of \mathbf{v}_q in step 7 of **Algorithm 3**.

To summarize, *Stage 2* of **Algorithm 3** alternatively updates \mathbf{D}_q , \mathbf{A}_q , and \mathbf{v}_q using (31), (36), and (38) until convergence. After obtaining the hybrid combiners and precoders for pilot transmission, the channel estimator (7) can be utilized to recover wireless channel matrices.

V. COMPUTATIONAL COMPLEXITY ANALYSIS AND KERNEL SELECTION

In this section, the computational complexities of the proposed algorithms are analyzed. Then, the kernel selection for the proposed channel estimator is discussed.

⁴This work assumes the digital precoder \mathbf{v}_q because it is technically feasible and commercially relevant for user equipment. In the analog precoding case, (38) should be modified as $\mathbf{v}_q^{\text{opt}} = \sqrt{P/N_T} \exp(j\angle \mathbf{c}_q)$. In the hybrid precoding case, as proved in [49], vector $\mathbf{v}_q^{\text{opt}}$ in (38) can be perfectly realized by the hybrid precoder with at least 2 RF chains.

A. Computational Complexity Analysis

Consider **Algorithm 1** at first. The EVDs for Σ_T and Σ_R require $\mathcal{O}(N_T^3 + N_R^3)$ FLOPS. The update of eigenvalues requires $\mathcal{O}(QN_{\text{RF}})$ FLOPS in total. In **Algorithm 2**, the linear search requires the sort operations with the complexity of $\mathcal{O}(N_T N_R \log_2(N_R))$, and calculating the objective requires $\mathcal{O}(N_T N_{\text{RF}})$ FLOPS. Thus, the overall complexity of **Algorithm 1** is $\mathcal{O}(N_T^3 + N_R^3 + N_T N_R \log_2(N_R) + (Q + N_T) N_{\text{RF}})$. The computational complexity of **Algorithm 3** is dominated by the alternating optimizations of \mathbf{D}_q , \mathbf{A}_q , and \mathbf{v}_q . In particular, their computations require $\mathcal{O}(Q(N_{\text{RF}}^2 N_R + N_{\text{RF}}^3))$, $\mathcal{O}(Q(N_{\text{RF}}^2 N_R + N_{\text{RF}}^3 + N_T N_R N_{\text{RF}}^2))$, and $\mathcal{O}(QN_T N_R N_{\text{RF}}^2)$ FLOPS, respectively. Assuming that the number of iterations is I_o , the overall computational complexity of **Algorithm 3** is $\mathcal{O}(I_o Q(N_{\text{RF}}^3 + N_T N_R N_{\text{RF}}^2))$.

It is worth noting that, **Algorithms 1~3** only rely on the given kernel Σ_h instead of the instantaneous channels or received pilots, thus they do not need to be implemented in real time. Since the channel covariance does not change so frequently, the designed observation matrices can be deployed online for channel estimation for a long time. Thereby, from the long-term perspective, the computational complexity of online deploying the proposed channel estimator is not dominated by the observation matrix design. This advantage allows the proposed 2DIF to achieve lower computational complexity compared with the existing channel estimators [50], [51]. Besides, since the proposed algorithms only aim to design precoders and combiners, they can naturally share the same running process as the current 5G new radio (NR) systems.

B. Kernel Selection

Selecting an appropriate covariance matrix, i.e., kernel Σ , is crucial for constructing a robust estimator. The kernel Σ dictates the shape and adaptability of the estimator, thereby influencing its performances to detect functional trends and provide precise predictions. Given the localized-correlation characteristic of MIMO channels, the ideal kernel should enhance the similarity between adjacent antennas while diminishing its impact as the distance increases. In this section, three kinds of kernels are recommended.

1) *Statistical Kernel*: Given the kernel's role as the prior covariance of channel $\mathbf{h} = \text{vec}(\mathbf{H})$, the optimal strategy is to utilize the actual covariance for channel estimation, i.e., $\Sigma_h = \text{E}(\text{vec}(\mathbf{H})\text{vec}(\mathbf{H})^H)$. Prior to deploying the proposed estimator online, it is feasible to train an approximation of Σ_h in advance by leveraging some existing channel models or channel datasets [44], [45], [46], [47]. Concretely, according to the law of large numbers, Σ_h can be trained by

$$\Sigma_h \approx \frac{1}{R} \sum_{r=1}^R \text{vec}(\mathbf{H}_r) \text{vec}(\mathbf{H}_r)^H, \quad (39)$$

where R is the number of channel realizations and \mathbf{H}_r is the r -th channel realization used for kernel training. As for the covariance matrices Σ_T and Σ_R that characterize the transmitter-side

and receive-side channel covariance, respectively, they can be obtained by $\Sigma_T = \frac{1}{RN_T} \sum_{r=1}^R \sum_{n=1}^{N_R} \mathbf{H}_r^T(n, :) \mathbf{H}_r^*(n, :)$ and $\Sigma_R = \frac{1}{RN_T} \sum_{r=1}^R \sum_{n=1}^{N_T} \mathbf{H}_r(:, n) \mathbf{H}_r^H(:, n)$.

2) *Artificial Kernels*: In practical scenarios where obtaining an explicit channel model or channel dataset is challenging, it is preferred to train an artificial kernel to replace Σ_h [52]. For simplicity, we assume that the uniform linear arrays (ULAs) are deployed at both the BS and the user, while it can be easily extended to the UPA case. The consistency between the artificial kernels and the statistical kernel is that both of them assign higher similarity to nearby antennas and decrease influence with distance. Given array parameters, the mutual coupling matrices \mathbf{C}_{rx} and \mathbf{C}_{tx} can be calculated or measured. Thus here we focus on characterizing the spatial correlations \mathbf{R}_{tx} and \mathbf{R}_{tx} , two artificial kernels are recommended [53]:

i) *Laplace kernel*: The Laplace kernel Σ_{La} is the most popular choice in Bayesian estimation. Let $\mathbf{n}_T = [-\frac{N_T-1}{2}, -\frac{N_T-3}{2}, \dots, \frac{N_T-1}{2}]^T$ and $\mathbf{n}_R = [-\frac{N_R-1}{2}, -\frac{N_R-3}{2}, \dots, \frac{N_R-1}{2}]^T$. Then, the Laplace kernels, which respectively characterize the spatial correlations at the user and the BS, can be modeled as

$$\mathbf{R}_{La,T} = \exp\left(-\eta^2 \frac{d^2}{\lambda^2} |\mathbf{1}_{N_T}^T \otimes \mathbf{n}_T - \mathbf{n}_T^T \otimes \mathbf{1}_{N_T}|^{\odot 2}\right), \quad (40)$$

$$\mathbf{R}_{La,R} = \exp\left(-\eta^2 \frac{d^2}{\lambda^2} |\mathbf{1}_{N_R}^T \otimes \mathbf{n}_R - \mathbf{n}_R^T \otimes \mathbf{1}_{N_R}|^{\odot 2}\right), \quad (41)$$

where $\eta > 0$ is an adjustable hyperparameter; and $\mathbf{Z}^{\odot 2}$ denotes the element-wise product of two matrices \mathbf{Z} . Thus, the overall kernel can be written as $\Sigma_{La} = ((\mathbf{C}_{tx}^{1/2})^T \mathbf{R}_{La,T} (\mathbf{C}_{tx}^{1/2})^*) \otimes ((\mathbf{C}_{rx}^{1/2})^T \Sigma_{La,R} (\mathbf{C}_{rx}^{1/2})^*)$.

ii) *Bessel kernel*: By exploiting the inherent periodic property of Bessel functions, the Bessel kernel, denoted as Σ_{Be} , is well-suited for recovering data with oscillatory patterns. The Bessel kernels, which respectively characterizes the spatial correlations at the user and the BS, can be modeled as

$$\mathbf{R}_{Be,T} = J_0\left(\eta \frac{d}{\lambda} |\mathbf{1}_{N_T}^T \otimes \mathbf{n}_T - \mathbf{n}_T^T \otimes \mathbf{1}_{N_T}|\right), \quad (42)$$

$$\mathbf{R}_{Be,R} = J_0\left(\eta \frac{d}{\lambda} |\mathbf{1}_{N_R}^T \otimes \mathbf{n}_R - \mathbf{n}_R^T \otimes \mathbf{1}_{N_R}|\right), \quad (43)$$

where J_0 is the zero-order Bessel function of the first kind and $\eta > 0$ is a hyperparameter. The overall kernel is written as $\Sigma_{Be} = ((\mathbf{C}_{tx}^{1/2})^T \mathbf{R}_{Be,T} (\mathbf{C}_{tx}^{1/2})^*) \otimes ((\mathbf{C}_{rx}^{1/2})^T \Sigma_{Be,R} (\mathbf{C}_{rx}^{1/2})^*)$.

The hyperparameter η plays a pivotal role in closely mirroring the real channel covariance, whose value can be determined by a maximum likelihood (ML) estimator. We assume that R channel realizations are utilized to train a kernel Σ , where Σ

can be either Σ_{La} or Σ_{Be} . Then, the estimation of η is written as

$$\eta^{\text{opt}} = \arg \max_{\eta > 0} \sum_{r=1}^R \ln(\mathbf{P}(\mathbf{y}_r | \eta)), \quad (44)$$

wherein the likelihood function is given by

$$\mathbf{P}(\mathbf{y}_r | \eta) = \frac{\exp\left(-\mathbf{y}_r^H (\mathbf{X}_r^H \Sigma \mathbf{X}_r + \mathbf{\Xi}_r)^{-1} \mathbf{y}_r\right)}{\pi^{QN_{RF}} \det(\mathbf{X}_r^H \Sigma \mathbf{X}_r + \mathbf{\Xi}_r)}, \quad (45)$$

in which $\mathbf{y}_r = \mathbf{X}_r^H \mathbf{h}_r + \mathbf{z}_r \in \mathbb{C}^{QN_{RF}}$ denotes the received pilot associated with the r -th channel realization \mathbf{h}_r for kernel training; \mathbf{X}_r is obtained based on the randomly generated precoders and combiners; and $\mathbf{\Xi}_r$ is the covariance of AWGN \mathbf{z}_r . One-dimensional search method is adopted to obtain η^{opt} .

3) *Adaptive Kernel*: In cases where the statistical kernel is unavailable and the artificial kernels are inaccurate, we introduce a novel adaptive kernel training strategy. It uses channels estimated by the 2DIF algorithm to adaptively update the channel kernel for implementing 2DIF. The method smoothly integrates the kernel training and channel estimation stages, eliminating the additional time period required for estimating the covariance matrix.

Consider T_f consecutive frames, where the inter-frame channels $\mathbf{H}_{t_f} \in \mathbb{C}^{N_R \times N_T}, \forall t_f \in \{1, 2, \dots, T_f\}$ are i.i.d distributed, following $\text{vec}(\mathbf{H}_{t_f}) \sim \mathcal{CN}(0, \Sigma_T \otimes \Sigma_R)$. Our strategy iteratively estimates the channels $\hat{\mathbf{H}}_{t_f}$ and updates the kernels $\hat{\Sigma}_T^{(t_f)}$ and $\hat{\Sigma}_R^{(t_f)}$, targeting at gradually improving the accuracy of $\hat{\mathbf{H}}_{t_f}$ and converging $\hat{\Sigma}_T^{(t_f)} \otimes \hat{\Sigma}_R^{(t_f)}$ to the real kernel $\Sigma_T \otimes \Sigma_R$. To be specific, at frame t_f , the channel \mathbf{H}_{t_f} is estimated by 2DIF using the kernel $\hat{\Sigma}_T^{(t_f-1)} \otimes \hat{\Sigma}_R^{(t_f-1)}$ trained during frames $1 \sim t_f - 1$:

$$\hat{\Sigma}_T^{(t_f-1)} \otimes \hat{\Sigma}_R^{(t_f-1)} \xrightarrow{2\text{DIF}} \hat{\mathbf{H}}_{t_f}. \quad (47)$$

Then, we accumulate the information of the newly estimated channel $\hat{\mathbf{H}}_{t_f}$ to update the kernels $\hat{\Sigma}_T^{(t_f)}$ and $\hat{\Sigma}_R^{(t_f)}$ as

$$\begin{cases} \hat{\Sigma}_T^{(t_f)} = \frac{t_f-1}{t_f} \hat{\Sigma}_T^{(t_f-1)} + \frac{1}{t_f N_R} \sum_{n=1}^{N_R} \hat{\mathbf{H}}_{t_f}^T(n, :) \hat{\mathbf{H}}_{t_f}^*(n, :) \\ \hat{\Sigma}_R^{(t_f)} = \frac{t_f-1}{t_f} \hat{\Sigma}_R^{(t_f-1)} + \frac{1}{t_f N_T} \sum_{n=1}^{N_T} \hat{\mathbf{H}}_{t_f}(:, n) \hat{\mathbf{H}}_{t_f}^H(:, n). \end{cases} \quad (48)$$

As a result, the proposed adaptive kernel training strategy proceeds as in (46), Shown at the bottom of the page. To trigger the working flow in (46), the initial kernels $\hat{\Sigma}_T^{(0)}$ and $\hat{\Sigma}_R^{(0)}$ are set as identity matrices: $\hat{\Sigma}_T^{(0)} = \mathbf{I}_{N_T}$ and $\hat{\Sigma}_R^{(0)} = \mathbf{I}_{N_R}$.

The proposed adaptive kernel learning method exhibits two key advantages. Firstly, as the initial kernels are identity matrices, the proposed method requires no prior information about the real kernels, making it applicable to general and practical communication systems. Secondly, the proposed method allows

$$\hat{\Sigma}_T^{(0)} \otimes \hat{\Sigma}_R^{(0)} \xrightarrow{2\text{DIF}} \hat{\mathbf{H}}_1 \xrightarrow{(48)} \hat{\Sigma}_T^{(1)} \otimes \hat{\Sigma}_R^{(1)} \xrightarrow{2\text{DIF}} \hat{\mathbf{H}}_2 \xrightarrow{(48)} \hat{\Sigma}_T^{(2)} \otimes \hat{\Sigma}_R^{(2)} \xrightarrow{2\text{DIF}} \dots \xrightarrow{2\text{DIF}} \hat{\mathbf{H}}_{T_f} \xrightarrow{(48)} \hat{\Sigma}_T^{(T_f)} \otimes \hat{\Sigma}_R^{(T_f)}. \quad (46)$$

frame 1
frame 2
frame T_f

the coexistence of channel estimation and kernel training within one frame, as it leverages the channels estimated by 2DIF itself to train the kernel. This integration eliminates the additional time period required for kernel learning, greatly simplifying the frame structure and protocol of 2DIF in practical systems. Besides, it is worth noting that, the Kronecker structure in **Lemma 1** should be viewed as a design and approximation tool rather than a strict requirement. For some non-separable channel models such as [54], the proposed schemes can hold effectiveness by adopting the kernel selection methods provided in this subsection.

VI. SIMULATION RESULTS

In this section, simulations are carried out to verify the effectiveness of the proposed 2DIF based channel estimation schemes.

A. Simulation Setup and Baselines

We consider a single-user MIMO system, where ULAs are equipped on both the BS and the user. The general correlated Rayleigh-fading channel model in (3) is considered to generate \mathbf{H} . In specific, the carrier frequency is set to 3.5 GHz. Following the setup in [41], the elements of the spatial correlation matrices \mathbf{R}_{rx} are calculated via (4), and the spatial correlation matrices \mathbf{R}_{tx} is obtained from the similar process. As a typical example to analyze antenna coupling effect [40], here we consider the dipole antennas with $\lambda/2$ length and $\lambda/100$ width for both transceivers. The antenna spacing is set to be $\lambda/8$. In this context, the spatial-scattering function $f_{\text{rx}}(\varphi, \theta)$ for the BS and that for the user $f_{\text{tx}}(\varphi, \theta)$ are given by $f_{\text{rx}}(\varphi, \theta) = f_{\text{tx}}(\varphi, \theta) = \frac{1.67}{2\pi} \cos^4(\theta)$ [43]. Given these array parameters, the mutual coupling matrices \mathbf{C}_{rx} and \mathbf{C}_{tx} are calculated using Matlab Antenna Toolbox [40]. Otherwise specifically specified, the numbers of transceiver antennas and RF chains are set as: $N_{\text{T}} = 4$, $N_{\text{R}} = 64$, and $N_{\text{RF}} = 4$, respectively. The signal-to-noise ratio (SNR) is defined as $\text{SNR} = \frac{P}{\sigma^2} \mathbb{E}(\|\mathbf{h}\|^2)$, whose default value is set to 10 dB. The evaluation criterion of estimation accuracy is the normalized mean square error (NMSE), which is expressed as $\text{NMSE} = \mathbb{E} \left(\frac{\|\mathbf{h} - \hat{\mathbf{h}}\|^2}{\|\mathbf{h}\|^2} \right)$. The number of channel realizations for kernel training is set to $R = 100$. The default value of pilot length is set to $Q = 48$.

To verify the effectiveness of the proposed 2DIF based channel estimator and TS-2DIF based channel estimator, the following seven schemes are simulated for comparison:

- **LS**: The LS channel estimator is feasible only when observation dimension is no less than the channel dimension, i.e., $QN_{\text{RF}} \geq N_{\text{T}}N_{\text{R}}$. To realize this, we assume the pilot length is $Q = \lceil N_{\text{T}}N_{\text{R}}/N_{\text{RF}} \rceil = 64$, and all combiners/precoders are generated by discrete Fourier transform (DFT) matrices.
- **MMSE**: Under the same setting of LS estimator, the classic MMSE estimator with DFT observation matrices is implemented to recover channel \mathbf{H} via (7).
- **AMP**: Utilizing the channel sparsity in angular domain, the approximate message passing (AMP) method proposed

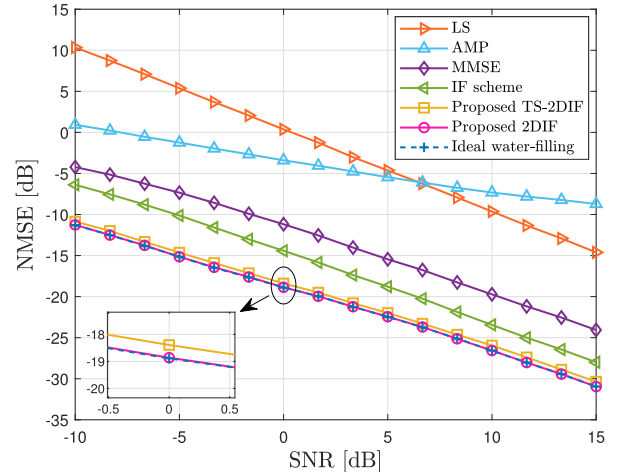


Fig. 3. The NMSE as a function of SNR for different schemes.

in [28] is implemented to estimate channel \mathbf{H} . The combiners and precoders are randomly generated from Gaussian random measurement matrices.

- **IF scheme**: By viewing the considered MIMO system as N_{T} independent SIMO systems, the IF-based channel estimator [36] can be utilized to recover \mathbf{H} in a column-by-column way. Note that, since IF scheme is only applicable to the single-RF-chain case, the pilot length used should be modified as $\lceil QN_{\text{RF}} \rceil = 192$ to ensure that it has the same number of observations with the 2DIF.
- **Proposed 2DIF**: The proposed 2DIF method in **Algorithm 1** is employed to design the precoders and combiners of MIMO system in Fig. 1(a). Based on these designed observation matrices/vectors, the MMSE estimator in (7) is employed to estimate channel \mathbf{H} .
- **Proposed TS-2DIF**: The proposed TS-2DIF method in **Algorithm 3** is employed to design the precoders and combiners of MIMO system in Fig. 1(b). The MMSE estimator in (7) is employed to recover channel \mathbf{H} .
- **Ideal water-filling**: To provide a fundamental performance limit, the ideal (but may not be practically achievable) observation matrices $\{\mathbf{X}_q\}_{q=1}^Q$ are directly obtained by solving (26) via water-filling method. Then, (7) is employed to recover channel \mathbf{H} .

B. Estimation Accuracy Under Statistical Kernel

In this subsection, we consider the ideal case when the statistical kernel $\Sigma_{\mathbf{h}} := \mathbb{E}(\mathbf{h}\mathbf{h}^H)$ can be trained thanks to the known channel models or datasets. Then, $\Sigma_{\mathbf{h}}$ is employed for all required estimators for the channel recovery.

Firstly, we plot the NMSE as a function of SNR in Fig. 3. One can observe that, thanks to the carefully designed observation matrices/vectors, the proposed 2DIF and TS-2DIF schemes remarkably outperform the benchmark schemes in estimation accuracy. The reason is that the proposed methods fully exploit the spatial correlations among the transceiver antennas for channel estimation. In particular, the NMSEs for the proposed 2DIF and TS-2DIF schemes are about 5 dB lower than that for the

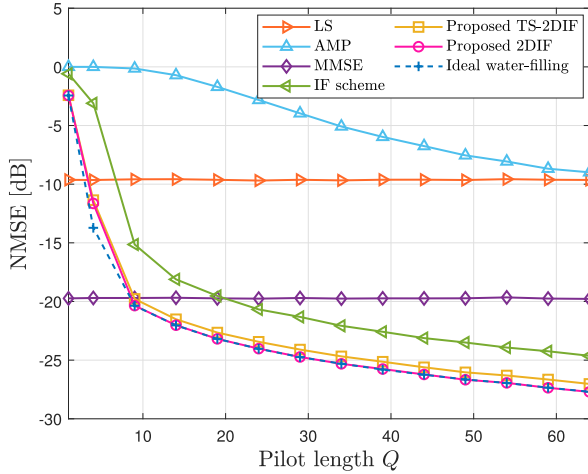


Fig. 4. The NMSE as a function of pilot length Q for different schemes.

IF scheme. It is because the IF schemes realizes the MIMO channel estimation by viewing it as N_T independent SIMO channel estimations, which ignores the spatial correlation of transmitter antennas. Besides, we note that the proposed 2DIF scheme achieves very similar performance to the ideal water-filling scheme. This phenomenon implies that the practical pilot allocation can behave almost the same as the theoretically-optimal “continuous” pilot allocation.

Then, the NMSE versus the number of pilots Q is provided in Fig. 4. One can find that the superiority of the proposed schemes still holds. Although the dimension of the estimated parameters is high as $N_T N_R = 256$, using a small number of pilots $Q = 20$, the NMSEs for the proposed schemes can be lower than -20 dB. In contrast, even if the pilot length is longer than $Q = 60$, the conventional AMP estimator is still unable to achieve such high accuracy. It indicates that utilizing the correlation of compact antennas is of great significance for high-accuracy channel estimation. In addition, observing Figs. 3 and 4, one can conclude that the TS-2DIF scheme can achieve almost the same estimation accuracy as the 2DIF. This indicates that, from the perspective of CSI acquisition, both hybrid MIMO structures in Fig. 1 have no obvious performance gap.

To observe the impact of spatial correlations on the estimation accuracy, we plot the NMSE as a function of the normalized antenna spacing d/λ in Fig. 5. One can observe that, as the antenna spacing decreases, the estimation accuracy of the proposed schemes becomes higher. It is because a smaller antenna spacing leads to stronger spatial correlations, which can provide more prior knowledge for Bayesian estimators. In this case, the more informative kernel allows the proposed schemes to realize more accurate channel estimation. As the antenna spacing increases, due to the reduced channel correlation, the proposed schemes gradually converge to the classical MMSE scheme. However, even if the antenna spacing is $\lambda/2$, the proposed 2DIF and TS-2DIF methods can still hold the superiority. This fact indicates that, for a conventional massive MIMO system, as long as its channels are not i.i.d. Rayleigh-fading (otherwise

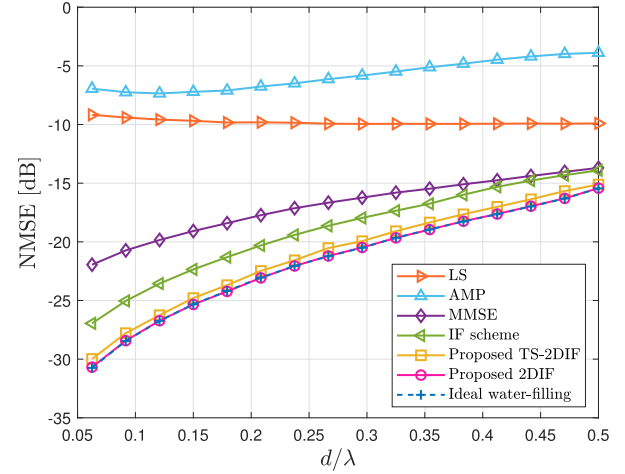


Fig. 5. The NMSE as a function of the normalized antenna spacing d/λ for different schemes.

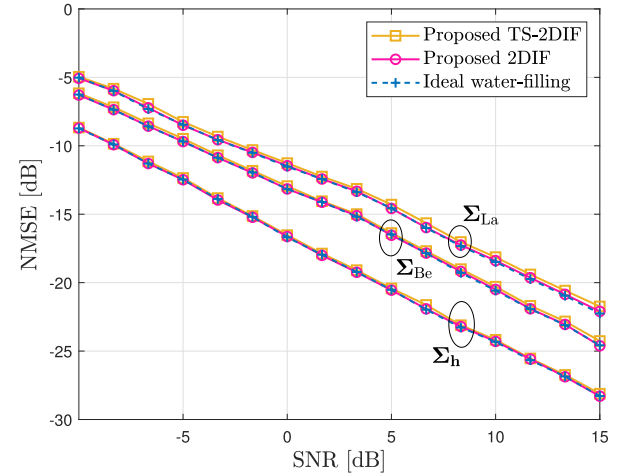


Fig. 6. The NMSE as a function of SNR for different kernels.

$\Sigma_h = \mathbf{I}_{N_T N_R}$), the non-diagonal kernel Σ_h with some structural properties can still contribute to the improvement of the estimation accuracy.

C. Estimation Accuracy Under Artificial Kernels

In practical scenarios where obtaining an explicit channel model or channel dataset is challenging, it is preferred to use an artificial kernel for channel estimation, as discussed in Subsection V-B. In this subsection, two popular artificial kernels, i.e., Laplace kernel Σ_{La} and Bessel kernel Σ_{Be} , are compared with the ideal statistical kernel Σ_h . We plot the NMSE as a function of the SNR in Fig. 6 and the NMSE as a function of the pilot length Q in Fig. 7, respectively.

One can observe that, for each type of kernel, the three proposed schemes exhibit very similar trends in estimation accuracy. This implies that our proposed estimators have the similar robustness for different kernels in channel estimation. Compared to the ideal kernel, the performance losses for both artificial kernels are acceptable. For examples, when $\text{SNR} = 10$

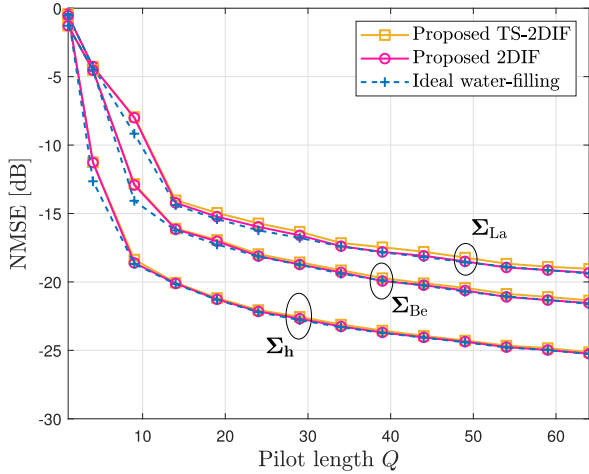
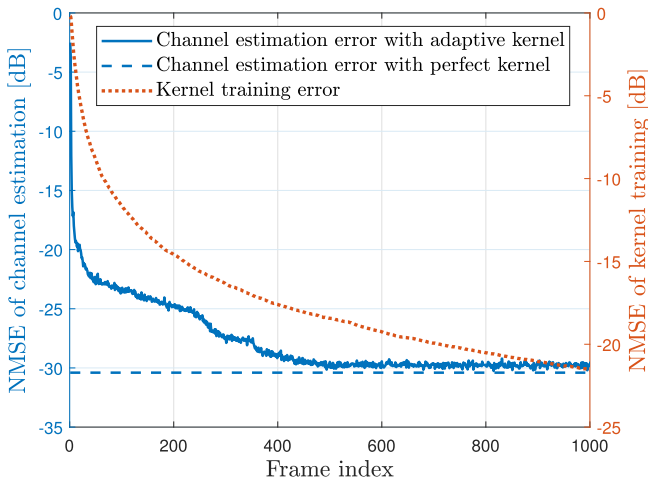

 Fig. 7. The NMSE as a function of pilot length Q for different kernels.


Fig. 8. The NMSE performance of the adaptive kernel training strategy.

dB, the NMSEs for Laplace kernel, Bessel kernel, and statistical kernel are about -18 dB, -20 dB, and -24 dB, respectively. When the pilot length is $Q = 19$, the NMSEs for these three kernels are about -15 dB, -17 dB, and -21 dB, respectively. We can conclude that, even if the real covariance (statistical kernel) is unknown, the proposed channel estimator can still hold their performance advantages by training artificial kernels. This fact encourages the potential applications of the proposed schemes in practice.

D. Estimation Accuracy Under Adaptive Kernel Training

We now evaluate the performance of the proposed adaptive kernel training strategy. Fig. 8 illustrates the NMSE of the estimated channel $\hat{\mathbf{H}}_{t_f}$ by 2DIF, accompanied by the NMSE of learned kernel $\hat{\Sigma}_T^{(t_f)} \otimes \hat{\Sigma}_R^{(t_f)}$ by adaptive training. The frame index t_f increases from 0 to 1000. The other settings are as follows: $N_T = 4$, $N_R = 64$, $d = \frac{\lambda}{8}$, $Q = 48$, SNR = 15 dB. It is observed that, as the frame index increases, the proposed 2DIF algorithm and adaptive kernel training method achieve

a mutual beneficial relationship. The newly estimated channels continuously improve the fidelity of learned kernels, while a more accurate kernel further decreases channel estimation error in future frames. Particularly, the NMSE of channel estimation rapidly declines below -20 dB after 16 frames of kernel training and closely approaches that achieved by the perfect kernel after 400 frames. Moreover, as the channel covariance matrix almost remains unchanged over time, we can keep using the learned kernel after adaptive training to achieve a near-optimal channel estimation performance (see frames 500~1000). This fact validates the feasibility and superiority of the proposed designs in practical situations where no prior information about the kernel is available.

VII. CONCLUSION

By fully exploiting the channel correlations among antennas, this work has proposed a generalized channel estimation framework for densifying MIMO systems, focusing on the design of observation matrices. By maximizing the MI between channels and received pilots, the 2DIF method has been proposed to design observation matrices through jointly optimizing the precoders and combiners. Subsequently, the TS-2DIF method has been proposed to extend the applicability of our framework to the typical hybrid MIMO whose analog combiner is phase-only controllable. Simulation results have validated the superiority of our proposed channel estimation schemes.

APPENDIX A PROOF OF LEMMA 1

Given the channel model in (3) and $\mathbf{h} \equiv \text{vec}(\mathbf{H})$, the vectorized channel can be rewritten as

$$\begin{aligned} \mathbf{h} &= \text{vec} \left(\mathbf{C}_{\text{rx}}^{1/2} \mathbf{R}_{\text{rx}}^{1/2} \mathbf{H}_{\text{iid}} \mathbf{R}_{\text{tx}}^{1/2} \mathbf{C}_{\text{tx}}^{1/2} \right) \\ &= \left(\left(\mathbf{R}_{\text{tx}}^{1/2} \mathbf{C}_{\text{tx}}^{1/2} \right)^T \otimes \left(\mathbf{C}_{\text{rx}}^{1/2} \mathbf{R}_{\text{rx}}^{1/2} \right) \right) \text{vec}(\mathbf{H}_{\text{iid}}), \quad (49) \end{aligned}$$

where the second equation holds since $\text{vec}(\mathbf{ABC}) = (\mathbf{C}^T \otimes \mathbf{A}) \text{vec}(\mathbf{B})$. Then, the covariance of channel \mathbf{h} can be derived as (52) at the bottom of the next page, where (a) holds since $\text{E}(\text{vec}(\mathbf{H}_{\text{iid}}) \text{vec}(\mathbf{H}_{\text{iid}})^H) = \mathbf{I}_{N_R N_T}$, (b) holds due to the commutative law of Kronecker product $(\mathbf{A} \otimes \mathbf{B})(\mathbf{C} \otimes \mathbf{D}) = (\mathbf{AC}) \otimes (\mathbf{BD})$; and (c) holds by defining

$$\Sigma_T = \left(\mathbf{C}_{\text{tx}}^{1/2} \right)^T \mathbf{R}_{\text{tx}}^* \left(\mathbf{C}_{\text{tx}}^{1/2} \right)^*, \quad (50)$$

$$\Sigma_R = \mathbf{C}_{\text{rx}}^{1/2} \mathbf{R}_{\text{rx}} \left(\mathbf{C}_{\text{rx}}^{1/2} \right)^H. \quad (51)$$

One can find that, the matrix Σ_T only depends on the spatial correlation matrix \mathbf{R}_{tx} and the mutual coupling matrix \mathbf{C}_{tx} at the user, while the matrix Σ_R is only associated with the spatial correlation matrix \mathbf{R}_{rx} and the mutual coupling matrix \mathbf{C}_{rx} at the BS. Thus, Σ_T and Σ_R can be viewed as the kernels that characterize the correlation among the transmitter antennas and

that among the receiver antennas, respectively. This completes the proof.

APPENDIX B
PROOF OF LEMMA 2

Using some matrix techniques, the MI $I(\mathbf{y}; \mathbf{h})$ can be rewritten as equation (53), shown at the bottom of the page, where (a) holds since $\det(\mathbf{I} + \mathbf{A}\mathbf{B}) = \det(\mathbf{I} + \mathbf{B}\mathbf{A})$ and $\Xi = \sigma^2 \text{blkdiag}(\mathbf{W}_1^H \mathbf{W}_1, \dots, \mathbf{W}_Q^H \mathbf{W}_Q)$; (b) holds according to the property that $(\mathbf{a} \otimes \mathbf{B}) \mathbf{C} (\mathbf{a}^H \otimes \mathbf{D}) = (\mathbf{a}\mathbf{a}^H) \otimes (\mathbf{B}\mathbf{C}\mathbf{D})$ if all dimensions meet the requirements of matrix multiplications. To find more insights, we perform singular value decomposition (SVD) on all $\{\mathbf{W}_q\}_{q=1}^Q$ and then substitute all decomposition formulas $\mathbf{W}_q = \mathbf{\Pi}_q \mathbf{\Omega}_q \mathbf{\Upsilon}_q^H$ into (53). It is evident that $\mathbf{W}_q (\mathbf{W}_q^H \mathbf{W}_q)^{-1} \mathbf{W}_q^H = \mathbf{\Pi}_q \mathbf{\Pi}_q^H$, thus the MI $I(\mathbf{y}; \mathbf{h})$ can be rewritten as

$$I(\mathbf{y}; \mathbf{h}) = \log_2 \det \left(\mathbf{I}_{N_R N_T} + \frac{1}{\sigma^2} \sum_{q=1}^Q ((\mathbf{v}_q^* \mathbf{v}_q^T) \otimes (\mathbf{\Pi}_q \mathbf{\Pi}_q^H)) \mathbf{\Sigma}_h \right). \quad (54)$$

Observing (54), one can find that the MI $I(\mathbf{y}; \mathbf{h})$ in (9) only relies on the orthogonal matrix $\mathbf{\Pi}_q \in \mathbb{C}^{N \times N_{\text{RF}}}$ decomposed from \mathbf{W}_q for all $q \in \{1, \dots, Q\}$, while it does not depend on any $\mathbf{\Omega}_q$ or $\mathbf{\Upsilon}_q$. It indicates that imposing $\mathbf{W}_q = \mathbf{\Pi}_q$ does not change the value of $I(\mathbf{y}; \mathbf{h})$. As a result, the orthogonality constraint $\mathbf{W}_q^H \mathbf{W}_q = \mathbf{\Pi}_q^H \mathbf{\Pi}_q = \mathbf{I}_{N_{\text{RF}}}$ can be safely introduced into the problem formulation regarding $I(\mathbf{y}; \mathbf{h})$, which completes the proof.

$$\begin{aligned} \mathbf{\Sigma}_h &= \mathbb{E}(\mathbf{h}\mathbf{h}^H) = \mathbb{E} \left(\left(\left(\mathbf{R}_{\text{tx}}^{1/2} \mathbf{C}_{\text{tx}}^{1/2} \right)^T \otimes \left(\mathbf{C}_{\text{rx}}^{1/2} \mathbf{R}_{\text{rx}}^{1/2} \right) \right) \text{vec}(\mathbf{H}_{\text{iid}}) (\text{vec}(\mathbf{H}_{\text{iid}}))^H \left(\left(\mathbf{R}_{\text{tx}}^{1/2} \mathbf{C}_{\text{tx}}^{1/2} \right)^* \otimes \left(\mathbf{C}_{\text{rx}}^{1/2} \mathbf{R}_{\text{rx}}^{1/2} \right)^H \right) \right) \\ &\stackrel{(a)}{=} \left(\left(\mathbf{R}_{\text{tx}}^{1/2} \mathbf{C}_{\text{tx}}^{1/2} \right)^T \otimes \left(\mathbf{C}_{\text{rx}}^{1/2} \mathbf{R}_{\text{rx}}^{1/2} \right) \right) \left(\left(\mathbf{R}_{\text{tx}}^{1/2} \mathbf{C}_{\text{tx}}^{1/2} \right)^* \otimes \left(\mathbf{C}_{\text{rx}}^{1/2} \mathbf{R}_{\text{rx}}^{1/2} \right)^H \right) \\ &\stackrel{(b)}{=} \left(\left(\mathbf{C}_{\text{tx}}^{1/2} \right)^T \mathbf{R}_{\text{tx}}^* \left(\mathbf{C}_{\text{tx}}^{1/2} \right)^* \right) \otimes \left(\mathbf{C}_{\text{rx}}^{1/2} \mathbf{R}_{\text{rx}} \left(\mathbf{C}_{\text{rx}}^{1/2} \right)^H \right) \stackrel{(c)}{=} \mathbf{\Sigma}_T \otimes \mathbf{\Sigma}_R. \end{aligned} \quad (52)$$

$$\begin{aligned} I(\mathbf{y}; \mathbf{h}) &\stackrel{(a)}{=} \log_2 \det \left(\mathbf{I}_{N_R N_T} + \frac{1}{\sigma^2} [\mathbf{v}_1^* \otimes \mathbf{W}_1, \dots, \mathbf{v}_Q^* \otimes \mathbf{W}_Q] \text{blkdiag} \left((\mathbf{W}_1^H \mathbf{W}_1)^{-1}, \dots, (\mathbf{W}_Q^H \mathbf{W}_Q)^{-1} \right) [\mathbf{v}_1^* \otimes \mathbf{W}_1, \dots, \mathbf{v}_Q^* \otimes \mathbf{W}_Q]^H \mathbf{\Sigma}_h \right) \\ &\stackrel{(b)}{=} \log_2 \det \left(\mathbf{I}_{N_R N_T} + \frac{1}{\sigma^2} \sum_{q=1}^Q ((\mathbf{v}_q^* \mathbf{v}_q^T) \otimes (\mathbf{W}_q (\mathbf{W}_q^H \mathbf{W}_q)^{-1} \mathbf{W}_q^H)) \mathbf{\Sigma}_h \right). \end{aligned} \quad (53)$$

$$\begin{aligned} f(\mathbf{v}_{t+1}, \mathbf{W}_{t+1}) &\stackrel{(a)}{=} \log_2 \det \left(\mathbf{I}_{N_{\text{RF}}} + \frac{1}{\sigma^2} \sum_{n_T=1}^{N_T} \sum_{n_R=1}^{N_R} \lambda_{t, n_T, n_R} (\mathbf{v}_{t+1}^T \otimes \mathbf{W}_{t+1}^H) (\mathbf{a}_{n_T} \otimes \mathbf{b}_{n_R}) (\mathbf{a}_{n_T}^H \otimes \mathbf{b}_{n_R}^H) (\mathbf{v}_{t+1} \otimes \mathbf{W}_{t+1}) \right) \\ &\stackrel{(b)}{=} \log_2 \det \left(\mathbf{I}_{N_{\text{RF}}} + \frac{1}{\sigma^2} \sum_{n_T=1}^{N_T} \sum_{n_R=1}^{N_R} \lambda_{t, n_T, n_R} |\mathbf{a}_{n_T}^H \mathbf{v}_{t+1}^*|^2 \mathbf{W}_{t+1}^H \mathbf{b}_{n_R} \mathbf{b}_{n_R}^H \mathbf{W}_{t+1} \right). \end{aligned} \quad (56)$$

APPENDIX C
PROOF OF MI INCREMENT $I(\bar{\mathbf{y}}_{t+1}; \mathbf{h}) - I(\bar{\mathbf{y}}_t; \mathbf{h})$

Using some matrix partition operations, the MI $I(\bar{\mathbf{y}}_{t+1}; \mathbf{h})$ can be rewritten as

$$\begin{aligned} I(\bar{\mathbf{y}}_{t+1}; \mathbf{h}) &\stackrel{(a)}{=} \log_2 \det \left(\mathbf{I}_{N_{\text{RF}} Q} + \frac{1}{\sigma^2} \bar{\mathbf{X}}_{t+1}^H \mathbf{\Sigma}_h \bar{\mathbf{X}}_{t+1} \right) \\ &= \log_2 \det \begin{bmatrix} \mathbf{I}_{N_{\text{RF}} t} + \frac{1}{\sigma^2} \bar{\mathbf{X}}_t^H \mathbf{\Sigma}_h \bar{\mathbf{X}}_t & \frac{1}{\sigma^2} \bar{\mathbf{X}}_t^H \mathbf{\Sigma}_h \mathbf{X}_{t+1} \\ \frac{1}{\sigma^2} \mathbf{X}_{t+1}^H \mathbf{\Sigma}_h \bar{\mathbf{X}}_t & \mathbf{I}_{N_{\text{RF}}} + \frac{1}{\sigma^2} \mathbf{X}_{t+1}^H \mathbf{\Sigma}_h \mathbf{X}_{t+1} \end{bmatrix} \\ &\stackrel{(b)}{=} \log_2 \det \begin{bmatrix} \mathbf{I}_{N_{\text{RF}} t} + \frac{1}{\sigma^2} \bar{\mathbf{X}}_t^H \mathbf{\Sigma}_h \bar{\mathbf{X}}_t & \frac{1}{\sigma^2} \bar{\mathbf{X}}_t^H \mathbf{\Sigma}_h \mathbf{X}_{t+1} \\ \mathbf{0}_{N_{\text{RF}} \times N_{\text{RF}} t} & \mathbf{I}_{N_{\text{RF}}} + \frac{1}{\sigma^2} \mathbf{X}_{t+1}^H \mathbf{\Sigma}_t \mathbf{X}_{t+1} \end{bmatrix} \\ &= I(\bar{\mathbf{y}}_t; \mathbf{h}) + \log_2 \det \left(\mathbf{I}_{N_{\text{RF}}} + \frac{1}{\sigma^2} \mathbf{X}_{t+1}^H \mathbf{\Sigma}_t \mathbf{X}_{t+1} \right), \end{aligned} \quad (55)$$

where (a) holds since according to **Lemma 2** and (b) holds by performing matrix triangularization. In particular, $\mathbf{\Sigma}_t$ is given by $\mathbf{\Sigma}_t = \mathbf{\Sigma}_h - \mathbf{\Sigma}_h \bar{\mathbf{X}}_t (\bar{\mathbf{X}}_t^H \mathbf{\Sigma}_h \bar{\mathbf{X}}_t + \sigma^2 \mathbf{I}_{N_{\text{RF}} t})^{-1} \bar{\mathbf{X}}_t^H \mathbf{\Sigma}_h$, which completes the proof.

APPENDIX D
PROOF OF LEMMA 3

The key idea of the proof is to rewrite the $\bar{\mathbf{X}}_t$ -related terms in (18) as $\mathbf{\Sigma}_h \bar{\mathbf{X}}_t = \mathbf{\Sigma}_h [\bar{\mathbf{X}}_{t-1}, \mathbf{X}_t]$ and

$$\bar{\mathbf{X}}_t^H \mathbf{\Sigma}_h \bar{\mathbf{X}}_t = \begin{bmatrix} \bar{\mathbf{X}}_{t-1}^H \mathbf{\Sigma}_h \bar{\mathbf{X}}_{t-1} & \bar{\mathbf{X}}_{t-1}^H \mathbf{\Sigma}_h \mathbf{X}_t \\ \mathbf{X}_t^H \mathbf{\Sigma}_h \bar{\mathbf{X}}_{t-1} & \mathbf{X}_t^H \mathbf{\Sigma}_h \mathbf{X}_t \end{bmatrix}. \quad (57)$$

Then, using the Schur's matrix inversion formula to expand the term $(\bar{\mathbf{X}}_t^H \mathbf{\Sigma}_h \bar{\mathbf{X}}_t + \sigma^2 \mathbf{I}_{N_{\text{RF}} t})^{-1}$ in (18), the following recursion formula of can be obtained:

$$\mathbf{\Sigma}_{t+1} = \mathbf{\Sigma}_t - \mathbf{\Sigma}_t \mathbf{X}_{t+1} (\mathbf{X}_{t+1}^H \mathbf{\Sigma}_t \mathbf{X}_{t+1} + \sigma^2 \mathbf{I}_{N_{\text{RF}}})^{-1} \mathbf{X}_{t+1}^H \mathbf{\Sigma}_t, \quad (58)$$

When $\mathbf{X}_{t+1} = \sqrt{P}\mathbf{U}_t(:, [1, \dots, N_{\text{RF}}])$, we have $\boldsymbol{\Sigma}_t \mathbf{X}_{t+1} = \mathbf{X}_{t+1} \text{diag}(\lambda_1(\boldsymbol{\Sigma}_t), \dots, \lambda_{N_{\text{RF}}}(\boldsymbol{\Sigma}_t))$ and $\mathbf{X}_{t+1}^H \boldsymbol{\Sigma}_t \mathbf{X}_{t+1} = P \text{diag}(\lambda_1(\boldsymbol{\Sigma}_t), \dots, \lambda_{N_{\text{RF}}}(\boldsymbol{\Sigma}_t))$. Thus, the following equality holds:

$$\boldsymbol{\Sigma}_{t+1} = \mathbf{U}_t \boldsymbol{\Lambda}_t \mathbf{U}_t^H - \mathbf{X}_{t+1} \text{diag}\left(\frac{\lambda_1^2(\boldsymbol{\Sigma}_t)}{P\lambda_1(\boldsymbol{\Sigma}_t) + \sigma^2}, \dots, \frac{\lambda_{N_{\text{RF}}}^2(\boldsymbol{\Sigma}_t)}{P\lambda_{N_{\text{RF}}}(\boldsymbol{\Sigma}_t) + \sigma^2}\right) \mathbf{X}_{t+1}^H. \quad (59)$$

Given that $\mathbf{X}_{t+1} \text{diag}\left(\frac{\lambda_1^2(\boldsymbol{\Sigma}_t)}{P\lambda_1(\boldsymbol{\Sigma}_t) + \sigma^2}, \dots, \frac{\lambda_{N_{\text{RF}}}^2(\boldsymbol{\Sigma}_t)}{P\lambda_{N_{\text{RF}}}(\boldsymbol{\Sigma}_t) + \sigma^2}\right) \mathbf{X}_{t+1}^H = \mathbf{U}_t \text{diag}\left(\frac{P\lambda_1^2(\boldsymbol{\Sigma}_t)}{P\lambda_1(\boldsymbol{\Sigma}_t) + \sigma^2}, \dots, \frac{P\lambda_{N_{\text{RF}}}^2(\boldsymbol{\Sigma}_t)}{P\lambda_{N_{\text{RF}}}(\boldsymbol{\Sigma}_t) + \sigma^2}, \underbrace{0, \dots, 0}_{N_{\text{R}}N_{\text{T}} - N_{\text{RF}}}\right) \mathbf{U}_t^H$

and $\boldsymbol{\Sigma}_t = \mathbf{U}_t \boldsymbol{\Lambda}_t \mathbf{U}_t^H$, the equality in (20) can be derived from (59), which completes the proof.

APPENDIX E PROOF OF COROLLARY 1

According to **Lemma 1** and equality $(\mathbf{A}\mathbf{B}\mathbf{A}^H) \otimes (\mathbf{C}\mathbf{D}\mathbf{C}^H) = (\mathbf{A} \otimes \mathbf{C})(\mathbf{B} \otimes \mathbf{D})(\mathbf{A}^H \otimes \mathbf{C}^H)$, the kernel $\boldsymbol{\Sigma}_{\mathbf{h}}$ can be decomposed as

$$\begin{aligned} \boldsymbol{\Sigma}_{\mathbf{h}} &= (\mathbf{U}_{\text{T}} \boldsymbol{\Lambda}_{\text{T}} \mathbf{U}_{\text{T}}^H) \otimes (\mathbf{U}_{\text{T}} \boldsymbol{\Lambda}_{\text{T}} \mathbf{U}_{\text{T}}^H) \\ &= \underbrace{(\mathbf{U}_{\text{T}} \otimes \mathbf{U}_{\text{R}})}_{\mathbf{U}_0} \underbrace{(\boldsymbol{\Lambda}_{\text{T}} \otimes \boldsymbol{\Lambda}_{\text{R}})}_{\text{Eigenvalue matrix}} (\mathbf{U}_{\text{T}}^H \otimes \mathbf{U}_{\text{R}}^H) \\ &= \sum_{n_{\text{T}}=1}^{N_{\text{T}}} \sum_{n_{\text{R}}=1}^{N_{\text{R}}} \alpha_{n_{\text{T}}} \beta_{n_{\text{R}}} (\mathbf{a}_{n_{\text{T}}} \otimes \mathbf{b}_{n_{\text{R}}}) (\mathbf{a}_{n_{\text{T}}}^H \otimes \mathbf{b}_{n_{\text{R}}}^H), \quad (60) \end{aligned}$$

Based on (60), one can verify without difficulty that (21) is exactly the eigenvalue decomposition of $\boldsymbol{\Sigma}_{\mathbf{h}}$, which completes the proof.

APPENDIX F PROOF OF LEMMA 4

Given the new constraints $\mathbf{v}_{t+1} \in \left\{ \sqrt{P} \mathbf{a}_{n_{\text{T}}}^* \right\}_{n_{\text{T}}=1}^{N_{\text{T}}}$ and $\mathbf{w}_{t+1,k} \in \left\{ \mathbf{b}_{n_{\text{R}}} \right\}_{n_{\text{R}}=1}^{N_{\text{RF}}}$ for all $k \in \{1, \dots, N_{\text{RF}}\}$, problem (19) can be reorganized as

$$\begin{aligned} \max_{\mathbf{v}_{t+1}, \mathbf{W}_{t+1}} & f(\mathbf{v}_{t+1}, \mathbf{W}_{t+1}) \\ \text{s.t.} & \mathbf{v}_{t+1} \in \left\{ \sqrt{P} \mathbf{a}_{n_{\text{T}}}^* \right\}_{n_{\text{T}}=1}^{N_{\text{T}}}, \\ & \mathbf{w}_{t+1,k} \in \left\{ \mathbf{b}_{n_{\text{R}}} \right\}_{n_{\text{R}}=1}^{N_{\text{RF}}}, \forall k \in \{1, \dots, N_{\text{RF}}\}, \\ & \mathbf{w}_{t+1,k} \neq \mathbf{w}_{t+1,k'}, \forall k \neq k', \quad (61) \end{aligned}$$

where the objective function is given in (56). Shown at bottom of the pervious page, in which (a) holds according to the definition in (22) and (b) holds by utilizing the property to the property that $(\mathbf{A} \otimes \mathbf{B})(\mathbf{C} \otimes \mathbf{D}) = (\mathbf{A}\mathbf{C}) \otimes (\mathbf{B}\mathbf{D})$. Note that, the constraint $\mathbf{w}_{t+1,k} \neq \mathbf{w}_{t+1,k'}$ for all $k \neq k'$ in (61) ensures the orthogonality of \mathbf{W}_{t+1} . Observing (61), one can find that our goal becomes finding optimal indexes n_{T} and $\{n_{\text{R},k}\}_{k=1}^{N_{\text{RF}}}$ that maximize the MI increment $f(\mathbf{v}_{t+1}, \mathbf{W}_{t+1})$. Assuming

that the optimal indexes are expressed by $n_{\text{T}}^{\text{opt}}$ and $\{n_{\text{R},k}^{\text{opt}}\}_{k=1}^{N_{\text{RF}}}$, the optimal precoder and the optimal combiner are

$$\mathbf{v}_{t+1}^{\text{opt}} = \sqrt{P} \mathbf{a}_{n_{\text{T}}^{\text{opt}}}^* \text{ and } \mathbf{W}_{t+1}^{\text{opt}} = \left[\mathbf{b}_{n_{\text{R},1}^{\text{opt}}}, \dots, \mathbf{b}_{n_{\text{R},N_{\text{RF}}}^{\text{opt}}} \right], \quad (62)$$

respectively. Then, we have

$$\mathbf{a}_{n_{\text{T}}}^H (\mathbf{v}_{t+1}^{\text{opt}})^* = \begin{cases} \sqrt{P}, & n_{\text{T}} = n_{\text{T}}^{\text{opt}} \\ 0, & \text{else,} \end{cases} \quad (63a)$$

$$\mathbf{b}_{n_{\text{R}}}^H \mathbf{W}_{t+1}^{\text{opt}} = \begin{cases} \mathbf{e}_{n_{\text{R}}}^T, & n_{\text{R}} \in \{n_{\text{R},k}^{\text{opt}}\}_{k=1}^{N_{\text{RF}}} \\ \mathbf{0}_{N_{\text{RF}}}^T, & \text{else,} \end{cases} \quad (63b)$$

where $\mathbf{e}_{n_{\text{R}}}$ denotes an N_{RF} -dimensional vector whose n_{R} -th entry is one and the other entries are zero. By substituting (63) into (61), the optimal MI increment $f(\mathbf{v}_{t+1}^{\text{opt}}, \mathbf{W}_{t+1}^{\text{opt}})$ can be expressed by

$$\begin{aligned} f(\mathbf{v}_{t+1}^{\text{opt}}, \mathbf{W}_{t+1}^{\text{opt}}) &= \log_2 \det \left(\mathbf{I}_{N_{\text{RF}}} + \frac{P}{\sigma^2} \sum_{n_{\text{R}}} \lambda_{t,n_{\text{T}}^{\text{opt}},n_{\text{R}}} (\mathbf{W}_{t+1}^{\text{opt}})^H \mathbf{b}_{n_{\text{R}}} \mathbf{b}_{n_{\text{R}}}^H \mathbf{W}_{t+1}^{\text{opt}} \right) \\ &= \log_2 \det \left(\mathbf{I}_{N_{\text{RF}}} + \frac{P}{\sigma^2} \text{diag} \left(\lambda_{t,n_{\text{T}}^{\text{opt}},n_{\text{R},1}^{\text{opt}}}, \dots, \lambda_{t,n_{\text{T}}^{\text{opt}},n_{\text{R},N_{\text{RF}}}^{\text{opt}}} \right) \right) \\ &= \sum_{k=1}^{N_{\text{RF}}} \log_2 \left(1 + \frac{P \lambda_{t,n_{\text{T}}^{\text{opt}},n_{\text{R},k}^{\text{opt}}}}{\sigma^2} \right), \quad (64) \end{aligned}$$

which only relies on the eigenvalues of $\boldsymbol{\Sigma}_t$. In this context, the problem becomes finding n_{T} and $\{n_{\text{R},k}\}_{k=1}^{N_{\text{RF}}}$ that maximize $f(\mathbf{v}_{t+1}, \mathbf{W}_{t+1})$, as formulated in (24). This completes the proof.

REFERENCES

- [1] Z. Zhang and M. Cui, "Two-dimensional ice filling based channel estimation in densifying MIMO systems," in *Proc. IEEE Int. Conf. Commun. (IEEE ICC'25)*, Jun. 2025, pp. 4750–4755.
- [2] C. Huang et al., "Holographic MIMO surfaces for 6G wireless networks: Opportunities, challenges, and trends," *IEEE Wireless Commun.*, vol. 27, no. 5, pp. 118–125, Oct. 2020.
- [3] R. Deng et al., "Reconfigurable holographic surfaces for ultra-massive MIMO in 6G: Practical design, optimization and implementation," *IEEE J. Sel. Areas Commun.*, vol. 41, no. 8, pp. 2367–2379, Aug. 2023.
- [4] Z. Zhang and L. Dai, "Pattern-division multiplexing for multi-user continuous-aperture MIMO," *IEEE J. Sel. Areas Commun.*, vol. 41, no. 8, pp. 2350–2366, Aug. 2023.
- [5] I. Kanbaz, O. Yurduseven, and M. Matthaiou, "Super-directive antenna arrays: How many elements do we need?" in *Proc. IEEE Wireless Commun. Netw. Conf. (IEEE WCNC'24)*, 2024, pp. 1–6.
- [6] Z. Zhang and L. Dai, "Reconfigurable intelligent surfaces for 6G: Nine fundamental issues and one critical problem," *Tsinghua Sci. Technol.*, vol. 28, no. 5, pp. 929–939, Oct. 2023.
- [7] K.-K. Wong and K.-F. Tong, "Fluid antenna multiple access," *IEEE Trans. Wireless Commun.*, vol. 21, no. 7, pp. 4801–4815, Jul. 2021.
- [8] D. González-Ovejero, G. Minatti, G. Chattopadhyay, and S. Maci, "Multibeam by metasurface antennas," *IEEE Trans. Antennas Propag.*, vol. 65, no. 6, pp. 2923–2930, Jun. 2017.
- [9] R.-B. Hwang, "Binary meta-hologram for a reconfigurable holographic metamaterial antenna," *Sci. Rep.*, vol. 10, no. 1, May 2020, Art. no. 8586.
- [10] C. Liaskos, S. Nie, A. Tsioliaridou, A. Pitsillides, S. Ioannidis, and I. Akyildiz, "A new wireless communication paradigm through software-controlled metasurfaces," *IEEE Commun. Mag.*, vol. 56, no. 9, pp. 162–169, Sep. 2018.
- [11] M. Liu et al., "Deeply subwavelength metasurface resonators for terahertz wavefront manipulation," *Adv. Opt. Mater.*, vol. 7, no. 21, 2019, Art. no. 1900736.

- [12] T. Gong et al., "Holographic MIMO communications: Theoretical foundations, enabling technologies, and future directions," *IEEE Commun. Surveys Tuts.*, vol. 26, no. 1, pp. 196–257, 1st Quart. 2024.
- [13] E. Basar et al., "Reconfigurable intelligent surfaces for 6G: Emerging hardware architectures, applications, and open challenges," *IEEE Veh. Technol. Mag.*, vol. 19, no. 3, pp. 27–47, Sep. 2024.
- [14] L. Wei et al., "Multi-user holographic MIMO surfaces: Channel modeling and spectral efficiency analysis," *IEEE J. Sel. Topics Signal Process.*, vol. 16, no. 5, pp. 1112–1124, Aug. 2022.
- [15] J. An, C. Yuen, C. Huang, M. Debbah, H. Vincent Poor, and L. Hanzo, "A tutorial on holographic MIMO communications—part I: Channel modeling and channel estimation," *IEEE Commun. Lett.*, vol. 27, no. 7, pp. 1664–1668, Jul. 2023.
- [16] J. An, C. Yuen, C. Huang, M. Debbah, H. V. Poor, and L. Hanzo, "A tutorial on holographic MIMO communications—Part II: Performance analysis and holographic beamforming," *IEEE Commun. Lett.*, vol. 27, no. 7, pp. 1669–1673, Jul. 2023.
- [17] Y. Liu, M. Zhang, T. Wang, A. Zhang, and M. Debbah, "Densifying MIMO: Channel modeling, physical constraints, and performance evaluation for holographic communications," *IEEE J. Sel. Areas Commun.*, vol. 42, no. 6, pp. 1504–1518, Jun. 2024.
- [18] M. Di Renzo, D. Dardari, and N. Decarli, "LoS MIMO-arrays vs. LoS MIMO-surfaces," in *Proc. 17th Eur. Conf. Antennas Propag. (EuCAP'23)*, 2023, pp. 1–5.
- [19] J. Xie, H. Yin, and L. Han, "A genetic algorithm based superdirective beamforming method under excitation power range constraints," 2023, *arXiv:2307.02063*.
- [20] M. Akrouf, V. Shyianov, F. Bellili, A. Mezghani, and R. W. Heath, "Super-wideband massive MIMO," *IEEE J. Sel. Areas Commun.*, vol. 41, no. 8, pp. 2414–2430, Aug. 2023.
- [21] D. Tse and P. Viswanath, *Fundamentals of Wireless Communication*. Cambridge, U.K.: Cambridge Univ. Press, 2005.
- [22] M. Cui, Z. Wu, Y. Lu, X. Wei, and L. Dai, "Near-field MIMO communications for 6G: Fundamentals, challenges, potentials, and future directions," *IEEE Commun. Mag.*, vol. 61, no. 1, pp. 40–46, Jan. 2022.
- [23] T. Zheng, J. Zhu, Q. Yu, Y. Yan, and L. Dai, "Coded beam training," *IEEE J. Sel. Areas Commun.*, vol. 43, no. 3, pp. 928–943, Mar. 2025.
- [24] S. M. Kay, *Fundamentals of Statistical Signal Processing: Estimation Theory*, Englewood Cliff s, NJ, USA: Prentice-Hall, Inc., 1993.
- [25] C. Huang, L. Liu, C. Yuen, and S. Sun, "Iterative channel estimation using LSE and sparse message passing for mmWave MIMO systems," *IEEE Trans. Signal Process.*, vol. 67, no. 1, pp. 245–259, Jan. 2019.
- [26] M. Ke, Z. Gao, Y. Wu, X. Gao, and R. Schober, "Compressive sensing-based adaptive active user detection and channel estimation: Massive access meets massive MIMO," *IEEE Trans. Signal Process.*, vol. 68, pp. 764–779, Jan. 2020.
- [27] M. Cui and L. Dai, "Channel estimation for extremely large-scale MIMO: Far-field or near-field?" *IEEE Trans. Commun.*, vol. 70, no. 4, pp. 2663–2677, Apr. 2022.
- [28] S. Rangan, P. Schniter, and A. K. Fletcher, "Vector approximate message passing," *IEEE Trans. Inf. Theory*, vol. 65, no. 10, pp. 6664–6684, Oct. 2019.
- [29] Y. Jin, J. Zhang, S. Jin, and B. Ai, "Channel estimation for cell-free mmWave massive MIMO through deep learning," *IEEE Trans. Veh. Technol.*, vol. 68, no. 10, pp. 10325–10329, Oct. 2019.
- [30] Z. Wan, Z. Gao, B. Shim, K. Yang, G. Mao, and M.-S. Alouini, "Compressive sensing based channel estimation for millimeter-wave full-dimensional MIMO with lens-array," *IEEE Trans. Veh. Technol.*, vol. 69, no. 2, pp. 2337–2342, Feb. 2020.
- [31] X. Ma, Z. Gao, F. Gao, and M. Di Renzo, "Model-driven deep learning based channel estimation and feedback for millimeter-wave massive hybrid MIMO systems," *IEEE J. Sel. Areas Commun.*, vol. 39, no. 8, pp. 2388–2406, Aug. 2021.
- [32] A. Pizzo and A. Lozano, "Mutual coupling in holographic MIMO: Physical modeling and information-theoretic analysis," *IEEE J. Sel. Areas Inf. Theory*, vol. 6, pp. 111–126, May 2025.
- [33] C. Williams and C. Rasmussen, "Gaussian processes for regression," in *Adv. Neural Inf. Process. Syst.*, vol. 8, 1995, pp. 514–520.
- [34] W. K. New, K.-K. Wong, H. Xu, F. R. Ghadi, R. Murch, and C.-B. Chae, "Channel estimation and reconstruction in fluid antenna system: Oversampling is essential," *IEEE J. Sel. Areas Commun.*, vol. 24, no. 1, pp. 309–322, Jan. 2025.
- [35] J. An, C. Yuen, C. Huang, M. Debbah, H. V. Poor, and L. Hanzo, "A tutorial on holographic MIMO communications—Part III: Open opportunities and challenges," *IEEE Commun. Lett.*, vol. 27, no. 7, pp. 1674–1678, Jul. 2023.
- [36] M. Cui, Z. Zhang, L. Dai, and K. Huang, "Ice-filling: Near-optimal channel estimation for dense array systems," *IEEE Trans. Wireless Commun.*, vol. 24, no. 10, pp. 8551–8564, Oct. 2025.
- [37] Y. Wu, J. P. Linnartz, J. Bergmans, and S. Attallah, "Effects of antenna mutual coupling on the performance of MIMO systems," in *Proc. 29th Symp. Inf. Theory in Benelux. Citeseer*, 2008, pp. 1–8.
- [38] K.-H. Chen and J.-F. Kiang, "Effect of mutual coupling on the channel capacity of MIMO systems," *IEEE Trans. Veh. Technol.*, vol. 65, no. 1, pp. 398–403, Jan. 2016.
- [39] N. Kolomvakis and E. Björnson, "Exploiting mutual coupling characteristics for channel estimation in holographic MIMO," in *Proc. IEEE Global Commun. Conf. (IEEE GLOBECOM'24)*, 2024, pp. 3570–3575.
- [40] MathWorks, "Antenna toolbox," Natick, MA, 2024. [Online]. Available: <https://www.mathworks.com/products/antenna.html>
- [41] A. Pizzo, L. Sanguinetti, and T. L. Marzetta, "Fourier plane-wave series expansion for holographic MIMO communications," *IEEE Trans. Wireless Commun.*, vol. 21, no. 9, pp. 6890–6905, Sep. 2022.
- [42] I. Gupta and A. Ksienski, "Effect of mutual coupling on the performance of adaptive arrays," *IEEE Trans. Antennas Propag.*, vol. 31, no. 5, pp. 785–791, Sep. 1983.
- [43] C. A. Balanis, *Antenna Theory: analysis and Design*. Hoboken, NJ, USA: John Wiley & Sons, 2016.
- [44] K. Werner and M. Jansson, "Estimating MIMO channel covariances from training data under the Kronecker model," *Signal Process.*, vol. 89, no. 1, pp. 1–13, Jan. 2009.
- [45] S. Park and R. W. Heath, "Spatial channel covariance estimation for the hybrid MIMO architecture: A compressive sensing-based approach," *IEEE Trans. Wireless Commun.*, vol. 17, no. 12, pp. 8047–8062, Dec. 2018.
- [46] K. Upadhyaya and S. A. Vorobyov, "Covariance matrix estimation for massive MIMO," *IEEE Signal Process. Lett.*, vol. 25, no. 4, pp. 546–550, Apr. 2018.
- [47] M. B. Khalilsarai, T. Yang, S. Haghghatshoar, and G. Caire, "Structured channel covariance estimation from limited samples in massive MIMO," in *Proc. IEEE Int. Conf. Commun. (IEEE ICC'20)*, Jun. 2020, pp. 1–7.
- [48] B. Clarke and A. Barron, "Information-theoretic asymptotics of Bayes methods," *IEEE Trans. Inf. Theory*, vol. 36, no. 3, pp. 453–471, May 1990.
- [49] E. Zhang and C. Huang, "On achieving optimal rate of digital precoder by rf-baseband codesign for MIMO systems," in *Proc. IEEE 80th Veh. Technol. Conf. (IEEE VTC'14-Fall)* 2014, pp. 1–5.
- [50] E. Vlachos, A. Kaushik, Y. C. Eldar, and G. C. Alexandropoulos, "Time-domain channel estimation for extremely large MIMO THz communication systems under dual-wideband fading conditions," 2023, *arXiv:2310.14745*.
- [51] E. Vlachos, G. C. Alexandropoulos, and J. Thompson, "Wideband MIMO channel estimation for hybrid beamforming millimeter wave systems via random spatial sampling," *IEEE J. Sel. Topics Signal Process.*, vol. 13, no. 5, pp. 1136–1150, Sep. 2019.
- [52] J. Zhu, X. Su, Z. Wan, L. Dai, and T. J. Cui, "The benefits of electromagnetic information theory for channel estimation," in *Proc. IEEE Int. Conf. Commun. (IEEE ICC'24)*, Jun. 2024, pp. 4869–4874.
- [53] N. Srinivas, A. Krause, S. M. Kakade, and M. W. Seeger, "Information-theoretic regret bounds for Gaussian process optimization in the bandit setting," *IEEE Trans. Inf. Theory*, vol. 58, no. 5, pp. 3250–3265, May 2012.
- [54] X. Zhang, S. Song, and K. B. Letaief, "Fundamental limits of non-centered non-separable channels and their application in holographic MIMO communications," *IEEE Trans. Inf. Theory*, vol. 71, no. 9, pp. 6870–6894, Sep. 2025.