# A Framework on Complex Matrix Derivatives with Special Structure Constraints for Wireless Systems

Xin Ju, Shiqi Gong, Nan Zhao, *Senior Member, IEEE*, Chengwen Xing, Arumugam Nallanathan, *Fellow, IEEE*, and Dusit Niyato, *Fellow, IEEE*

*Abstract*—Matrix-variate optimization plays a central role in advanced wireless system designs. In this paper, we aim to explore optimal solutions of matrix variables under two special structure constraints using complex matrix derivatives, including diagonal structure constraints and constant modulus constraints, both of which are closely related to the state-of-the-art wireless applications. Specifically, for diagonal structure constraints mostly considered in the uplink multi-user single-input multiple-output (MU-SIMO) system and the amplitude-adjustable intelligent reflecting surface (IRS)-aided multiple-input multiple-output (MIMO) system, the capacity maximization problem, the mean-squared error (MSE) minimization problem and their variants are rigorously investigated. By leveraging complex matrix derivatives, the optimal solutions of these problems are directly obtained in closed forms. Nevertheless, for constant modulus constraints with the intrinsic nature of element-wise decomposability, which are often seen in the hybrid analog-digital MIMO system and the fully-passive IRS-aided MIMO system, we firstly explore inherent structures of the element-wise phase derivatives associated with different optimization problems. Then, we propose a novel alternating optimization (AO) algorithm with the aid of several arbitrary feasible solutions, which avoids the complicated matrix inversion and matrix factorization involved in conventional element-wise iterative algorithms. Numerical simulations reveal that the proposed algorithm can dramatically reduce the computational complexity without loss of system performance.

*Index Terms*—Complex matrix derivatives, special structure constraints, matrix-variate optimization, hybrid analog-digital system, intelligent reflecting surface.

## I. INTRODUCTION

Multi-antenna technology opens a new era for wireless communications due to its effective utilization of limited spatial resources [1]–[3]. From the mathematical viewpoint, the deployment of multi-antenna arrays at transceivers generally leads to matrix-variate optimization problems [4]–[6]. Specifically, in the typical multiple-input multiple-output (MIMO) communication systems, the transmit beamformer optimization and the receive equalizer optimization can be both modeled as matrix-variate optimization problems [7]–[9]. Compared to scalar-variate optimization, matrix-variate optimization is generally more challenging to tackle because

it inherently involves complex matrix operations, including matrix determinant, inversion, matrix decomposition and so on. In fact, with the development of wireless communications, many matrix-variate optimization problems with special structure constraints such as symmetric, diagonal and constant modulus structure constraints are emerging, which are closely related to the state-of-the-art wireless systems equipped with multi-antenna transceiver antenna arrays.

In general, a structure of the matrix variable strongly depends on three factors, namely, network architectures, frequency bands and communication demands. First, the distributed network architecture has been studied, since its involved distributed antenna arrays are capable of increasing spatial diversity gain and extending communication coverage, as compared with the centralized counterpart [10]. In this distributed network, the corresponding matrix variable usually has a diagonal structure. Second, for high-frequency millimeter wave (mmWave) and terahertz (THz) communications, a hybrid analog-digital transceiver structure has been regarded as an economic and effective way to achieve a large array gain [11], [12], in which the analog beamforming matrix is usually subject to the nonconvex and intractable constant modulus constraints. Third, for smartly reconfiguring the wireless environment in a cost-effective manner, intelligent reflecting surfaces (IRSs) composed of a large number of passive reflecting elements have attracted a lot of attention recently [13]–[15]. Considering different levels of hardware implementation, there are two main types of IRS structures that are widely studied, i.e., the amplitude-adjustable IRS and the fully-passive IRS. Note that the reflection matrices of these two IRSs can be mathematically modeled as diagonal matrices. In particular, for the fully-passive IRS, the corresponding reflection matrix is additionally subject to the nonconvex constant modulus constraint. Building upon the above discussions, it is clear that matrix-variate optimization problems with special structure constraints have been widely considered in the state-of-the-art wireless systems. Therefore, it is essential to develop a framework for optimization algorithms with guaranteed performance and low complexity for the matrix-variate optimization.

Currently, there have been many common popular algorithms for solving matrix-variate optimization problems, such as the Karush-Kuhn-Tucker (KKT)-based algorithm [16]–[19], the block coordinate descent (BCD) algorithm [20] and the majorization-minimization (MM)-based algorithm [21]. It is well-known that for convex matrix-variate problems, the KKT-based algorithm is able to directly derive the optimal structures of matrices. Generally, the symmetric structure constraints of matrix variables can be implicitly satisfied by the derived optimal closed-form solutions [22]. Moreover,

X. Ju and C. Xing are with the School of Information and Electronics, Beijing Institute of Technology, Beijing 100081, China (E-mails: xinjubit@gmail.com; chengwenxing@ieee.org).

S. Gong is with the School of Cyberspace Science and Technology, Beijing Institute of Technology, Beijing 100081, China (E-mail: gsqyx@163.com).

N. Zhao is with the School of Information and Communication Engineering, Dalian University of Technology, Dalian 116024, China (E-mail: zhaonan@dlut.edu.cn).

Arumugam Nallanathan is with the QueenMary University of London, E14NS London, U.K. (E-mail: a.nallanathan@qmul.ac.uk).

Dusit Niyato is with the School of Computer Science and Engineering, Nanyang Technological University, Singapore 639798 (E-mail: dniyato@ntu.edu.sg).

TABLE I
COMPLEX MATRIX DERIVATIVES UNDER DIAGONAL STRUCTURE CONSTRAINTS

| Function Type | Derivative w.r.t. $\Lambda_\Theta$ | Derivative w.r.t. $\Lambda_\Theta^*$ |
|---|---|---|
| $f_{\mathrm{D,TL}} = \mathrm{Tr}(\Lambda_\Theta^{\mathrm{H}} M) + \mathrm{Tr}(\Lambda_\Theta M^{\mathrm{H}})$ | $\mathrm{Diag}\{M^{\mathrm{H}}\}$ | $\mathrm{Diag}\{M\}$ |
| $f_{\mathrm{D,TQ}} = \mathrm{Tr}(\Lambda_\Theta^{\mathrm{H}} W \Lambda_\Theta)$ | $\mathrm{Diag}\{\Lambda_\Theta^{\mathrm{H}} W\}$ | $\mathrm{Diag}\{W \Lambda_\Theta\}$ |
| $f_{\mathrm{D,TI}} = \mathrm{Tr}\left((I_M + \Phi\Lambda_\Theta)^{-1}\right)$ | $-\mathrm{Diag}\left\{\Phi^{\frac{1}{2}}(I_M + \Phi^{\frac{1}{2}}\Lambda_\Theta\Phi^{\frac{1}{2}})^{-2}\Phi^{\frac{1}{2}}\right\}$ | $0$ |
| $f_{\mathrm{D,LD}} = \log|I_M + \Phi\Lambda_\Theta|$ | $\mathrm{Diag}\left\{\Phi^{\frac{1}{2}}(I_M + \Phi^{\frac{1}{2}}\Lambda_\Theta\Phi^{\frac{1}{2}})^{-1}\Phi^{\frac{1}{2}}\right\}$ | $0$ |

for diagonal structure constraints, this algorithm considers applying the first-order derivative to each diagonal element to obtain the optimal solution. Furthermore, in terms of the intractable constant modulus constraints, dual variables are usually introduced and iteratively optimized by the subgradient method to satisfy complementary slackness conditions, thereby potentially suffering from high iteration complexity [23]. In contrast, the BCD algorithm is always adopted to solve highly nonconvex problems caused by strongly-coupled matrix variables. Specifically, under the BCD framework, the original matrix-variate optimization problem can be decomposed into multiple low-dimensional subproblems, each of which needs to be iteratively optimized until convergence. In order to ensure at least local convergence of the BCD algorithm [20], each subproblem is required to have a unique optimal solution. Nevertheless, considering that subproblems may be nonconvex and thus hard to globally solve, the MM-based algorithm has attained extensive attention, whose core idea is to construct a tractable surrogate function to locally approximate the original nonconvex subproblem. Unfortunately, the derivation of the surrogate function usually involves high-complexity matrix manipulations and also needs to be iteratively carried out to achieve a close approximation.

Obviously, the KKT-based algorithm based on complex matrix derivatives generally achieves the lowest complexity among the three types of algorithms [24]. Nonetheless, its application range is relatively limited as compared to the BCD and MM-based algorithms. Note that the implementation of the latter two algorithms depends on the specific wireless system and may have high computational complexity, especially for large-scale arrays. To circumvent these issues, in this paper, we aim to develop a unified framework for matrix-variate optimization with two special structure constraints, namely, diagonal structure and constant modulus constraints. For each considered case, the novel low-complexity algorithm with guaranteed performance is proposed. The main contributions of our work are further summarized as follows.

- Firstly, we consider the diagonal structure constraints often seen in the uplink multi-user single-input multiple-output (MU-SIMO) system and the amplitude-adjustable IRS-aided MIMO system, which are always involved in the capacity maximization problem, mean squared error (MSE) minimization problem and their variants. We propose complex matrix derivatives associated with diagonal structures, based on which the optimal solutions of these matrix-variate problems are directly obtained in closed forms. Furthermore, the above study is extended to the case of block-diagonal structure constraints.
- Secondly, in terms of constant modulus constraints mostly adopted in the hybrid analog-digital MIMO system and the fully-passive IRS-aided MIMO system, we propose the element-wise phase derivatives inspired by their

element-wise decomposability nature. For different classical matrix-variate optimization problems, it is revealed that the element-wise phase derivatives can be classified into the following two general forms, i.e., the linear form and the conjugate linear form.
- Finally, by exploring inherent structures of the element-wise phase derivatives, we develop a novel alternating optimization (AO) algorithm with the aid of several arbitrary feasible solutions for the matrix-variate optimization under constant modulus constraints. Note that the computational complexity of the proposed AO algorithm sharply decreases, since it avoids the complicated matrix inversion and matrix factorization involved in the conventional element-wise iterative algorithm. Moreover, we demonstrate that the proposed algorithm is able to achieve almost the same performance as the existing benchmark schemes.

*Notation:* Scalars, vectors and matrices are represented by non-bold, bold lowercase, and bold uppercase letters, respectively. The notations $A^{\mathrm{T}}, A^*, A^{\mathrm{H}}, A^{-1}, \mathrm{Tr}(A)$ and $|A|$ denote the transpose, conjugate, hermitian, inversion, trace and determinant of the complex matrix $A$, respectively. $\mathrm{Diag}\{A\}$ denotes a vector whose elements are diagonal elements of matrix $A$, and $\mathrm{Blockdiag}(\{A_k\}_{k=1}^K)$ is a block diagonal matrix with diagonal sub-matrices of $A_k$'s. Moreover, the $i$th row and the $j$th column of $A$ are denoted as $[A]_{i,:}$ and $[A]_{:,j}$, respectively, the element in the $i$th row and the $j$th column is denoted as $[A]_{i,j}$. $\frac{\mathrm{d}f}{\mathrm{d}a}$ and $\frac{\partial f}{\partial a}$ denote the differential and the partial derivative of $f$ with respect to $a$, respectively. $\odot$ denotes the Hadamard product and $(a)^+ = \max\{0, a\}$. $\Re\{a\}$ and $\Im\{a\}$ denote the real and imaginary parts of a complex variable $a$, respectively. The symbol $\mathrm{Phase}\{a\}$ denotes the phase of a complex scalar $a$ and $-\pi < \mathrm{Phase}\{a\} \leq \pi$. Lastly, the word "with respect to" is abbreviated as "w.r.t.".

## II. DIAGONAL STRUCTURE CONSTRAINTS

In this section, we firstly provide some fundamental properties of complex matrix derivatives associated with diagonal structures for several types of objective functions. Based on these properties, we then obtain the optimal solutions of a series of optimization problems in the uplink MU-SIMO system and the amplitude-adjustable IRS-aided MIMO system in closed forms. Moreover, the above study is extended to the case of block-diagonal matrix variables.

### A. Mathematical Preliminaries

At the beginning, some fundamental definitions for diagonal matrices are provided, which are the basis of the following analysis. For any two diagonal matrices $\Lambda_{\Theta_1} \in \mathbb{C}^{M \times M}$ and $\Lambda_{\Theta_2} \in \mathbb{C}^{M \times M}$, we have

$$\mathrm{Diag}\{\Lambda_{\Theta_1} + \Lambda_{\Theta_2}\} = \mathrm{Diag}\{\Lambda_{\Theta_1}\} + \mathrm{Diag}\{\Lambda_{\Theta_2}\}. \quad (1)$$
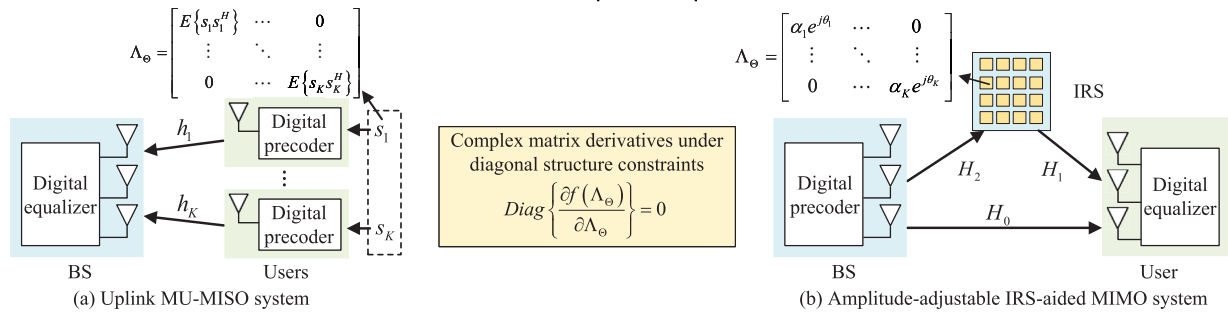
Fig. 1. A diagram of application scenarios associated with diagonal matrix variables.

For an arbitrary square matrix $\boldsymbol{\Phi} \in \mathbb{C}^{M \times M}$, we have

$$\text{Diag}\{\boldsymbol{\Phi}\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}\} = \text{Diag}\{\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}\boldsymbol{\Phi}\}. \tag{2}$$

Moreover, the following equality holds for any complex matrices $\boldsymbol{N} \in \mathbb{C}^{N \times M}$ and $\boldsymbol{M} \in \mathbb{C}^{M \times N}$, i.e.,

$$\text{Diag}\{\boldsymbol{N}\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}\boldsymbol{M}\} = \boldsymbol{\Phi}\text{Diag}\{\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}\} \text{ with } \boldsymbol{\Phi} = \boldsymbol{N} \odot \boldsymbol{M}^{\text{T}}. \tag{3}$$

Together with the operation of $\text{Diag}\{\cdot\}$, the vectorization of the first-order derivative of a scalar-valued function $f(\boldsymbol{\Lambda}_{\boldsymbol{\Theta}})$ w.r.t the complex diagonal matrix $\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}$ can be defined as [24]

$$\text{Diag}\left\{\frac{\partial f(\boldsymbol{\Lambda}_{\boldsymbol{\Theta}})}{\partial \boldsymbol{\Lambda}_{\boldsymbol{\Theta}}}\right\} = \left[\frac{\partial f(\boldsymbol{\Lambda}_{\boldsymbol{\Theta}})}{\partial [\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}]_{1,1}}, \cdots, \frac{f(\boldsymbol{\Lambda}_{\boldsymbol{\Theta}})}{\partial [\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}]_{N,N}}\right]^{\text{T}},$$

$$\text{Diag}\left\{\frac{\partial f(\boldsymbol{\Lambda}_{\boldsymbol{\Theta}})}{\partial \boldsymbol{\Lambda}_{\boldsymbol{\Theta}}^{*}}\right\} = \left[\frac{\partial f(\boldsymbol{\Lambda}_{\boldsymbol{\Theta}})}{\partial [\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}]_{1,1}^{*}}, \cdots, \frac{f(\boldsymbol{\Lambda}_{\boldsymbol{\Theta}})}{\partial [\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}]_{N,N}^{*}}\right]^{\text{T}}. \tag{4}$$

In the sequel, we mainly concern about complex matrix derivatives w.r.t. diagonal matrices for four common objective functions, including the trace-linear function, the trace-quadratic function, the trace-inverse function and the log-determinant function.

*1) Trace-Linear Function:* For an arbitrary complex matrix $\boldsymbol{M} \in \mathbb{C}^{M \times M}$, the differential of a trace-linear function $f_{\text{D,TL}} = \text{Tr}(\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}^{\text{H}}\boldsymbol{M}) + \text{Tr}(\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}\boldsymbol{M}^{\text{H}})$ w.r.t. $\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}$ can be obtained as $\text{d}(f_{\text{D,TL}}) = \text{Tr}(\text{d}(\boldsymbol{\Lambda}_{\boldsymbol{\Theta}})\boldsymbol{M}^{\text{H}})$. Based on this, the corresponding first-order derivative w.r.t. $\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}$ can be obtained as

$$\text{Diag}\left\{\frac{\partial f_{\text{D,TL}}}{\partial \boldsymbol{\Lambda}_{\boldsymbol{\Theta}}}\right\} = \text{Diag}\{\boldsymbol{M}^{\text{H}}\}. \tag{5}$$

The first-order derivative of $f_{\text{D,TL}}$ w.r.t. $\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}^{*}$ is also given by

$$\text{Diag}\left\{\frac{\partial f_{\text{D,TL}}}{\partial \boldsymbol{\Lambda}_{\boldsymbol{\Theta}}^{*}}\right\} = \text{Diag}\{\boldsymbol{M}\}. \tag{6}$$

*2) Trace-Quadratic Function:* For a Hermitian matrix $\boldsymbol{W} \in \mathbb{C}^{M \times M}$, the differential of a trace-quadratic function $f_{\text{D,TQ}} = \text{Tr}(\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}^{\text{H}}\boldsymbol{W}\boldsymbol{\Lambda}_{\boldsymbol{\Theta}})$ w.r.t. $\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}$ and $\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}^{*}$ are respectively calculated as $\text{d}(f_{\text{D,TQ}}) = \text{Tr}(\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}^{\text{H}}\boldsymbol{W}\text{d}(\boldsymbol{\Lambda}_{\boldsymbol{\Theta}})), \text{d}(f_{\text{D,TQ}}) = \text{Tr}(\text{d}(\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}^{\text{H}})\boldsymbol{W}\boldsymbol{\Lambda}_{\boldsymbol{\Theta}})$. Then, the corresponding first-order derivatives w.r.t. $\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}$ and $\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}^{*}$ can be obtained as

$$\text{Diag}\left\{\frac{\partial f_{\text{D,TQ}}}{\partial \boldsymbol{\Lambda}_{\boldsymbol{\Theta}}}\right\} = \text{Diag}\{\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}^{\text{H}}\boldsymbol{W}\},$$

$$\text{Diag}\left\{\frac{\partial f_{\text{D,TQ}}}{\partial \boldsymbol{\Lambda}_{\boldsymbol{\Theta}}^{*}}\right\} = \text{Diag}\{\boldsymbol{W}\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}\}. \tag{7}$$

*3) Trace-Inverse Function:* For a positive semi-definite matrix $\boldsymbol{\Phi} \in \mathbb{C}^{M \times M}$, the differential of a trace-inverse function $f_{\text{D,TI}} = \text{Tr}\left((\boldsymbol{I}_M + \boldsymbol{\Phi}\boldsymbol{\Lambda}_{\boldsymbol{\Theta}})^{-1}\right)$ w.r.t. $\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}$ is given by

$\text{d}(f_{\text{D,TI}}) = -\text{Tr}\left((\boldsymbol{I}_M + \boldsymbol{\Phi}\boldsymbol{\Lambda}_{\boldsymbol{\Theta}})^{-2}\boldsymbol{\Phi}\text{d}\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}\right)$. The corresponding first-order derivative w.r.t. $\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}$ is then given by

$$\text{Diag}\left\{\frac{\partial f_{\text{D,TI}}}{\partial \boldsymbol{\Lambda}_{\boldsymbol{\Theta}}}\right\} = -\text{Diag}\left\{(\boldsymbol{I}_M + \boldsymbol{\Phi}\boldsymbol{\Lambda}_{\boldsymbol{\Theta}})^{-2}\boldsymbol{\Phi}\right\} \tag{8}$$

$$= -\text{Diag}\left\{\boldsymbol{\Phi}^{\frac{1}{2}}(\boldsymbol{I}_M + \boldsymbol{\Phi}^{\frac{1}{2}}\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}\boldsymbol{\Phi}^{\frac{1}{2}})^{-2}\boldsymbol{\Phi}^{\frac{1}{2}}\right\}.$$

*4) Log-Determinant Function:* Considering a log-determinant function $f_{\text{D,LD}} = \log|\boldsymbol{I}_M + \boldsymbol{\Phi}\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}|$, its differential w.r.t. $\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}$ can be derived as $\text{d}(f_{\text{D,LD}}) = \text{Tr}\left((\boldsymbol{I}_M + \boldsymbol{\Phi}\boldsymbol{\Lambda}_{\boldsymbol{\Theta}})^{-1}\boldsymbol{\Phi}\text{d}\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}\right)$, based on which the following first-order derivative w.r.t. $\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}$ holds.

$$\text{Diag}\left\{\frac{\partial f_{\text{D,LD}}}{\partial \boldsymbol{\Lambda}_{\boldsymbol{\Theta}}}\right\} = \text{Diag}\left\{(\boldsymbol{I}_M + \boldsymbol{\Phi}\boldsymbol{\Lambda}_{\boldsymbol{\Theta}})^{-1}\boldsymbol{\Phi}\right\} \tag{9}$$

$$= \text{Diag}\left\{\boldsymbol{\Phi}^{\frac{1}{2}}(\boldsymbol{I}_M + \boldsymbol{\Phi}^{\frac{1}{2}}\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}\boldsymbol{\Phi}^{\frac{1}{2}})^{-1}\boldsymbol{\Phi}^{\frac{1}{2}}\right\}.$$

In conclusion, the complex matrix derivatives of the above four types of objective functions are summarized in Table I. Exploiting these fundamental properties, some classical wireless applications are investigated in the following subsection.

*B. Specific Wireless Applications*

*1) Uplink MU-SIMO System:* As shown in Fig. 1(a), we firstly consider an uplink distributed MU-SIMO system, where $K$ single-antenna users transmit independent data streams to the BS equipped with $N_t$ antennas [10]. The received signal at the BS can be expressed as

$$\boldsymbol{y} = \boldsymbol{H}\boldsymbol{s} + \boldsymbol{n} \text{ with } \boldsymbol{H} = [\boldsymbol{h}_1, \cdots, \boldsymbol{h}_K] \in \mathbb{C}^{N_t \times K}, \tag{10}$$

where $\boldsymbol{h}_k \in \mathbb{C}^{N_t \times 1}$ denotes the channel between the BS and the $k$th user, and $\boldsymbol{s} = [s_1, \cdots, s_K]^{\text{T}} \in \mathbb{C}^{K \times 1}$ is the transmit signal, $\boldsymbol{n}$ is the additive noise obeying Gaussian distribution with zero mean and covariance matrix $\mathbb{E}\{\boldsymbol{n}\boldsymbol{n}^{\text{H}}\} = \boldsymbol{\Sigma}$. It is worth noting that the covariance matrix of $\boldsymbol{s}$ is diagonal, since the transmit data streams of $K$ users are independent of each other, i.e., $\mathbb{E}\{\boldsymbol{s}\boldsymbol{s}^{\text{H}}\} = \boldsymbol{\Lambda}_{\boldsymbol{\Theta}}$. Then, the uplink capacity maximization problem is formulated as

**Prob.1:** $\max\limits_{\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}} \log\left|\boldsymbol{I}_{N_t} + \boldsymbol{\Sigma}^{-1}\boldsymbol{H}\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}\boldsymbol{H}^{\text{H}}\right|$

$$\text{s.t. } \text{Tr}(\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}) \leq P, \ 0 \leq [\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}]_{k,k} \leq P_k, \ \forall k. \tag{11}$$

By recalling (9), the first-order derivative of the objective function of **Prob.1** w.r.t. $\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}$ is given by

$$\text{Diag}\left\{\frac{\partial \log|\boldsymbol{I}_{N_t} + \boldsymbol{\Sigma}^{-1}\boldsymbol{H}\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}\boldsymbol{H}^{\text{H}}|}{\partial \boldsymbol{\Lambda}_{\boldsymbol{\Theta}}}\right\} \tag{12}$$

$$= \text{Diag}\left\{\boldsymbol{H}^{\text{H}}\boldsymbol{\Sigma}^{-\frac{1}{2}}\left(\boldsymbol{I}_{N_t} + \boldsymbol{\Sigma}^{-\frac{1}{2}}\boldsymbol{H}\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}\boldsymbol{H}^{\text{H}}\boldsymbol{\Sigma}^{-\frac{1}{2}}\right)^{-1}\boldsymbol{\Sigma}^{-\frac{1}{2}}\boldsymbol{H}\right\}.$$

In order to derive the optimal $\boldsymbol{\Lambda_\Theta}$, we present the KKT optimality conditions of **Prob.1** as follows [25]:

$$\begin{cases} \mathrm{Diag}\Big\{\boldsymbol{H}^{\mathrm{H}}\boldsymbol{\Sigma}^{-\frac{1}{2}}\Big(\boldsymbol{I}_{N_t}+\boldsymbol{\Sigma}^{-\frac{1}{2}}\boldsymbol{H}\boldsymbol{\Lambda_\Theta}\boldsymbol{H}^{\mathrm{H}}\boldsymbol{\Sigma}^{-\frac{1}{2}}\Big)^{-1} \\ \qquad\qquad \times\boldsymbol{\Sigma}^{-\frac{1}{2}}\boldsymbol{H}\Big\}=\mu\mathrm{Diag}\{\boldsymbol{I}_K\}-[\psi_1,\cdots,\psi_K]^{\mathrm{T}}, & (13\mathrm{a}) \\ \mu\left(\mathrm{Tr}(\boldsymbol{\Lambda_\Theta})-P\right)=0, & (13\mathrm{b}) \\ \psi_k\left([\boldsymbol{\Lambda_\Theta}]_{k,k}-P_k\right)=0,\ \forall k, & (13\mathrm{c}) \end{cases}$$

where $\mu$ is the Lagrange multiplier associated with the sum power constraint and $\psi_k$'s correspond to the power constraints imposed on each user. We firstly define $\boldsymbol{J}_k=\boldsymbol{I}_K+\left(\boldsymbol{H}^{\mathrm{H}}\boldsymbol{\Sigma}^{-1}\boldsymbol{H}-\overline{\boldsymbol{H}}_k\right)\left(\boldsymbol{\Lambda_\Theta}-\overline{\boldsymbol{\Lambda_\Theta}}_k\right)$, where $\overline{\boldsymbol{H}}_k\in\mathbb{C}^{K\times K}$ is an all-zero matrix except for its $k$th column being $\left[\boldsymbol{H}^{\mathrm{H}}\boldsymbol{\Sigma}^{-1}\boldsymbol{H}\right]_{:,k}$ and $\overline{\boldsymbol{\Lambda_\Theta}}_k\in\mathbb{C}^{K\times K}$ is an all-zero matrix except for its $(k,k)$th element being $[\boldsymbol{\Lambda_\Theta}]_{k,k}$. Then, the left-hand side of (13a) can be rewritten as

$$\mathrm{Diag}\Big\{\boldsymbol{H}^{\mathrm{H}}\boldsymbol{\Sigma}^{-\frac{1}{2}}\Big(\boldsymbol{I}_{N_t}+\boldsymbol{\Sigma}^{-\frac{1}{2}}\boldsymbol{H}\boldsymbol{\Lambda_\Theta}\boldsymbol{H}^{\mathrm{H}}\boldsymbol{\Sigma}^{-\frac{1}{2}}\Big)^{-1}\boldsymbol{\Sigma}^{-\frac{1}{2}}\boldsymbol{H}\Big\}$$
$$\overset{(a_1)}{=}\mathrm{Diag}\Big\{\left(\boldsymbol{I}_K+\boldsymbol{H}^{\mathrm{H}}\boldsymbol{\Sigma}^{-1}\boldsymbol{H}\boldsymbol{\Lambda_\Theta}\right)^{-1}\boldsymbol{H}^{\mathrm{H}}\boldsymbol{\Sigma}^{-1}\boldsymbol{H}\Big\}$$
$$\overset{(a_2)}{=}-\frac{[\boldsymbol{\Lambda_\Theta}]_{k,k}\mathrm{Diag}\left\{\boldsymbol{J}_k^{-1}\overline{\boldsymbol{H}}_k\boldsymbol{J}_k^{-1}\boldsymbol{H}^{\mathrm{H}}\boldsymbol{\Sigma}^{-1}\boldsymbol{H}\right\}}{1+[\boldsymbol{\Lambda_\Theta}]_{k,k}\mathrm{Tr}(\boldsymbol{J}_k^{-1}\overline{\boldsymbol{H}}_k)}$$
$$+\mathrm{Diag}\left\{\boldsymbol{J}_k^{-1}\boldsymbol{H}^{\mathrm{H}}\boldsymbol{\Sigma}^{-1}\boldsymbol{H}\right\}, \qquad (14)$$

where $(a_1)$ holds based on the matrix inversion lemma $\boldsymbol{N}(\boldsymbol{I}+\boldsymbol{M}\boldsymbol{N})^{-1}=(\boldsymbol{I}+\boldsymbol{N}\boldsymbol{M})^{-1}\boldsymbol{N}$ and $(a_2)$ is attained using the Sherman Morrison formula, i.e., $(\boldsymbol{M}+\boldsymbol{N})^{-1}=\boldsymbol{M}^{-1}-\frac{\boldsymbol{M}^{-1}\boldsymbol{N}\boldsymbol{M}^{-1}}{1+\mathrm{Tr}(\boldsymbol{M}^{-1}\boldsymbol{N})}$ for a full-rank matrix $\boldsymbol{M}$ and a rank-one matrix $\boldsymbol{N}$. By substituting (14) into (13a), the optimal $[\boldsymbol{\Lambda_\Theta}]_{k,k}$'s is derived in the following water-filling form [26].

$$[\boldsymbol{\Lambda_\Theta}]_{k,k}=\begin{cases} \left(\frac{y_k}{\left[\mathrm{Diag}\left\{\boldsymbol{J}_k^{-1}\overline{\boldsymbol{H}}_k\boldsymbol{J}_k^{-1}\boldsymbol{H}^{\mathrm{H}}\boldsymbol{\Sigma}^{-1}\boldsymbol{H}\right\}\right]_k+\mathrm{Tr}(\boldsymbol{J}_k^{-1}\overline{\boldsymbol{H}}_k)y_k}\right)^+, \\ \qquad\qquad \text{if }[\boldsymbol{\Lambda_\Theta}]_{k,k}\leq P_k, \\ P_k, \qquad \text{if }[\boldsymbol{\Lambda_\Theta}]_{k,k}>P_k,\ \forall k, \end{cases} \quad (15)$$

where $y_k=\left[\mathrm{Diag}\left\{\boldsymbol{J}_k^{-1}\boldsymbol{H}^{\mathrm{H}}\boldsymbol{\Sigma}^{-1}\boldsymbol{H}\right\}\right]_k-\mu$. Moreover, since $\mathrm{Tr}(\boldsymbol{\Lambda_\Theta})$ is monotonically decreasing w.r.t $\mu>0$, the optimal $\mu$ satisfying (13b) can be found via the bisection search.

In addition, MSE is a widely used performance metric, which reflects the accuracy of the desired signals that can be recovered from the noise corrupted observations. Accordingly, the MSE minimization problem is formulated as

**Prob.2:** $\displaystyle\min_{\boldsymbol{\Lambda_\Theta}}\ \mathrm{Tr}\left[\left(\boldsymbol{I}_{N_t}+\boldsymbol{\Sigma}^{-1}\boldsymbol{H}\boldsymbol{\Lambda_\Theta}\boldsymbol{H}^{\mathrm{H}}\right)^{-1}\right]$

$\qquad\qquad$ s.t. $\mathrm{Tr}(\boldsymbol{\Lambda_\Theta})\leq P,\ 0\leq[\boldsymbol{\Lambda_\Theta}]_{k,k}\leq P_k,\ \forall k.$ (16)

By recalling (8), the first-order derivative of the objective function of **Prob.2** w.r.t. $\boldsymbol{\Lambda_\Theta}$ is given by

$$\mathrm{Diag}\left\{\frac{\partial\mathrm{Tr}\left[\left(\boldsymbol{I}_{N_t}+\boldsymbol{\Sigma}^{-1}\boldsymbol{H}\boldsymbol{\Lambda_\Theta}\boldsymbol{H}^{\mathrm{H}}\right)^{-1}\right]}{\partial\boldsymbol{\Lambda_\Theta}}\right\} \qquad (17)$$
$$=-\mathrm{Diag}\left\{\boldsymbol{H}^{\mathrm{H}}\boldsymbol{\Sigma}^{-\frac{1}{2}}\Big(\boldsymbol{I}_{N_t}+\boldsymbol{\Sigma}^{-\frac{1}{2}}\boldsymbol{H}\boldsymbol{\Lambda_\Theta}\boldsymbol{H}^{\mathrm{H}}\boldsymbol{\Sigma}^{-\frac{1}{2}}\Big)^{-2}\boldsymbol{\Sigma}^{-\frac{1}{2}}\boldsymbol{H}\right\}.$$

Based on (17), the KKT optimality conditions of **Prob.2** can be formulated as

$$\begin{cases} \mathrm{Diag}\Big\{\boldsymbol{H}^{\mathrm{H}}\boldsymbol{\Sigma}^{-\frac{1}{2}}\Big(\boldsymbol{I}_{N_t}+\boldsymbol{\Sigma}^{-\frac{1}{2}}\boldsymbol{H}\boldsymbol{\Lambda_\Theta}\boldsymbol{H}^{\mathrm{H}}\boldsymbol{\Sigma}^{-\frac{1}{2}}\Big)^{-2} \\ \qquad\qquad \times\boldsymbol{\Sigma}^{-\frac{1}{2}}\boldsymbol{H}\Big\}=\mu\mathrm{Diag}\{\boldsymbol{I}_K\}-[\psi_1,\cdots,\psi_K]^{\mathrm{T}}, & (18\mathrm{a}) \\ \mu\left(\mathrm{Tr}(\boldsymbol{\Lambda_\Theta})-P\right)=0, & (18\mathrm{b}) \\ \psi_k\left([\boldsymbol{\Lambda_\Theta}]_{k,k}-P_k\right)=0,\ \forall k, & (18\mathrm{c}) \end{cases}$$

where $\mu$ and $\psi_k$'s are defined similarly to **Prob.1**. The left-hand side of (18a) can be rewritten as

$$\mathrm{Diag}\Big\{\boldsymbol{H}^{\mathrm{H}}\boldsymbol{\Sigma}^{-\frac{1}{2}}\Big(\boldsymbol{I}_{N_t}+\boldsymbol{\Sigma}^{-\frac{1}{2}}\boldsymbol{H}\boldsymbol{\Lambda_\Theta}\boldsymbol{H}^{\mathrm{H}}\boldsymbol{\Sigma}^{-\frac{1}{2}}\Big)^{-2}\boldsymbol{\Sigma}^{-\frac{1}{2}}\boldsymbol{H}\Big\}$$
$$\overset{(b)}{=}\mathrm{Diag}\left\{\left(\boldsymbol{J}_k^{-1}-\frac{[\boldsymbol{\Lambda_\Theta}]_{k,k}\boldsymbol{J}_k^{-1}\overline{\boldsymbol{H}}_k\boldsymbol{J}_k^{-1}}{1+[\boldsymbol{\Lambda_\Theta}]_{k,k}\mathrm{Tr}(\boldsymbol{J}_k^{-1}\overline{\boldsymbol{H}}_k)}\right)^2\boldsymbol{H}^{\mathrm{H}}\boldsymbol{\Sigma}^{-1}\boldsymbol{H}\right\}$$
$$=\frac{[\boldsymbol{\Lambda_\Theta}]_{k,k}^2\mathrm{Diag}\left\{\boldsymbol{J}_k^{-1}\overline{\boldsymbol{H}}_k\boldsymbol{J}_k^{-2}\overline{\boldsymbol{H}}_k\boldsymbol{J}_k^{-1}\boldsymbol{H}^{\mathrm{H}}\boldsymbol{\Sigma}^{-1}\boldsymbol{H}\right\}}{\left(1+[\boldsymbol{\Lambda_\Theta}]_{k,k}\mathrm{Tr}(\boldsymbol{J}_k^{-1}\overline{\boldsymbol{H}}_k)\right)^2}$$
$$-\frac{[\boldsymbol{\Lambda_\Theta}]_{k,k}\mathrm{Diag}\left\{\boldsymbol{J}_k^{-2}\overline{\boldsymbol{H}}_k\boldsymbol{J}_k^{-1}\boldsymbol{H}^{\mathrm{H}}\boldsymbol{\Sigma}^{-1}\boldsymbol{H}\right\}}{1+[\boldsymbol{\Lambda_\Theta}]_{k,k}\mathrm{Tr}(\boldsymbol{J}_k^{-1}\overline{\boldsymbol{H}}_k)}$$
$$-\frac{[\boldsymbol{\Lambda_\Theta}]_{k,k}\mathrm{Diag}\left\{\boldsymbol{J}_k^{-1}\overline{\boldsymbol{H}}_k\boldsymbol{J}_k^{-2}\boldsymbol{H}^{\mathrm{H}}\boldsymbol{\Sigma}^{-1}\boldsymbol{H}\right\}}{1+[\boldsymbol{\Lambda_\Theta}]_{k,k}\mathrm{Tr}(\boldsymbol{J}_k^{-1}\overline{\boldsymbol{H}}_k)}$$
$$+\mathrm{Diag}\left\{\boldsymbol{J}_k^{-2}\boldsymbol{H}^{\mathrm{H}}\boldsymbol{\Sigma}^{-1}\boldsymbol{H}\right\}, \qquad (19)$$

where $(b)$ holds similarly to (14). Substituting (19) into (18a), we have

$$\widetilde{a}_k[\boldsymbol{\Lambda_\Theta}]_{k,k}^2+\widetilde{b}_k[\boldsymbol{\Lambda_\Theta}]_{k,k}+x_k=0,\ \forall k, \qquad (20)$$

where

$$\widetilde{a}_k=\mathrm{Tr}^2(\boldsymbol{J}_k^{-1}\overline{\boldsymbol{H}}_k)x_k+\left[\mathrm{Diag}\left\{\boldsymbol{J}_k^{-1}\overline{\boldsymbol{H}}_k\boldsymbol{J}_k^{-2}\overline{\boldsymbol{H}}_k\boldsymbol{J}_k^{-1}\boldsymbol{H}^{\mathrm{H}}\boldsymbol{\Sigma}^{-1}\right.\right.$$
$$\left.\left.\boldsymbol{H}\right\}\right]_k-\mathrm{Tr}(\boldsymbol{J}_k^{-1}\overline{\boldsymbol{H}}_k)[\mathrm{Diag}\left\{\boldsymbol{J}_k^{-2}\overline{\boldsymbol{H}}_k\boldsymbol{J}_k^{-1}\boldsymbol{H}^{\mathrm{H}}\boldsymbol{\Sigma}^{-1}\boldsymbol{H}\right\}]_k$$
$$-\mathrm{Tr}(\boldsymbol{J}_k^{-1}\overline{\boldsymbol{H}}_k)[\mathrm{Diag}\left\{\boldsymbol{J}_k^{-1}\overline{\boldsymbol{H}}_k\boldsymbol{J}_k^{-2}\boldsymbol{H}^{\mathrm{H}}\boldsymbol{\Sigma}^{-1}\boldsymbol{H}\right\}]_k,$$
$$\widetilde{b}_k=2\mathrm{Tr}(\boldsymbol{J}_k^{-1}\overline{\boldsymbol{H}}_k)x_k-[\mathrm{Diag}\left\{\boldsymbol{J}_k^{-2}\overline{\boldsymbol{H}}_k\boldsymbol{J}_k^{-1}\boldsymbol{H}^{\mathrm{H}}\boldsymbol{\Sigma}^{-1}\boldsymbol{H}\right\}]_k$$
$$-[\mathrm{Diag}\left\{\boldsymbol{J}_k^{-1}\overline{\boldsymbol{H}}_k\boldsymbol{J}_k^{-2}\boldsymbol{H}^{\mathrm{H}}\boldsymbol{\Sigma}^{-1}\boldsymbol{H}\right\}]_k,$$
$$x_k=\left[\mathrm{Diag}\left\{\boldsymbol{J}_k^{-2}\boldsymbol{H}^{\mathrm{H}}\boldsymbol{\Sigma}^{-1}\boldsymbol{H}\right\}\right]_k-\mu,\ \forall k. \qquad (21)$$

Then, the optimal $[\boldsymbol{\Lambda_\Theta}]_{k,k}$'s of **Prob.2** can be attained using quadratic formula as follows:

$$[\boldsymbol{\Lambda_\Theta}]_{k,k}=\begin{cases} \left(\frac{-\widetilde{b}\pm\sqrt{\widetilde{b}^2-4[\boldsymbol{x}_k]_k\widetilde{a}}}{2\widetilde{a}}\right)^+, & \text{if }[\boldsymbol{\Lambda_\Theta}]_{k,k}\leq P_k, \\ P_k, & \text{if }[\boldsymbol{\Lambda_\Theta}]_{k,k}>P_k, \end{cases}\ \forall k. \quad (22)$$

Similarly, $\mu$ can be obtained using the bisection search.

*2) Amplitude-Adjustable IRS-aided MIMO System:* Hereafter, we consider the state-of-the-art amplitude-adjustable IRS-aided point-to-point MIMO system as shown in Fig. 1(b), where a $N_t$-antenna BS serves a $N_r$-antenna user with the aid of a $K$-element amplitude-adjustable IRS. The received signal at the user can be expressed as

$$\boldsymbol{y}=\boldsymbol{H}\boldsymbol{s}+\boldsymbol{n}\quad\text{with}\quad\boldsymbol{H}=\boldsymbol{H}_0+\boldsymbol{H}_1\boldsymbol{\Lambda_\Theta}\boldsymbol{H}_2, \quad (23)$$

where $\boldsymbol{H}_0\in\mathbb{C}^{N_r\times N_t}$, $\boldsymbol{H}_1\in\mathbb{C}^{N_r\times K}$ and $\boldsymbol{H}_2\in\mathbb{C}^{K\times N_t}$ represent the BS-user direct channel, the IRS-user channel and the BS-IRS channel, respectively. $\boldsymbol{\Lambda_\Theta}\in\mathbb{C}^{K\times K}$ denotes the diagonal

IRS reflection matrix whose each diagonal element represents the adjustable amplitude and phase of the corresponding reflecting element, which usually satisfies $\left|[\mathbf{\Lambda}_\Theta]_{i,i}\right| \leq 1$ or $\mathrm{Tr}(\mathbf{\Lambda}_\Theta \mathbf{\Lambda}_\Theta^{\mathrm{H}}) \leq K$ [27]. Similar to Sec. II-B1, we firstly consider the following capacity maximization problem.

**Prob.3:** $\max_{\mathbf{\Lambda}_\Theta} \ \log \left|\mathbf{\Sigma}^{-1}\mathbf{H}\mathbf{H}^{\mathrm{H}} + \mathbf{I}_{N_r}\right|$

$$\text{s.t. } \mathrm{Tr}(\mathbf{\Lambda}_\Theta \mathbf{\Lambda}_\Theta^{\mathrm{H}}) \leq K. \tag{24}$$

By recalling (7) and (9), the first-order derivative of the objective function of **Prob.3** w.r.t. $\mathbf{\Lambda}_\Theta$ is given by

$$\mathrm{Diag}\left\{\frac{\partial \log\left|\mathbf{\Sigma}^{-1}\mathbf{H}\mathbf{H}^{\mathrm{H}} + \mathbf{I}_{N_r}\right|}{\partial \mathbf{\Lambda}_\Theta}\right\}$$

$$=\mathrm{Diag}\left\{\mathbf{H}_2\mathbf{H}_2^{\mathrm{H}}\mathbf{\Lambda}_\Theta^{\mathrm{H}}\mathbf{H}_1^{\mathrm{H}}\left[\mathbf{\Sigma}^{-1}\mathbf{H}\mathbf{H}^{\mathrm{H}} + \mathbf{I}_{N_r}\right]^{-1}\mathbf{\Sigma}^{-1}\mathbf{H}_1\right\}$$

$$+ \mathrm{Diag}\left\{\mathbf{H}_2\mathbf{H}_0^{\mathrm{H}}\left[\mathbf{\Sigma}^{-1}\mathbf{H}\mathbf{H}^{\mathrm{H}} + \mathbf{I}_{N_r}\right]^{-1}\mathbf{\Sigma}^{-1}\mathbf{H}_1\right\}. \tag{25}$$

Unfortunately, even though the first-order derivative is derived, it is still difficult to derive the optimal closed-form solution from (25), since its involved quadratic term w.r.t. $\mathbf{\Lambda}_\Theta$ appears in an inverse form. As a remedy, we intend to solve it based on problem transformation. Specifically, via introducing a series of auxiliary variables, **Prob.3** is equivalently transformed into

**Prob.4:** $\min_{\mathbf{G}_\Lambda, \mathbf{\Lambda}_\Theta, \mathbf{W}} \ \mathrm{Tr}\left(\mathbf{W}\left[\mathbf{G}_\Lambda \mathbf{H} - \mathbf{I}_{N_t}\right]\left[\mathbf{G}_\Lambda \mathbf{H} - \mathbf{I}_{N_t}\right]^{\mathrm{H}}\right)$

$$+ \mathrm{Tr}(\mathbf{W}\mathbf{G}_\Lambda\mathbf{\Sigma}\mathbf{G}_\Lambda^{\mathrm{H}}) - \log|\mathbf{W}|$$

$$\text{s.t. } \mathrm{Tr}(\mathbf{\Lambda}_\Theta \mathbf{\Lambda}_\Theta^{\mathrm{H}}) \leq K. \tag{26}$$

The equivalence between **Prob.3** and **Prob.4** is built based on the idea of weighted MSE minimization (WMMSE) [28], [29]. Then, **Prob.4** can be efficiently solved via the AO among $\mathbf{G}_\Lambda$, $\mathbf{\Lambda}_\Theta$ and $\mathbf{W}$. Specifically, both optimal $\mathbf{G}_\Lambda$ and $\mathbf{W}$ can be directly derived by taking the first-order derivatives of the objective function of **Prob.4** w.r.t. $\mathbf{G}_\Lambda$ and $\mathbf{W}$ to zeros, i.e., $\mathbf{G}_\Lambda = \mathbf{H}^{\mathrm{H}}\left(\mathbf{\Sigma} + \mathbf{H}\mathbf{H}^{\mathrm{H}}\right)^{-1}$, $\mathbf{W} = \left(\mathbf{I}_{N_t} - \mathbf{G}_\Lambda\mathbf{H}\right)^{-1}$. Then, the optimization problem w.r.t. $\mathbf{\Lambda}_\Theta$ can be written as

**Prob.5:** $\min_{\mathbf{\Lambda}_\Theta} \ \mathrm{Tr}\left(\mathbf{W}\left[\mathbf{G}_\Lambda \mathbf{H} - \mathbf{I}_{N_t}\right]\left[\mathbf{G}_\Lambda \mathbf{H} - \mathbf{I}_{N_t}\right]^{\mathrm{H}}\right)$

$$\text{s.t. } \mathrm{Tr}(\mathbf{\Lambda}_\Theta \mathbf{\Lambda}_\Theta^{\mathrm{H}}) \leq K. \tag{27}$$

The first-order derivative of the objective function of **Prob.5** w.r.t. $\mathbf{\Lambda}_\Theta$ is given by

$$\mathrm{Diag}\left\{\frac{\partial \mathrm{Tr}\left(\mathbf{W}\left[\mathbf{G}_\Lambda \mathbf{H} - \mathbf{I}_{N_t}\right]\left[\mathbf{G}_\Lambda \mathbf{H} - \mathbf{I}_{N_t}\right]^{\mathrm{H}}\right)}{\partial \mathbf{\Lambda}_\Theta}\right\}$$

$$=\mathrm{Diag}\{\mathbf{H}_2\left(\mathbf{H}^{\mathrm{H}}\mathbf{G}_\Lambda^{\mathrm{H}} - \mathbf{H}_2\right)\mathbf{W}\mathbf{G}_\Lambda\mathbf{H}_1\}, \tag{28}$$

based on which the KKT optimality conditions of **Prob.5** can be formulated as

$$\begin{cases} \mathrm{Diag}\{\mathbf{H}_2\left(\mathbf{H}^{\mathrm{H}}\mathbf{G}_\Lambda^{\mathrm{H}} - \mathbf{H}_2\right)\mathbf{W}\mathbf{G}_\Lambda\mathbf{H}_1\} = -\mu\mathrm{Diag}\{\mathbf{\Lambda}_\Theta^{\mathrm{H}}\}, & (29a) \\ \mu\left(\mathrm{Tr}(\mathbf{\Lambda}_\Theta \mathbf{\Lambda}_\Theta^{\mathrm{H}}) - K\right) = 0, & (29b) \end{cases}$$

where $\mu$ is the dual variable associated with the amplitude constraint. Based on (29a), the optimal $\mathbf{\Lambda}_\Theta$ is derived as

$$\mathrm{Diag}\{\mathbf{\Lambda}_\Theta^{\mathrm{H}}\} = (\mu\mathbf{I}_K + \mathbf{\Phi})^{-1}\mathbf{a}, \tag{30}$$

where $\mathbf{\Phi} = \mathbf{H}_2\mathbf{H}_2^{\mathrm{H}} \odot \left(\mathbf{H}_1^{\mathrm{H}}\mathbf{G}^{\mathrm{H}}\mathbf{W}\mathbf{G}\mathbf{H}_1\right)^{\mathrm{T}}$ and $\mathbf{a} = \mathrm{Diag}\{\mathbf{H}_2\mathbf{W}\mathbf{G}_\Lambda\mathbf{H}_1 - \mathbf{H}_2\mathbf{H}_0^{\mathrm{H}}\mathbf{G}_\Lambda^{\mathrm{H}}\mathbf{W}\mathbf{G}_\Lambda\mathbf{H}_1\}$. Moreover, since $\mathrm{Tr}(\mathbf{\Lambda}_\Theta \mathbf{\Lambda}_\Theta^{\mathrm{H}})$ is monotonically decreasing w.r.t. $\mu$, the optimal $\mu$ satisfying (29b) is found via the bisection search.

In addition, we formulate the MSE minimization problem for the amplitude-adjustable IRS-aided MIMO system as

**Prob.6:** $\min_{\mathbf{\Lambda}_\Theta} \ \mathrm{Tr}\left(\left[\mathbf{\Sigma}^{-1}\mathbf{H}\mathbf{H}^{\mathrm{H}} + \mathbf{I}_{N_r}\right]^{-1}\right)$

$$\text{s.t. } \mathrm{Tr}(\mathbf{\Lambda}_\Theta \mathbf{\Lambda}_\Theta^{\mathrm{H}}) \leq K. \tag{31}$$

Similar to **Prob.3**, **Prob.6** is difficult to solve since it involves a quadratic term w.r.t. $\mathbf{\Lambda}_\Theta$ appears in an inverse form. Fortunately, it can also be equivalently transformed into the WMMSE minimization problem **Prob.4** by setting $\mathbf{W} = \mathbf{I}_{N_t}$.

### C. Extension to Block-Diagonal Structure Constraints

In this subsection, we extend the complex matrix derivative to the block-diagonal matrix, which is essentially a kind of bidiagonal matrix. In a general multi-antenna MU-MIMO uplink system, the received signal from $K$ users at the BS can be written as [30]

$$\mathbf{y} = \mathbf{H}\mathbf{s} + \mathbf{n} \ \text{ with } \ \mathbf{H} = [\mathbf{H}_1, \cdots, \mathbf{H}_K], \tag{32}$$

where $\mathbf{H}_k \in \mathbb{C}^{N_t \times N_r}$ denotes the channel between the BS and the $k$th user, $\mathbf{s}_k \in \mathbb{C}^{N_r \times 1}$ denotes the transmitted data stream of the $k$th user, and all data streams $\mathbf{s}_k$'s are stacked into the vector $\mathbf{s} \in \mathbb{C}^{N_r K \times 1}$, i.e., $\mathbf{s} = [\mathbf{s}_1^{\mathrm{T}}, \cdots, \mathbf{s}_K^{\mathrm{T}}]^{\mathrm{T}}$. Accordingly, the covariance matrix of $\mathbf{s}$ is a block-diagonal matrix, i.e.,

$$\mathbf{Q} = \mathbb{E}\{\mathbf{s}\mathbf{s}^{\mathrm{H}}\} = \mathrm{Blockdiag}\left(\{\mathbf{Q}_k\}_{k=1}^K\right), \tag{33}$$

where $\mathbf{Q}_k = \mathbb{E}\{\mathbf{s}_k\mathbf{s}_k^{\mathrm{H}}\}$ is the transmit covariance matrix of $\mathbf{s}_k$. Hereafter, we mainly consider the capacity maximization problem under the general user grouping power constraints, which is formulated as

**Prob.7:** $\max_{\mathbf{Q}} \ \log\left|\mathbf{I}_{N_t} + \mathbf{\Sigma}^{-1}\mathbf{H}\mathbf{Q}\mathbf{H}^{\mathrm{H}}\right|$

$$\text{s.t. } \sum_{k \in \phi_n} \mathrm{Tr}(\mathbf{Q}_k) \leq P_n, \mathbf{Q}_k \succeq \mathbf{0}, \ \forall k, n, \tag{34}$$

where $\phi_n$ is the set of user indices in the $n$-th user group and $P_n$ is the corresponding available transmit power. Generally, we have $\bigcup_{n=1}^N \phi_n = \{1, 2, \cdots, K\}$ and $\phi_n \cap \phi_m = \varnothing, \forall m \neq n$. Specifically, we set $n = 1, \cdots, K$, $\phi_n = \{k\}$ in the MU-MIMO uplink system, and $n = 1$, $\phi_n = \{1, 2, \cdots, K\}$ in the virtual MU-MIMO uplink system based on the uplink-downlink duality [17]. The differential of the objective function of **Prob.7** w.r.t. $\mathbf{Q}$ is given by $\mathrm{d}\left(\log\left|\mathbf{I}_{N_t} + \mathbf{\Sigma}^{-1}\mathbf{H}\mathbf{Q}\mathbf{H}^{\mathrm{H}}\right|\right) = \mathrm{Tr}\left(\left(\mathbf{I}_{N_t} + \mathbf{\Sigma}^{-\frac{1}{2}}\mathbf{H}\mathbf{Q}\mathbf{H}^{\mathrm{H}}\mathbf{\Sigma}^{-\frac{1}{2}}\right)^{-1}\mathbf{\Sigma}^{-\frac{1}{2}}\mathbf{H}\mathrm{d}(\mathbf{Q})\mathbf{H}^{\mathrm{H}}\mathbf{\Sigma}^{-\frac{1}{2}}\right)$. The corresponding first-order derivative w.r.t $\mathbf{Q}$ is then derived as

$$\frac{\partial \log\left|\mathbf{I}_{N_t} + \mathbf{\Sigma}^{-1}\mathbf{H}\mathbf{Q}\mathbf{H}^{\mathrm{H}}\right|}{\partial \mathbf{Q}}$$

$$=\mathrm{Blockdiag}\left(\left\{\mathbf{H}_k^{\mathrm{H}}\mathbf{\Sigma}^{-\frac{1}{2}}\left(\mathbf{I}_{N_t} + \mathbf{\Sigma}^{-\frac{1}{2}}\mathbf{H}\mathbf{Q}\mathbf{H}^{\mathrm{H}}\mathbf{\Sigma}^{-\frac{1}{2}}\right)^{-1}\right.$$

$$\left. \times \mathbf{\Sigma}^{-\frac{1}{2}}\mathbf{H}_k\right\}_{k=1}^K\right). \tag{35}$$

TABLE II
ELEMENT-WISE PHASE DERIVATIVES UNDER THE CONSTANT MODULUS CONSTRAINT

| Function Type | Element-Wise Phase Derivative w.r.t. $[\mathbf{\Theta}]_{i,j}, \forall i,j$ |
|---|---|
| $f_{\mathrm{C,TL}} = \mathrm{Tr}(\mathbf{B}^{\mathrm{H}}\mathbf{X}) + \mathrm{Tr}(\mathbf{B}\mathbf{X}^{\mathrm{H}})$ | $-2\Im\{[\mathbf{B}^*]_{i,j}[\mathbf{X}]_{i,j}\}$ |
| $f_{\mathrm{C,TQ}} = \mathrm{Tr}(\mathbf{X}\mathbf{\Pi}\mathbf{X}^{\mathrm{H}}\mathbf{\Phi})$ | $-2\Im\{[(\mathbf{\Phi}\mathbf{X}\mathbf{\Pi})^*]_{i,j}[\mathbf{X}]_{i,j}\}$ |
| $f_{\mathrm{C,TI}} = \mathrm{Tr}\left((\mathbf{\Phi} + \mathbf{X}^{\mathrm{H}}\mathbf{\Pi}\mathbf{X})^{-1}\right)$ | $2\Im\left\{\left[(\mathbf{\Pi}\mathbf{X}(\mathbf{\Phi} + \mathbf{X}^{\mathrm{H}}\mathbf{\Pi}\mathbf{X})^{-2})^*\right]_{i,j}[\mathbf{X}]_{i,j}\right\}$ |
| $f_{\mathrm{C,LD}} = \log|\mathbf{\Phi} + \mathbf{X}^{\mathrm{H}}\mathbf{\Pi}\mathbf{X}|$ | $-2\Im\left\{\left[(\mathbf{\Pi}\mathbf{X}(\mathbf{\Phi} + \mathbf{X}^{\mathrm{H}}\mathbf{\Pi}\mathbf{X})^{-1})^*\right]_{i,j}[\mathbf{X}]_{i,j}\right\}$ |

Following that, the KKT optimality conditions are given by

$$
\begin{cases}
\mathrm{Blockdiag}\left(\left\{\mathbf{H}_k^{\mathrm{H}}\mathbf{\Sigma}^{-\frac{1}{2}}\left(\mathbf{I}_{N_t} + \mathbf{\Sigma}^{-\frac{1}{2}}\mathbf{H}\mathbf{Q}\mathbf{H}^{\mathrm{H}}\mathbf{\Sigma}^{-\frac{1}{2}}\right)^{-1}\right.\right. \\
\qquad \left.\left.\times\mathbf{\Sigma}^{-\frac{1}{2}}\mathbf{H}_k\right\}_{k=1}^K\right) = \mathrm{Blockdiag}\left(\{\mu_k\mathbf{I}_{N_r} - \mathbf{\Psi}_k\}_{k=1}^K\right), \quad (36\mathrm{a}) \\
\mu_n\left(\sum_{k\in\phi_n}\mathrm{Tr}(\mathbf{Q}_k) - P_n\right) = 0, \mu_k = \mu_n, \forall k\in\phi_n, \quad (36\mathrm{b}) \\
\mathrm{Tr}(\mathbf{\Psi}_k\mathbf{Q}_k) = 0, \ \forall k, \quad (36\mathrm{c})
\end{cases}
$$

where $\mu_k$ and $\mathbf{\Psi}_k$ are the Lagrange multipliers associated with the transmit power constraint and the positive semi-definite constraint at the $k$th user, respectively. Then, (36a) can be rewritten in terms of $\mathbf{Q}_k$ as follows:

$$
\mathbf{H}_k^{\mathrm{H}}\mathbf{\Sigma}^{-\frac{1}{2}}\left(\mathbf{I}_{N_t} + \mathbf{\Sigma}^{-\frac{1}{2}}\mathbf{H}\mathbf{Q}\mathbf{H}^{\mathrm{H}}\mathbf{\Sigma}^{-\frac{1}{2}}\right)^{-1}\mathbf{\Sigma}^{-\frac{1}{2}}\mathbf{H}_k
$$
$$
= \mathbf{H}_k^{\mathrm{H}}\mathbf{\Sigma}^{-\frac{1}{2}}\mathbf{L}_k^{-\frac{1}{2}}\left(\mathbf{I}_{N_t} + \mathbf{L}_k^{-\frac{1}{2}}\mathbf{\Sigma}^{-\frac{1}{2}}\mathbf{H}_k\mathbf{Q}_k\mathbf{H}_k^{\mathrm{H}}\mathbf{\Sigma}^{-\frac{1}{2}}\mathbf{L}_k^{-\frac{1}{2}}\right)^{-1}
$$
$$
\times \mathbf{L}_k^{-\frac{1}{2}}\mathbf{\Sigma}^{-\frac{1}{2}}\mathbf{H}_k = \mu_k\mathbf{I}_{N_r} - \mathbf{\Psi}_k, \ \forall k, \quad (37)
$$

where $\mathbf{L}_k = \mathbf{I}_{N_t} + \sum_{j\neq k}\mathbf{\Sigma}^{-\frac{1}{2}}\mathbf{H}_j\mathbf{Q}_j\mathbf{H}_j^{\mathrm{H}}\mathbf{\Sigma}^{-\frac{1}{2}}$. Thus, based on eigenspace alignment, the optimal $\mathbf{Q}_k$'s can be derived as that in [16, Theorem 1], i.e., $\mathbf{Q}_k = \mathbf{V}_{\mathcal{H}_k}\mathbf{\Lambda}_{\mathbf{Q}_k}\mathbf{V}_{\mathcal{H}_k}^{\mathrm{H}}, \ \forall k$, where $\mathbf{\Lambda}_{\mathbf{Q}_k}$ is a diagonal matrix, each diagonal element of which has a water-filling form, and $\mathbf{V}_{\mathcal{H}_k}$ is an unitary matrix coming from the singular value decomposition (SVD) represented as $\mathbf{L}_k^{-\frac{1}{2}}\mathbf{\Sigma}^{-\frac{1}{2}}\mathbf{H}_k = \mathbf{U}_{\mathcal{H}_k}\mathbf{\Lambda}_{\mathcal{H}_k}\mathbf{V}_{\mathcal{H}_k}^{\mathrm{H}}$ with $\mathbf{\Lambda}_{\mathcal{H}_k}\searrow$, where $\mathbf{\Lambda}_{\mathcal{H}_k}\searrow$ implies that the diagonal elements of $\mathbf{\Lambda}_{\mathcal{H}_k}$ are arranged in descending order. Similarly, $\mu_k$ that satisfying (36b) can be obtained by the bisection search.

*Remark 1:* Based on the above discussions, we can conclude that globally optimal solutions of several classical optimization problems in the state-of-the-art wireless systems can be directly obtained with low complexity by the proposed complex matrix derivatives under diagonal structure constraints. In addition, for optimization problems that not satisfy diagonal structure constraints directly, the proposed algorithm is also able to obtain an approximate solution by further exploring the inherent structure of the optimal solution. For example, the optimal matrix variables in the point-to-point MIMO system operating at high SNR conditions and the MU-MISO downlink system employing the BD-ZF strategy [31] are both validated to be approximately diagonal.

## III. CONSTANT MODULUS CONSTRAINTS

Different from diagonal structure constraints, constant modulus constraints are imposed on matrix variables in an element-wise manner, which makes the optimization problem challenging to directly solve using complex matrix derivatives. Motivated by this fact, we firstly provide some mathematical preliminaries for the element-wise phase derivatives of several widely adopted objective functions. Then, we investigate specific optimization problems in both the hybrid analog-digital

MIMO system and the fully-passive IRS-aided MIMO system. In order to avoid complicated matrix inversion and matrix factorization, a novel AO algorithm with the aid of several arbitrary feasible solutions is proposed.

### A. Mathematical Preliminaries

We firstly introduce a complex matrix variable $\mathbf{X}\in\mathbb{C}^{N\times M}$ subject to constant modulus constraints as follows:

$$
[\mathbf{X}]_{i,j} = e^{j\theta_{i,j}}, \ \forall i,j \text{ and } \mathrm{Tr}(\mathbf{X}\mathbf{X}^{\mathrm{H}}) = NM, \quad (38)
$$

where $\theta_{i,j}\in[0,2\pi]$ denotes the phase of $[\mathbf{X}]_{i,j}$. Based on (38), the first-order derivative w.r.t. the constant modulus constrained $\mathbf{X}$ can be replaced by the first-order derivative w.r.t. the corresponding unconstrained phase matrix $\mathbf{\Theta}$, where $[\mathbf{\Theta}]_{i,j} = \theta_{i,j}, \forall i,j$. Accordingly, the element-wise phase derivatives of the function $f(\mathbf{X})$ w.r.t. $\mathbf{\Theta}$ can be defined as

$$
\left[\frac{\partial f(\mathbf{X})}{\partial\mathbf{\Theta}}\right]_{i,j} = \frac{\partial f(\mathbf{X})}{\partial[\mathbf{\Theta}]_{i,j}}, \ \forall i,j. \quad (39)
$$

Similar to Sec. II-A, we also consider the element-wise phase derivatives for four common objective functions, i.e.,

*1) Trace-Linear Function:* Since the phase $[\mathbf{\Theta}]_{i,j}$'s are real scalar, for arbitrary complex matrix $\mathbf{B}\in\mathbb{C}^{N\times M}$, the element-wise phase derivatives of a trace-linear function $f_{\mathrm{C,TL}} = \mathrm{Tr}(\mathbf{B}^{\mathrm{H}}\mathbf{X}) + \mathrm{Tr}(\mathbf{B}\mathbf{X}^{\mathrm{H}})$ w.r.t. $[\mathbf{\Theta}]_{i,j}$'s can be obtained as

$$
\left[\frac{\partial f_{\mathrm{C,TL}}}{\partial\mathbf{\Theta}}\right]_{i,j} = \left[\frac{\partial\sum_m\sum_n[\mathbf{B}^*]_{n,m}[\mathbf{X}]_{n,m}}{\partial\mathbf{\Theta}}\right]_{i,j}
$$
$$
+ \left[\partial\frac{\sum_m\sum_n[\mathbf{B}]_{n,m}[\mathbf{X}^*]_{n,m}}{\partial\mathbf{\Theta}}\right]_{i,j}
$$
$$
= j[\mathbf{B}^*]_{i,j}[\mathbf{X}]_{i,j} - j[\mathbf{B}]_{i,j}[\mathbf{X}^*]_{i,j}
$$
$$
= -2\Im\{[\mathbf{B}^*]_{i,j}[\mathbf{X}]_{i,j}\}, \ \forall i,j. \quad (40)
$$

*2) Trace-Quadratic Function:* For arbitrary Hermitian matrices $\mathbf{\Pi}\in\mathbb{C}^{M\times M}$ and $\mathbf{\Phi}\in\mathbb{C}^{N\times N}$, the element-wise phase derivatives of a trace-quadratic function $f_{\mathrm{C,TQ}} = \mathrm{Tr}(\mathbf{X}\mathbf{\Pi}\mathbf{X}^{\mathrm{H}}\mathbf{\Phi})$ w.r.t. $[\mathbf{\Theta}]_{i,j}$'s are given by

$$
\left[\frac{\partial f_{\mathrm{C,TQ}}}{\partial\mathbf{\Theta}}\right]_{i,j} = j[(\mathbf{\Phi}\mathbf{X}\mathbf{\Pi})^*]_{i,j}[\mathbf{X}]_{i,j} - j[(\mathbf{\Phi}\mathbf{X}\mathbf{\Pi})]_{i,j}[\mathbf{X}^*]_{i,j}
$$
$$
= -2\Im\{[(\mathbf{\Phi}\mathbf{X}\mathbf{\Pi})^*]_{i,j}[\mathbf{X}]_{i,j}\}, \ \forall i,j. \quad (41)
$$

*3) Trace-Inverse Function:* Regarding a trace-inverse function $f_{\mathrm{C,TI}} = \mathrm{Tr}\left((\mathbf{\Phi} + \mathbf{X}^{\mathrm{H}}\mathbf{\Pi}\mathbf{X})^{-1}\right)$, we have the following element-wise phase derivatives w.r.t. $[\mathbf{\Theta}]_{i,j}$'s.

$$
\frac{\partial f_{\mathrm{C,TI}}}{\partial\mathbf{\Theta}} = -j\left[(\mathbf{\Phi} + \mathbf{X}^{\mathrm{H}}\mathbf{\Pi}\mathbf{X})^{-2}\mathbf{X}^{\mathrm{H}}\mathbf{\Pi}\right]_{i,j}[\mathbf{X}]_{i,j}
$$
$$
+ j\left[(\mathbf{\Pi}\mathbf{X}(\mathbf{\Phi} + \mathbf{X}^{\mathrm{H}}\mathbf{\Pi}\mathbf{X})^{-2})\right]_{i,j}[\mathbf{X}^*]_{i,j} \quad (42)
$$
$$
= 2\Im\left\{\left[(\mathbf{\Pi}\mathbf{X}(\mathbf{\Phi} + \mathbf{X}^{\mathrm{H}}\mathbf{\Pi}\mathbf{X})^{-2})^*\right]_{i,j}[\mathbf{X}]_{i,j}\right\}, \ \forall i,j.
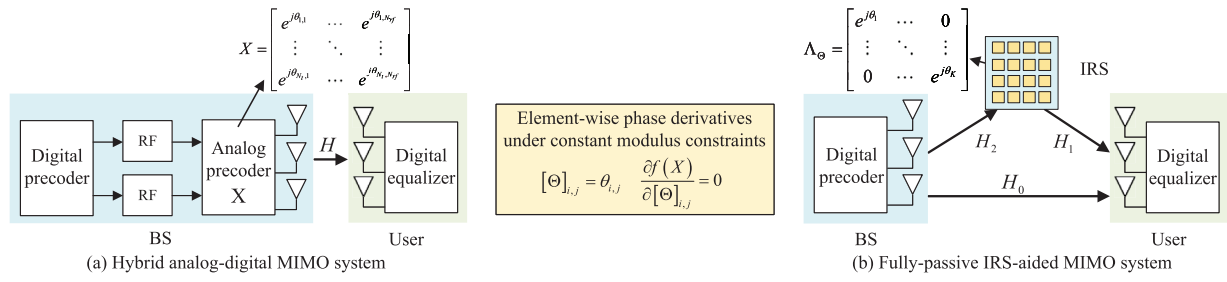$$

Fig. 2. A diagram of application scenarios associated with constant modulus matrix variables.

*4) Log-Determinant Function:* Similarly, the element-wise phase derivatives of a log-determinant function $f_{\text{C,LD}} = \log|\boldsymbol{\Phi} + \boldsymbol{X}^{\text{H}}\boldsymbol{\Pi}\boldsymbol{X}|$ w.r.t. $[\boldsymbol{\Theta}]_{i,j}$'s are given by

$$\left[\frac{\partial f_{\text{C,LD}}}{\partial \boldsymbol{\Theta}}\right]_{i,j} = -2\Im\left\{\left[\left(\boldsymbol{\Pi}\boldsymbol{X}(\boldsymbol{\Phi} + \boldsymbol{X}^{\text{H}}\boldsymbol{\Pi}\boldsymbol{X})^{-1}\right)^*\right]_{i,j}\right.$$
$$\left. \times [\boldsymbol{X}]_{i,j}\right\}, \ \forall i,j. \quad (43)$$

The element-wise phase derivatives of the above four types of objective functions are summarized in Table II. In the following subsection, several state-of-the-art wireless applications will be investigated in detail based on the above fundamental properties.

*B. Specific Wireless Applications*

In Fig. 2, there are two typical wireless applications associated with constant modulus constraints, i.e., the analog beamforming optimization in the hybrid analog-digital MIMO system and the phase shift optimization in the fully-passive IRS-aided MIMO system, which are elaborated as follows.

*1) Hybrid Analog-Digital MIMO System:* As shown in Fig. 2(a), we firstly consider the hybrid analog-digital beamforming design in the downlink point-to-point MIMO system, where the BS equipped with $N_t$ antennas and $N_{rf}$ radio-frequency (RF) chains transfers $N_s$ data streams to a $N_r$-antenna user with the fully-digital equalizer [12]. Then, the received signal at the user can be expressed as

$$\boldsymbol{y} = \boldsymbol{G}\boldsymbol{H}\boldsymbol{X}\boldsymbol{F}_{\text{D}}\boldsymbol{s} + \boldsymbol{G}\boldsymbol{n}, \quad (44)$$

where $\boldsymbol{G} \in \mathbb{C}^{N_s \times N_r}$ denotes the fully-digital receive equalizer, $\boldsymbol{H} \in \mathbb{C}^{N_r \times N_t}$ denotes the channel between the BS and the user, $\boldsymbol{X} \in \mathbb{C}^{N_t \times N_{rf}}$ and $\boldsymbol{F}_{\text{D}} \in \mathbb{C}^{N_{rf} \times N_s}$ are the constant modulus analog beamformer and the digital beamformer, respectively. $\boldsymbol{s} \in \mathbb{C}^{N_s \times 1}$ is the transmit data streams with unit covariance matrix, i.e., $\mathbb{E}\{\boldsymbol{s}\boldsymbol{s}^{\text{H}}\} = \boldsymbol{I}_{N_s}$. $\boldsymbol{n} \in \mathbb{C}^{N_r \times 1}$ is the additive Gaussian noise with zero mean and covariance matrix $\mathbb{E}\{\boldsymbol{n}\boldsymbol{n}^{\text{H}}\} = \boldsymbol{\Sigma}$. Based on (44), the MSE matrix is given by

$$\boldsymbol{E}_{\text{MSE}} = \mathbb{E}\left[(\hat{\boldsymbol{s}} - \boldsymbol{s})(\hat{\boldsymbol{s}} - \boldsymbol{s})^{\text{H}}\right]$$
$$= (\boldsymbol{G}\boldsymbol{H}\boldsymbol{X}\boldsymbol{F}_{\text{D}} - \boldsymbol{I}_{N_s})(\boldsymbol{G}\boldsymbol{H}\boldsymbol{X}\boldsymbol{F}_{\text{D}} - \boldsymbol{I}_{N_s})^{\text{H}} + \boldsymbol{G}\boldsymbol{\Sigma}\boldsymbol{G}^{\text{H}}$$
$$\overset{(c)}{=} \left[\boldsymbol{I}_{N_s} + \boldsymbol{F}_{\text{D}}^{\text{H}}\boldsymbol{X}^{\text{H}}\boldsymbol{H}^{\text{H}}\boldsymbol{\Sigma}^{-1}\boldsymbol{H}\boldsymbol{X}\boldsymbol{F}_{\text{D}}\right]^{-1}, \quad (45)$$

where $\hat{\boldsymbol{s}}$ denotes the estimated signal and $(c)$ holds based on the optimal unconstrained Wiener filter $\boldsymbol{G} = \boldsymbol{F}_{\text{D}}^{\text{H}}\boldsymbol{X}^{\text{H}}\boldsymbol{H}^{\text{H}}\left(\boldsymbol{H}\boldsymbol{X}\boldsymbol{F}_{\text{D}}\boldsymbol{F}_{\text{D}}^{\text{H}}\boldsymbol{X}^{\text{H}}\boldsymbol{H}^{\text{H}} + \boldsymbol{R}_{\text{n}}\right)^{-1}$ [32]. Without loss of generality, we usually assume $\boldsymbol{F}_{\text{D}}\boldsymbol{F}_{\text{D}}^{\text{H}} \approx \gamma^2 \boldsymbol{I}_{N_{rf}}$ for large-scale MIMO systems [20]. Under this assumption, the capacity maximization problem of the hybrid analog-digital MIMO system can be formulated as

**Prob.8:** $\max_{\boldsymbol{X}} \ \log|\boldsymbol{I}_{N_{rf}} + \boldsymbol{X}^{\text{H}}\boldsymbol{\Pi}\boldsymbol{X}|$

$$\text{s.t. } |[\boldsymbol{X}]_{i,j}| = 1, \ \forall i,j, \quad (46)$$

where $\boldsymbol{\Pi} = \gamma^2 \boldsymbol{H}^H \boldsymbol{\Sigma}^{-1}\boldsymbol{H}$ represents the effective signal-to-noise ratio (SNR). According to the KKT optimality conditions, the element-wise phase derivatives of the objective function of **Prob.8** w.r.t. $[\boldsymbol{X}]_{i,j}$'s must equal zeros at the optimal $[\boldsymbol{X}]_{i,j}$'s. Specifically, we recall (43) to obtain

$$\underbrace{\Im\left\{\left[\left(\boldsymbol{\Pi}\boldsymbol{X}(\boldsymbol{I}_{N_{rf}} + \boldsymbol{X}^{\text{H}}\boldsymbol{\Pi}\boldsymbol{X})^{-1}\right)^*\right]_{i,j}[\boldsymbol{X}]_{i,j}\right\}}_{\triangleq g_{i,j}(\boldsymbol{X})} = 0, \ \forall i,j. \quad (47)$$

Define $\boldsymbol{A}_j = \boldsymbol{I}_{N_t} + \widetilde{\boldsymbol{X}}_j\widetilde{\boldsymbol{X}}_j^{\text{H}}\boldsymbol{\Pi}$ and $n_j = 1 + [\boldsymbol{X}]_{:,j}^{\text{H}}\boldsymbol{\Pi}\boldsymbol{A}_j^{-1}[\boldsymbol{X}]_{:,j}$, where $\widetilde{\boldsymbol{X}}_j$ denotes the sub-matrix of $\boldsymbol{X}$ with the $j$th column removed. Then, $\left[\boldsymbol{\Pi}\boldsymbol{X}(\boldsymbol{I}_{N_{rf}} + \boldsymbol{X}^{\text{H}}\boldsymbol{\Pi}\boldsymbol{X})^{-1}\right]_{i,j}$'s can be rewritten as

$$\left[\boldsymbol{\Pi}\boldsymbol{X}(\boldsymbol{I}_{N_{rf}} + \boldsymbol{X}^{\text{H}}\boldsymbol{\Pi}\boldsymbol{X})^{-1}\right]_{i,j} \quad (48)$$
$$\overset{(d)}{=} \left[\boldsymbol{\Pi}\left(\boldsymbol{A}_j^{-1} - \frac{\boldsymbol{A}_j^{-1}[\boldsymbol{X}]_{:,j}[\boldsymbol{X}]_{:,j}^{\text{H}}\boldsymbol{\Pi}\boldsymbol{A}_j^{-1}}{n_j}\right)\boldsymbol{X}\right]_{i,j}$$
$$= \frac{1}{n_j}\sum_{l \neq i}[\boldsymbol{\Pi}\boldsymbol{A}_j^{-1}]_{l,i}^{\text{H}}[\boldsymbol{X}]_{l,j} + \frac{1}{n_j}[\boldsymbol{\Pi}\boldsymbol{A}_j^{-1}]_{i,i}^{\text{H}}[\boldsymbol{X}]_{i,j}, \ \forall i,j,$$

where $(d)$ holds similarly to (14). By substituting (48) into (47), we have

$$\Im\left\{\left[\left(\boldsymbol{\Pi}\boldsymbol{X}(\boldsymbol{I}_{N_{rf}} + \boldsymbol{X}^{\text{H}}\boldsymbol{\Pi}\boldsymbol{X})^{-1}\right)^*\right]_{i,j}[\boldsymbol{X}]_{i,j}\right\} \quad (49a)$$
$$= \Im\left\{\left(\frac{1}{n_j}\sum_{l \neq i}[\boldsymbol{\Pi}\boldsymbol{A}_j^{-1}]_{l,i}^{\text{H}}[\boldsymbol{X}]_{l,j}\right)^*[\boldsymbol{X}]_{i,j}\right.$$
$$\left. + \left(\frac{1}{n_j}[\boldsymbol{\Pi}\boldsymbol{A}_j^{-1}]_{i,i}^{\text{H}}\right)^*\right\} = 0, \ \forall i,j. \quad (49b)$$

Since the following equality holds, i.e.,

$$\boldsymbol{\Pi}\boldsymbol{A}_j^{-1} = \boldsymbol{\Pi}(\boldsymbol{I}_{N_t} + \widetilde{\boldsymbol{X}}_j\widetilde{\boldsymbol{X}}_j^{\text{H}}\boldsymbol{\Pi})^{-1}$$
$$= (\boldsymbol{I}_{N_t} + \boldsymbol{\Pi}\widetilde{\boldsymbol{X}}_j\widetilde{\boldsymbol{X}}_j^{\text{H}})^{-1}\boldsymbol{\Pi} = (\boldsymbol{\Pi}\boldsymbol{A}_j^{-1})^{\text{H}}, \ \forall j, \quad (50)$$

we can conclude that $\boldsymbol{\Pi}\boldsymbol{A}_j^{-1}$ is a Hermitian matrix. As such, it is readily inferred that $n_j$ and $\frac{1}{n_j}[\boldsymbol{\Pi}\boldsymbol{A}_j^{-1}]_{i,i}^{\text{H}}$ are both real scalars. Recall the definition of $[\boldsymbol{X}]_{i,j} = e^{j\theta_{i,j}}, \forall i,j$, it follows from (49b) that the optimal $[\boldsymbol{\Theta}]_{i,j}$'s to **Prob.8** are derived as

$$[\boldsymbol{\Theta}]_{i,j} = \text{Phase}\left\{\sum_{l \neq i}[\boldsymbol{\Pi}\boldsymbol{A}_j^{-1}]_{l,i}^{\text{H}}[\boldsymbol{X}]_{l,j}\right\}$$
$$\text{or } \pi + \text{Phase}\left\{\sum_{l \neq i}[\boldsymbol{\Pi}\boldsymbol{A}_j^{-1}]_{l,i}^{\text{H}}[\boldsymbol{X}]_{l,j}\right\}, \ \forall i,j. \quad (51)$$

In addition, the MSE minimization problem is expressed as

**Prob.9:** $\min_{\boldsymbol{X}} \ \text{Tr}\left(\left(\boldsymbol{I}_{N_{rf}} + \boldsymbol{X}^{\text{H}}\boldsymbol{\Pi}\boldsymbol{X}\right)^{-1}\right)$

$$\text{s.t. } |[\boldsymbol{X}]_{i,j}| = 1, \ \forall i,j. \quad (52)$$

Similarly, to find the optimal solution of **Prob.9**, the element-wise phase derivatives of the objective function w.r.t. $[\boldsymbol{X}]_{i,j}$'s must equal zeros, that is,

$$\underbrace{\Im\left\{\left[\left(\boldsymbol{\Pi}\boldsymbol{X}\left(\boldsymbol{I}_{N_{rf}}+\boldsymbol{X}^{\mathrm{H}}\boldsymbol{\Pi}\boldsymbol{X}\right)^{-2}\right)^*\right]_{i,j}[\boldsymbol{X}]_{i,j}\right\}}_{\triangleq g_{i,j}(\boldsymbol{X})}=0,\ \forall i,j. \quad (53)$$

The terms $\left[\boldsymbol{\Pi}\boldsymbol{X}\left(\boldsymbol{I}_{N_{rf}}+\boldsymbol{X}^{\mathrm{H}}\boldsymbol{\Pi}\boldsymbol{X}\right)^{-2}\right]_{i,j}$'s can be rewritten as

$$\left[\boldsymbol{\Pi}\boldsymbol{X}(\boldsymbol{I}_{N_{rf}}+\boldsymbol{X}^{\mathrm{H}}\boldsymbol{\Pi}\boldsymbol{X})^{-2}\right]_{i,j}$$
$$\overset{(e_1)}{=}\left[\boldsymbol{\Pi}\left(\boldsymbol{A}_j^{-1}-\frac{\boldsymbol{A}_j^{-1}[\boldsymbol{X}]_{:,j}[\boldsymbol{X}]_{:,j}^{\mathrm{H}}\boldsymbol{\Pi}\boldsymbol{A}_j^{-1}}{n_j}\right)^2\boldsymbol{X}\right]_{i,j}$$
$$\overset{(e_2)}{=}\frac{1}{n_j}[\boldsymbol{\Pi}\boldsymbol{A}_j^{-2}]_{:,i}^{\mathrm{H}}[\boldsymbol{X}]_{:,j}-\frac{m_j}{n_j^2}[\boldsymbol{\Pi}\boldsymbol{A}_j^{-1}]_{:,i}^{\mathrm{H}}[\boldsymbol{X}]_{:,j}$$
$$\overset{(e_3)}{=}\frac{A_{1,i,j}^*+B_{1,i,j}^*[\boldsymbol{X}]_{i,j}^2+C_{1,i,j}^*[\boldsymbol{X}]_{i,j}}{n_j^2},\ \forall i,j, \quad (54)$$

where $(e_1)$ holds due to the same reasons as (14), $(e_2)$ holds by defining $m_j=[\boldsymbol{X}]_{:,j}^{\mathrm{H}}\boldsymbol{\Pi}\boldsymbol{A}_j^{-2}[\boldsymbol{X}]_{:,j}$, which is a real scalar and this can be proved similarly to (50). $(e_3)$ is obtained by rewriting $n_j$ and $m_j$ in terms of $[\boldsymbol{X}]_{i,j}$ as

$$n_j=\zeta_j^n+2\Re\left\{\eta_j^n[\boldsymbol{X}]_{i,j}^*\right\},m_j=\zeta_j^m+2\Re\left\{\eta_j^m[\boldsymbol{X}]_{i,j}^*\right\}, \quad (55)$$

where $\zeta_j^n=1+[\boldsymbol{\Pi}\boldsymbol{A}_j^{-1}]_{i,i}+\Re\left\{\sum_{p\neq i,q\neq i}[\boldsymbol{X}]_{p,j}^*[\boldsymbol{\Pi}\boldsymbol{A}_j^{-1}]_{p,q}\right.$ $\left.[\boldsymbol{X}]_{q,j}\right\}$, $\zeta_j^m=[\boldsymbol{\Pi}\boldsymbol{A}_j^{-2}]_{i,i}+\Re\left\{\sum_{p\neq i,q\neq i}[\boldsymbol{X}]_{p,j}^*[\boldsymbol{\Pi}\boldsymbol{A}_j^{-2}]_{p,q}\right.$ $\left.[\boldsymbol{X}]_{q,j}\right\}$, $\eta_j^n=\sum_{l\neq i}[\boldsymbol{\Pi}\boldsymbol{A}_j^{-1}]_{l,i}[\boldsymbol{X}]_{l,j}^*$, $\eta_j^m=\sum_{l\neq i}[\boldsymbol{\Pi}\boldsymbol{A}_j^{-2}]_{l,i}[\boldsymbol{X}]_{l,j}^*$, $\forall i,j$. Moreover, $A_{1,i,j}$, $B_{1,i,j}$ and $C_{1,i,j}$ in (54) are defined as

$$A_{1,i,j}=(\eta_j^n)^*[\boldsymbol{\Pi}\boldsymbol{A}_j^{-2}]_{i,i}^*-(\eta_j^m)^*[\boldsymbol{\Pi}\boldsymbol{A}_j^{-1}]_{i,i}^*$$
$$+(\zeta_j^n)^*(\eta_j^m)^*-(\zeta_j^m)^*(\eta_j^n)^*,$$
$$B_{1,i,j}=\eta_j^n[\boldsymbol{\Pi}\boldsymbol{A}_j^{-2}]_{i,i}^*-\eta_j^m[\boldsymbol{\Pi}\boldsymbol{A}_j^{-1}]_{i,i}^*,$$
$$C_{1,i,j}=(\zeta_j^n)^*[\boldsymbol{\Pi}\boldsymbol{A}_j^{-2}]_{i,i}^*-(\zeta_j^m)^*[\boldsymbol{\Pi}\boldsymbol{A}_j^{-1}]_{i,i}^*$$
$$+\eta_j^n(\eta_j^m)^*-\eta_j^m(\eta_j^n)^*,\ \forall i,j. \quad (56)$$

By substituting (54) into (53), we have

$$\Im\left\{\left[\left(\boldsymbol{\Pi}\boldsymbol{X}(\boldsymbol{I}_{N_{rf}}+\boldsymbol{X}^{\mathrm{H}}\boldsymbol{\Pi}\boldsymbol{X})^{-2}\right)^*\right]_{i,j}[\boldsymbol{X}]_{i,j}\right\} \quad (57a)$$
$$=\Im\left\{\frac{A_{1,i,j}[\boldsymbol{X}]_{i,j}+B_{1,i,j}[\boldsymbol{X}]_{i,j}^*+C_{1,i,j}}{n_j^2}\right\}=0,\ \forall i,j. \quad (57b)$$

Then, since $\Im\{\frac{a}{b}\}=0$ is equal to $\Im\{a\}=0$ for a real scalar $b$, (57b) can be further simplified as

$$\Im\left\{A_{1,i,j}e^{j\theta_{i,j}}+B_{1,i,j}e^{-j\theta_{i,j}}+C_{1,i,j}\right\} \quad (58a)$$
$$\overset{(f_1)}{=}|A_{1,i,j}|\sin(\theta_{i,j}+\alpha_{i,j})-|B_{1,i,j}|\sin(\theta_{i,j}-\beta_{i,j})+\Im\{C_{1,i,j}\}$$
$$\overset{(f_2)}{=}\sqrt{z_{1,i,j}^2+z_{2,i,j}^2}\sin(\theta_{i,j}+\phi_{i,j})+\Im\{C_{1,i,j}\}=0,\ \forall i,j, \quad (58b)$$

where $(f_1)$ is obtained using Euler's formula with $A_{1,i,j}=|A_{1,i,j}|e^{j\alpha_{i,j}}$ and $B_{1,i,j}=|B_{1,i,j}|e^{j\beta_{i,j}}$, $(f_2)$ holds due to the sum-to-product trigonometric identity with $z_{1,i,j}=$

$|A_{1,i,j}|\cos(\alpha_{i,j})+|B_{1,i,j}|\cos(\beta_{i,j})$, $z_{2,i,j}=|A_{1,i,j}|\sin(\alpha_{i,j})+|B_{1,i,j}|\sin(\beta_{i,j})$ and

$$\phi_{i,j}=\begin{cases}\arctan\left(\frac{z_{2,i,j}}{z_{1,i,j}}\right), & \text{if } z_{1,i,j}\geq 0,\\ \pi-\arctan\left(\frac{z_{2,i,j}}{z_{1,i,j}}\right), & \text{if } z_{1,i,j}<0,\end{cases}\forall i,j. \quad (59)$$

Based on (58b), the optimal $[\boldsymbol{\Theta}]_{i,j}$'s to **Prob.9** are obtained as

$$[\boldsymbol{\Theta}]_{i,j}=-\phi_{i,j}+\arcsin\left(-\frac{\Im\{C_{1,i,j}\}}{\sqrt{z_{1,i,j}^2+z_{2,i,j}^2}}\right) \quad (60)$$
$$\text{or }\pi-\phi_{i,j}+\arcsin\left(-\frac{\Im\{C_{1,i,j}\}}{\sqrt{z_{1,i,j}^2+z_{2,i,j}^2}}\right),\ \forall i,j.$$

*Remark 2:* Since the capacity maximization problem w.r.t. each user's analog beamformer $\boldsymbol{X}_u$ in the uplink hybrid analog-digital MU-MIMO system can be formulated in a similar form to **Prob.8** that for the single-user case by separating $\boldsymbol{X}_u$ from other $\boldsymbol{X}_i,\forall i\neq u$ [33], its optimal solution can still be obtained according to (51) by modifying $\boldsymbol{\Pi}$ as $\boldsymbol{H}_u^H\left(\frac{1}{\gamma^2}\boldsymbol{\Sigma}+\sum_{l\neq u}\boldsymbol{H}_l\boldsymbol{X}_l\boldsymbol{X}_l^{\mathrm{H}}\boldsymbol{H}_l^{\mathrm{H}}\right)^{-1}\boldsymbol{H}_u$. Moreover, for the MSE minimization problem in the uplink hybrid analog-digital MU-MIMO system, the corresponding element-wise phase derivative also has a similar form to (53) for the single-user counterpart, thereby leading to the optimal solution obtained by (60).

Next, we consider a general WMMSE minimization problem often studied in the downlink hybrid analog-digital MU-MIMO system for both capacity maximization and MSE minimization problems, which can be formulated as [21]

**Prob.10:** $$\min_{\{\boldsymbol{X}_u\}}\sum_{u=1}^{U}\left(\mathrm{Tr}(\boldsymbol{\Phi}_u\boldsymbol{X}_u\boldsymbol{\Pi}_u\boldsymbol{X}_u^{\mathrm{H}})-\mathrm{Tr}(\boldsymbol{B}_u^{\mathrm{H}}\boldsymbol{X}_u)\right.$$
$$-\mathrm{Tr}(\boldsymbol{B}_u\boldsymbol{X}_u^{\mathrm{H}}))$$
$$\text{s.t. }|[\boldsymbol{X}_u]_{i,j}|=1,\ \forall u,i,j, \quad (61)$$

where $\boldsymbol{X}_u$'s are the analog beamformer for the $u$th user. $\boldsymbol{\Phi}_u\in\mathbb{C}^{N_t\times N_t}$'s, $\boldsymbol{\Pi}_u\in\mathbb{C}^{N_{rf}\times N_{rf}}$'s and $\boldsymbol{B}_u\in\mathbb{C}^{N_t\times N_{rf}}$'s denote the corresponding effective channel covariance matrix, the digital beamforming covariance matrix and the cascade channel, respectively, which are mathematically modeled as

$$\boldsymbol{\Phi}_u=\begin{cases}\boldsymbol{H}_{u,u}^{\mathrm{H}}\boldsymbol{G}_{\boldsymbol{\Theta},u}^{\mathrm{H}}\boldsymbol{W}_u\boldsymbol{G}_{\boldsymbol{\Theta},u}\boldsymbol{H}_{u,u},\text{for capacity max problem,}\\ \boldsymbol{H}_{u,u}^{\mathrm{H}}\boldsymbol{G}_{\boldsymbol{\Theta},u}^{\mathrm{H}}\boldsymbol{G}_{\boldsymbol{\Theta},u}\boldsymbol{H}_{u,u}, \text{ for MSE min problem,}\end{cases}$$
$$\boldsymbol{\Pi}_u=\boldsymbol{F}_{\mathrm{D,u}}\boldsymbol{F}_{\mathrm{D,u}}^{\mathrm{H}}, \quad (62)$$
$$\boldsymbol{B}_u=\begin{cases}\boldsymbol{H}_{u,u}^{\mathrm{H}}\boldsymbol{G}_{\boldsymbol{\Theta},u}^{\mathrm{H}}\boldsymbol{W}_u\boldsymbol{F}_{\mathrm{D,u}}^{\mathrm{H}}, \text{ for capacity max problem,}\\ \boldsymbol{H}_{u,u}^{\mathrm{H}}\boldsymbol{G}_{\boldsymbol{\Theta},u}^{\mathrm{H}}\boldsymbol{F}_{\mathrm{D,u}}^{\mathrm{H}}, \text{ for MSE min problem,}\end{cases}\forall u.$$

Generally, the optimal $[\boldsymbol{X}_u]_{i,j}$'s can be obtained when the element-wise phase derivatives of the objective function of **Prob.10** w.r.t. $[\boldsymbol{X}_u]_{i,j}$'s equal zeros, i.e.,

$$\underbrace{\Im\left\{\left([(\boldsymbol{\Phi}\boldsymbol{X}_u\boldsymbol{\Pi})^*]_{i,j}-[\boldsymbol{B}^*]_{i,j}\right)[\boldsymbol{X}_u]_{i,j}\right\}}_{\triangleq g_{i,j}(\boldsymbol{X}_u)}=0,\ \forall u,i,j, \quad (63)$$

where $[(\boldsymbol{\Phi}\boldsymbol{X}_u\boldsymbol{\Pi})]_{i,j}$'s can be rewritten as

$$[(\boldsymbol{\Phi}\boldsymbol{X}_u\boldsymbol{\Pi})]_{i,j}=\mathrm{Tr}\left(\boldsymbol{X}_u[\boldsymbol{\Pi}]_{:,j}[\boldsymbol{\Phi}]_{i,:}\right)$$

$$=[\boldsymbol{X}_u]_{i,j}\left[([\boldsymbol{\Pi}]_{:,j}[\boldsymbol{\Phi}]_{i,:})\right]_{j,i} \tag{64}$$

$$+\underbrace{\sum_{m\neq i}\sum_{n\neq j}[\boldsymbol{X}_u]_{m,n}\left[([\boldsymbol{\Pi}]_{:,j}[\boldsymbol{\Phi}]_{i,:})\right]_{n,m}}_{\triangleq s_{u,i,j}},\ \forall u,i,j.$$

By substituting (64) into (63), we have

$$\Im\left\{\left([[\boldsymbol{\Phi}\boldsymbol{X}_u\boldsymbol{\Pi}]^*]_{i,j}-[\boldsymbol{B}^*]_{i,j}\right)[\boldsymbol{X}_u]_{i,j}\right\} \tag{65a}$$

$$=\Im\left\{\left(s_{i,j}-[\boldsymbol{B}]_{i,j}\right)^*[\boldsymbol{X}_u]_{i,j}+\left[([\boldsymbol{\Pi}]_{:,j}[\boldsymbol{\Phi}]_{i,:})\right]_{j,i}^*\right\}=0,\ \forall u,i,j. \tag{65b}$$

The last term $\left[([\boldsymbol{\Pi}]_{:,j}[\boldsymbol{\Phi}]_{i,:})\right]_{j,i}^*$ is a real scalar since it satisfies $\left[([\boldsymbol{\Pi}]_{:,j}[\boldsymbol{\Phi}]_{i,:})\right]_{j,i}^*=[\boldsymbol{\Pi}]_{j,j}[\boldsymbol{\Phi}]_{i,i}=([\boldsymbol{\Pi}]_{j,j}[\boldsymbol{\Phi}]_{i,i})^*$. Thus, the optimal $[\boldsymbol{\Theta}_u]_{i,j}$'s satisfying (65b) for **Prob.10** are given by

$$[\boldsymbol{\Theta}_u]_{i,j}=\text{Phase}\left\{s_{i,j}-[\boldsymbol{B}]_{i,j}\right\}$$
$$\text{or }\pi+\text{Phase}\left\{s_{i,j}-[\boldsymbol{B}]_{i,j}\right\},\ \forall u,i,j. \tag{66}$$

*2) Fully-Passive IRS-aided MIMO System:* In the fully-passive IRS-aided point-to-point MIMO system as shown in Fig. 2(b), the received signal at the user can be written as

$$\boldsymbol{y}=\boldsymbol{H}\boldsymbol{s}+\boldsymbol{n}=\left(\boldsymbol{H}_0+\boldsymbol{H}_1\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}\boldsymbol{H}_2\right)\boldsymbol{s}+\boldsymbol{n}, \tag{67}$$

where $\boldsymbol{H}_0\in\mathbb{C}^{N_r\times N_t}$, $\boldsymbol{H}_1\in\mathbb{C}^{N_r\times K}$ and $\boldsymbol{H}_2\in\mathbb{C}^{K\times N_t}$ represent the BS-user direct channel, the IRS-user channel and the BS-IRS channel, respectively. $\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}\in\mathbb{C}^{K\times K}$ is the diagonal IRS reflection matrix subject to both diagonal structure constraints and constant modulus constraints [13]. The phase shift vector $\boldsymbol{\theta}$ corresponding to the IRS reflection matrix $\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}$ is then defined as

$$[\boldsymbol{\theta}]_i=\theta_i,\ [\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}]_{i,i}=e^{j\theta_i},\ \forall i. \tag{68}$$

Firstly, we consider the classical capacity maximization problem in the fully-passive IRS-aided MIMO system as follows:

**Prob.11:** $\max_{\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}}\ \log\left|\boldsymbol{\Sigma}^{-1}\boldsymbol{H}\boldsymbol{H}^{\text{H}}+\boldsymbol{I}_{N_r}\right|$

$$\text{s.t. }|[\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}]_{i,i}|=1,\ \forall i. \tag{69}$$

By leveraging the KKT optimality conditions, the element-wise phase derivatives of the objective function of **Prob.11** w.r.t. $[\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}]_{i,i}$'s equal zeros at the optimal $[\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}]_{i,i}$'s, i.e.,

$$\underbrace{\Im\left\{\left[\boldsymbol{H}_2\boldsymbol{H}^{\text{H}}\boldsymbol{\Sigma}^{-1}(\boldsymbol{\Sigma}^{-1}\boldsymbol{H}\boldsymbol{H}^{\text{H}}+\boldsymbol{I}_{N_r})^{-1}\boldsymbol{H}_1\right]_{i,i}[\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}]_{i,i}\right\}}_{\triangleq g_{i,i}(\boldsymbol{\Lambda}_{\boldsymbol{\Theta}})}=0,\ \forall i. \tag{70}$$

The left-hand side of (70) can be further rewritten as

$$\Im\left\{\left[\boldsymbol{H}_2\boldsymbol{H}^{\text{H}}\boldsymbol{\Sigma}^{-1}\left(\boldsymbol{\Sigma}^{-1}\boldsymbol{H}\boldsymbol{H}^{\text{H}}+\boldsymbol{I}_{N_r}\right)^{-1}\boldsymbol{H}_1\right]_{i,i}[\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}]_{i,i}\right\}$$

$$=\Im\left\{\left[\boldsymbol{H}_2\boldsymbol{H}^{\text{H}}\boldsymbol{\Sigma}^{-\frac{1}{2}}\left(\boldsymbol{\Sigma}^{-\frac{1}{2}}\boldsymbol{H}\boldsymbol{H}^{\text{H}}\boldsymbol{\Sigma}^{-\frac{1}{2}}+\boldsymbol{I}_{N_r}\right)^{-1}\right.\right.$$
$$\left.\left.\times\boldsymbol{\Sigma}^{-\frac{1}{2}}\boldsymbol{H}_1\right]_{i,i}[\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}]_{i,i}\right\}$$

$$\overset{(g_1)}{=}\Im\left\{\text{Tr}\left[\boldsymbol{\Gamma}_i\left(\boldsymbol{M}_i^{\text{H}}+e^{-j\theta_i}\boldsymbol{\Gamma}_i^{\text{H}}\right)\left(\boldsymbol{M}_i\boldsymbol{M}_i^{\text{H}}+e^{-j\theta_i}\boldsymbol{M}_i\boldsymbol{\Gamma}_i^{\text{H}}\right.\right.\right.$$
$$\left.\left.\left.+e^{j\theta_i}\boldsymbol{\Gamma}_i\boldsymbol{M}_i^{\text{H}}+\boldsymbol{\Gamma}_i\boldsymbol{\Gamma}_i^{\text{H}}+\boldsymbol{I}_{N_r}\right)^{-1}\right]e^{j\theta_i}\right\}$$

$$\overset{(g_2)}{=}\Im\left\{\text{Tr}\left[e^{j\theta_i}\boldsymbol{\Gamma}_i\boldsymbol{M}_i^{\text{H}}\left(\boldsymbol{\Phi}_i+e^{-j\theta_i}\boldsymbol{M}_i\boldsymbol{\Gamma}_i^{\text{H}}e^{j\theta_i}\boldsymbol{\Gamma}_i\boldsymbol{M}_i^{\text{H}}\right)^{-1}\right]+c_i\right\}$$

$$\overset{(g_3)}{=}\Im\left\{\text{Tr}\left[e^{j\theta_i}\boldsymbol{u}_i\boldsymbol{v}_i^{\text{H}}\left(\boldsymbol{\Psi}_i-\boldsymbol{a}_i\boldsymbol{a}_i^{\text{H}}\right)^{-1}\right]+c_i\right\}$$

$$\overset{(g_4)}{=}\Im\left\{e^{j\theta_i}\boldsymbol{v}_i^{\text{H}}\left(\boldsymbol{\Psi}_i^{-1}+\frac{\boldsymbol{\Psi}_i^{-1}\boldsymbol{a}_i\boldsymbol{a}_i^{\text{H}}\boldsymbol{\Psi}_i^{-1}}{1-\boldsymbol{a}_i^{\text{H}}\boldsymbol{\Psi}_i^{-1}\boldsymbol{a}_i}\right)\boldsymbol{u}_i+c_i\right\},\ \forall i, \tag{71}$$

where $(g_1)$ holds by rewriting $\boldsymbol{\Sigma}^{-\frac{1}{2}}\boldsymbol{H}$ in terms of $[\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}]_{i,i}$ as

$$\boldsymbol{\Sigma}^{-\frac{1}{2}}\boldsymbol{H}=[\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}]_{i,i}\boldsymbol{\Gamma}_i+\underbrace{\boldsymbol{\Sigma}^{-\frac{1}{2}}\boldsymbol{H}_0+\sum_{n\neq i}[\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}]_{n,n}\boldsymbol{\Gamma}_n}_{\triangleq\boldsymbol{M}_i}, \tag{72}$$

where $\boldsymbol{\Gamma}_i=[\boldsymbol{\Sigma}^{-\frac{1}{2}}\boldsymbol{H}_1]_{:,i}[\boldsymbol{H}_2]_{i,:}$. The equality $(g_2)$ in (71) is obtained by defining $\boldsymbol{\Phi}_i=\boldsymbol{M}_i\boldsymbol{M}_i^{\text{H}}+\boldsymbol{\Gamma}_i\boldsymbol{\Gamma}_i^{\text{H}}+\boldsymbol{I}_{N_r}$, $c_i=\text{Tr}\left[\boldsymbol{\Gamma}_i\boldsymbol{\Gamma}_i^{\text{H}}\left(\boldsymbol{\Phi}_i+e^{-j\theta_i}\boldsymbol{M}_i\boldsymbol{\Gamma}_i^{\text{H}}+e^{j\theta_i}\boldsymbol{\Gamma}_i\boldsymbol{M}_i^{\text{H}}\right)^{-1}\right]$, where $c_i$ is a real scalar independent of the optimal $[\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}]_{i,i}$. The equality $(g_3)$ holds based on $\boldsymbol{a}_i=e^{j\theta_i}\boldsymbol{v}_i-\boldsymbol{u}_i$, $\boldsymbol{\Psi}_i=\boldsymbol{\Phi}_i+\boldsymbol{v}_i\boldsymbol{v}_i^{\text{H}}+\boldsymbol{u}_i\boldsymbol{u}_i^{\text{H}}$, where $\boldsymbol{v}_i$ and $\boldsymbol{u}_i$ come from the SVD of the rank-one matrix $\boldsymbol{M}_i\boldsymbol{\Gamma}_i^{\text{H}}$, i.e., $\boldsymbol{M}_i\boldsymbol{\Gamma}_i^{\text{H}}=\boldsymbol{v}_i\boldsymbol{u}_i^{\text{H}}$. The equality $(g_4)$ holds similarly to $(a_2)$ in (14). Then, by substituting (71) into (70), we have

$$\Im\left\{\left[\boldsymbol{H}_2\boldsymbol{H}^{\text{H}}\boldsymbol{\Sigma}^{-1}(\boldsymbol{\Sigma}^{-1}\boldsymbol{H}\boldsymbol{H}^{\text{H}}+\boldsymbol{I}_{N_r})^{-1}\boldsymbol{H}_1\right]_{i,i}[\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}]_{i,i}\right\} \tag{73a}$$

$$=\Im\left\{\frac{A_{2,i}e^{j\theta_i}+C_{2,i}}{D_{2,i}}+c_i\right\}=0,\ \forall i, \tag{73b}$$

where

$$A_{2,i}=\boldsymbol{v}_i^{\text{H}}\boldsymbol{\Psi}_i^{-1}\boldsymbol{u}_i, \tag{74}$$

$$C_{2,i}=\boldsymbol{v}_i^{\text{H}}\boldsymbol{\Psi}_i^{-1}\boldsymbol{u}_i\boldsymbol{u}_i^{\text{H}}\boldsymbol{\Psi}_i^{-1}\boldsymbol{v}_i-\boldsymbol{v}_i^{\text{H}}\boldsymbol{\Psi}_i^{-1}\boldsymbol{v}_i\boldsymbol{u}_i^{\text{H}}\boldsymbol{\Psi}_i^{-1}\boldsymbol{u}_i,$$

$$D_{2,i}=2\text{Re}\left\{\boldsymbol{v}_i^{\text{H}}\boldsymbol{\Psi}_i^{-1}\boldsymbol{u}_ie^{j\theta_i}\right\}+1-\boldsymbol{v}_i^{\text{H}}\boldsymbol{\Psi}_i^{-1}\boldsymbol{v}_i-\boldsymbol{u}_i^{\text{H}}\boldsymbol{\Psi}_i^{-1}\boldsymbol{u}_i,\ \forall i.$$

Similar to (49b), since $C_{2,i}$'s, $D_{2,i}$'s and $c_i$'s are real scalars, the optimal $[\boldsymbol{\theta}]_i$'s satisfying (73b) for **Prob.11** are given by

$$[\boldsymbol{\theta}]_i=\text{Phase}\left\{A_{2,i}^*\right\}\ \text{or }\pi+\text{Phase}\left\{A_{2,i}^*\right\},\ \forall i. \tag{75}$$

Additionally, we consider the MSE minimization problem in the fully-passive IRS-aided point-to-point MIMO system, which is formulated as

**Prob.12:** $\min_{\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}}\ \text{Tr}\left[\left(\boldsymbol{\Sigma}^{-1}\boldsymbol{H}\boldsymbol{H}^{\text{H}}+\boldsymbol{I}_{N_r}\right)^{-1}\right]$

$$\text{s.t. }|[\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}]_{i,i}|=1,\ \forall i. \tag{76}$$

Since the element-wise phase derivatives of the objective function of **Prob.12** w.r.t. $[\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}]_{i,i}$'s equal zeros at the optimal solution, we have

$$\underbrace{\Im\left\{\left[\boldsymbol{H}_2\boldsymbol{H}^{\text{H}}\boldsymbol{\Sigma}^{-1}(\boldsymbol{\Sigma}^{-1}\boldsymbol{H}\boldsymbol{H}^{\text{H}}+\boldsymbol{I}_{N_r})^{-2}\boldsymbol{H}_1\right]_{i,i}[\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}]_{i,i}\right\}}_{\triangleq g_{i,i}(\boldsymbol{\Lambda}_{\boldsymbol{\Theta}})}=0,\ \forall i. \tag{77}$$

The left-hand side of (77) can be further rewritten as

$$\Im\left\{\left[\boldsymbol{H}_2\boldsymbol{H}^{\text{H}}\boldsymbol{\Sigma}^{-1}\left(\boldsymbol{\Sigma}^{-1}\boldsymbol{H}\boldsymbol{H}^{\text{H}}+\boldsymbol{I}_{N_r}\right)^{-2}\boldsymbol{H}_1\right]_{i,i}[\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}]_{i,i}\right\}$$

$$\overset{(h)}{=}\Im\left\{e^{j\theta_i}\boldsymbol{v}_i^{\text{H}}\left(\boldsymbol{\Psi}_i^{-1}+\frac{\boldsymbol{\Psi}_i^{-1}\boldsymbol{a}_i\boldsymbol{a}_i^{\text{H}}\boldsymbol{\Psi}_i^{-1}}{1-\boldsymbol{a}_i^{\text{H}}\boldsymbol{\Psi}_i^{-1}\boldsymbol{a}_i}\right)\right.$$
$$\left.\times\left(\boldsymbol{\Psi}_i^{-1}+\frac{\boldsymbol{\Psi}_i^{-1}\boldsymbol{a}_i\boldsymbol{a}_i^{\text{H}}\boldsymbol{\Psi}_i^{-1}}{1-\boldsymbol{a}_i^{\text{H}}\boldsymbol{\Psi}_i^{-1}\boldsymbol{a}_i}\right)\boldsymbol{u}_i+c_i\right\}$$

$$=\Im\left\{\frac{e^{j\theta_i}\left(d_i\eta_{1,i}+d_i\eta_{2,i}+d_i^2\varepsilon_i+b_{1,i}-r_{0,i}\lambda_i\right)}{\left(2\text{Re}\left\{\lambda_ie^{j\theta_i}\right\}+d\right)^2}\right.$$

$$+ \frac{e^{-\jmath\theta_i}\left((\lambda_i^*)^2\varepsilon_i + b_{2,i} - r_{0,i}\lambda_i^*\right)}{\left(2\mathrm{Re}\left\{\lambda_i e^{\jmath\theta_i}\right\} + d\right)^2} \tag{78}$$

$$+ \left. \frac{2d_i\lambda_i^*\varepsilon_i + \lambda_i^*\eta_{1,i} + \lambda_i^*\eta_{2,i} - b_{0,i} - r_{0,i}d_i}{\left(2\mathrm{Re}\left\{\lambda_i e^{\jmath\theta_i}\right\} + d\right)^2} + c_i\right\}, \ \forall i,$$

where $(h)$ holds similarly to (71) and we have

$$\lambda_i = \boldsymbol{v}_i^{\mathrm{H}}\boldsymbol{\Psi}_i^{-1}\boldsymbol{u}_i, \varepsilon_i = \boldsymbol{v}_i^{\mathrm{H}}\boldsymbol{\Psi}_i^{-2}\boldsymbol{u}_i, d_i = 1 - \boldsymbol{v}_i^{\mathrm{H}}\boldsymbol{\Psi}_i^{-1}\boldsymbol{v}_i - \boldsymbol{u}_i^{\mathrm{H}}\boldsymbol{\Psi}_i^{-1}\boldsymbol{u}_i,$$

$$\eta_{1,i} = \boldsymbol{v}_i^{\mathrm{H}}\boldsymbol{\Psi}_i^{-2}\left(\boldsymbol{u}_i\boldsymbol{u}_i^{\mathrm{H}} + \boldsymbol{v}_i\boldsymbol{v}_i^{\mathrm{H}}\right)\boldsymbol{\Psi}_i^{-1}\boldsymbol{u}_i,$$

$$\eta_{2,i} = \boldsymbol{v}_i^{\mathrm{H}}\boldsymbol{\Psi}_i^{-1}\left(\boldsymbol{u}_i\boldsymbol{u}_i^{\mathrm{H}} + \boldsymbol{v}_i\boldsymbol{v}_i^{\mathrm{H}}\right)\boldsymbol{\Psi}_i^{-2}\boldsymbol{u}_i,$$

$$b_{0,i} = \boldsymbol{v}_i^{\mathrm{H}}\boldsymbol{\Psi}_i^{-1}\left(\boldsymbol{u}_i\boldsymbol{u}_i^{\mathrm{H}} + \boldsymbol{v}_i\boldsymbol{v}_i^{\mathrm{H}}\right)\boldsymbol{\Psi}_i^{-2}\boldsymbol{v}_i\boldsymbol{u}_i^{\mathrm{H}}\boldsymbol{\Psi}_i^{-1}\boldsymbol{u}_i$$

$$+ \boldsymbol{v}_i^{\mathrm{H}}\boldsymbol{\Psi}_i^{-1}\boldsymbol{v}_i\boldsymbol{u}_i^{\mathrm{H}}\boldsymbol{\Psi}_i^{-2}\left(\boldsymbol{u}_i\boldsymbol{u}_i^{\mathrm{H}} + \boldsymbol{v}_i\boldsymbol{v}_i^{\mathrm{H}}\right)\boldsymbol{\Psi}_i^{-1}\boldsymbol{u}_i,$$

$$b_{1,i} = \boldsymbol{v}_i^{\mathrm{H}}\boldsymbol{\Psi}_i^{-1}\left(\boldsymbol{u}_i\boldsymbol{u}_i^{\mathrm{H}} + \boldsymbol{v}_i\boldsymbol{v}_i^{\mathrm{H}}\right)\boldsymbol{\Psi}_i^{-2}\left(\boldsymbol{u}_i\boldsymbol{u}_i^{\mathrm{H}} + \boldsymbol{v}_i\boldsymbol{v}_i^{\mathrm{H}}\right)\boldsymbol{\Psi}_i^{-1}\boldsymbol{u}_i$$

$$+ \boldsymbol{v}_i^{\mathrm{H}}\boldsymbol{\Psi}_i^{-1}\left(\boldsymbol{\Psi}_i^{-1}\boldsymbol{u}_i\boldsymbol{v}_i^{\mathrm{H}}\boldsymbol{\Psi}_i^{-2}\boldsymbol{v}_i^{\mathrm{H}}\boldsymbol{u}_i + \boldsymbol{v}_i\boldsymbol{u}_i^{\mathrm{H}}\boldsymbol{\Psi}_i^{-2}\boldsymbol{u}_i^{\mathrm{H}}\boldsymbol{v}_i\right)\boldsymbol{\Psi}_i^{-1}\boldsymbol{u}_i,$$

$$b_{2,i} = \boldsymbol{v}_i^{\mathrm{H}}\boldsymbol{\Psi}_i^{-1}\boldsymbol{v}_i\boldsymbol{u}_i^{\mathrm{H}}\boldsymbol{\Psi}_i^{-2}\boldsymbol{v}_i\boldsymbol{u}_i^{\mathrm{H}}\boldsymbol{\Psi}_i^{-1}\boldsymbol{u}_i, \tag{79}$$

$$r_{0,i} = \boldsymbol{v}_i^{\mathrm{H}}\boldsymbol{\Psi}_i^{-1}\boldsymbol{v}_i\boldsymbol{u}_i^{\mathrm{H}}\boldsymbol{\Psi}_i^{-2}\boldsymbol{u}_i + \boldsymbol{v}_i^{\mathrm{H}}\boldsymbol{\Psi}_i^{-2}\boldsymbol{v}_i\boldsymbol{u}_i^{\mathrm{H}}\boldsymbol{\Psi}_i^{-1}\boldsymbol{u}_i, \quad \forall i.$$

By substituting (78) into (77), we have

$$\Im\left\{\left[\boldsymbol{H}_2\boldsymbol{H}^{\mathrm{H}}\boldsymbol{\Sigma}^{-1}(\boldsymbol{\Sigma}^{-1}\boldsymbol{H}\boldsymbol{H}^{\mathrm{H}} + \boldsymbol{I}_{N_r})^{-2}\boldsymbol{H}_1\right]_{i,i}[\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}]_{i,i}\right\} \tag{80a}$$

$$= \Im\left\{\frac{A_{3,i}e^{\jmath\theta_i} + B_{3,i}e^{-\jmath\theta_i} + C_{3,i}}{D_{3,i}} + c_i\right\} = 0, \ \forall i, \tag{80b}$$

where

$$A_{3,i} = d_i\eta_{1,i} + d_i\eta_{2,i} + d_i^2\varepsilon_i + b_{1,i} - r_{0,i}\lambda_i,$$

$$B_{3,i} = (\lambda_i^*)^2\varepsilon_i + b_{2,i} - r_{0,i}\lambda_i^*,$$

$$C_{3,i} = 2d_i\lambda_i^*\varepsilon_i + \lambda_i^*\eta_{1,i} + \lambda_i^*\eta_{2,i} - b_{0,i} - r_{0,i}d_i,$$

$$D_{3,i} = \left(2\mathrm{Re}\left\{\lambda_i e^{\jmath\theta_i}\right\} + d_i\right)^2, \ \forall i. \tag{81}$$

It is noted that (80b) has the same form as (57b) for **Prob.9**. Thus, the optimal $[\boldsymbol{\theta}]_i$'s to **Prob.12** can be obtained similarly.

Similarly, the capacity maximization and MSE minimization problems in the uplink fully-passive IRS-aided MU-MIMO system have the same forms as these in the single-user case, thus their corresponding optimal solutions can be attained directly. Moreover, it is essential to investigate the general WMMSE minimization problem mostly considered in the downlink fully-passive IRS-aided MU-MIMO system, which is formulated as

**Prob.13:** $\min_{\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}} \mathrm{Tr}(\boldsymbol{\Phi}\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}\boldsymbol{\Pi}\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}^{\mathrm{H}}) - \mathrm{Tr}(\boldsymbol{B}^{\mathrm{H}}\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}) - \mathrm{Tr}(\boldsymbol{B}\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}^{\mathrm{H}})$

$$\text{s.t. } |[\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}]_{i,i}| = 1, \ \forall i. \tag{82}$$

Since the element-wise phase derivatives of the objective function w.r.t. $[\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}]_{i,i}$'s are zeros at the optimal solution, i.e.,

$$\underbrace{\Im\left\{\left[(\boldsymbol{\Phi}\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}\boldsymbol{\Pi})^*\right]_{i,i}[\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}]_{i,i}\right\} - \Im\left\{[\boldsymbol{B}^*]_{i,i}[\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}]_{i,i}\right\}}_{\triangleq g_{i,i}(\boldsymbol{\Lambda}_{\boldsymbol{\Theta}})} = 0, \ \forall i, \tag{83}$$

referring to (54), the optimal solution to **Prob.13** can be easily derived as

$$[\boldsymbol{\theta}]_i = \mathrm{Phase}\left\{s_i - [\boldsymbol{B}]_{i,i}\right\}$$

$$\text{or } \pi + \mathrm{Phase}\left\{s_i - [\boldsymbol{B}]_{i,i}\right\}, \ \forall i, \tag{84}$$

where $s_i = \sum_{m \neq i}[\boldsymbol{\Lambda}_{\boldsymbol{\Theta}}]_{m,m}\left[([\boldsymbol{\Pi}]_{:,i}[\boldsymbol{\Phi}]_{i,:})\right]_{m,m}$.

---

**Algorithm 1** A Novel AO Algorithm for Solving Problems under Constant Modulus Constraints

---

**Initialize:** Arbitrary five feasible solutions $\widehat{\boldsymbol{X}}_m^{(0)}$, $m = 1, \cdots, 5$; iteration index $t = 0$; convergence threshold $\epsilon$.

1: **repeat**
2:     **for** $i = 1$ to $N_t$, $j = 1$ to $N_{rf}$ **do**
3:         Calculate $g_{i,j}(\widehat{\boldsymbol{X}}_m^{(t)})$, $m = 1, \cdots, 5$ and the auxiliary vector $\boldsymbol{w}_{i,j}^{(t)}$ as in (87).
4:         Update $\left[\boldsymbol{\Theta}^{(t)}\right]_{i,j}$ as in (89) and obtain $[\boldsymbol{X}^{(t)}]_{i,j}$ according to (38).
5:     **end for**
6:     Update $\widehat{\boldsymbol{X}}_m^{(t+1)} = \boldsymbol{X}^{(t)}$ for an arbitrary $m \in [1, 5]$.
7:     $t = t + 1$.
8: **until** The increment/decrement of the objective function value between two consecutive iterations is less than $\epsilon$.
9: **return** $\boldsymbol{X}$.

---

In a nutshell, a series of optimization problems in the wireless systems associated with constant modulus constraints are investigated in this subsection, whose optimal solutions are available using the proposed element-wise phase derivatives.

### C. A Novel AO Algorithm

It follows from Sec. III-B that these element-wise phase derivatives associated with **Prob.8**~**Prob.13** can be mainly classified into two forms. Moreover, each form always has two zero-derivative points. The optimal one in these two zero-derivative points can be determined by an elegant function, as summarized in the following proposition.

**Proposition 1.** *For different types of objective functions, the element-wise phase derivatives under constant modulus constraints can be mainly summarized as the following two general forms, i.e., the linear form and the conjugate linear form, which are shown as*

$$g_{i,j}(\boldsymbol{X}) = \begin{cases} \Im\left\{A_{i,j}^{\mathrm{L}}[\boldsymbol{X}]_{i,j} + C_{i,j}^{\mathrm{L}}\right\}, \forall i, j, \textit{for trace-linear,} \\ \qquad \textit{trace-quadratic, log-determinant functions,} \tag{85a} \\ \Im\left\{A_{i,j}^{\mathrm{CL}}[\boldsymbol{X}]_{i,j} + B_{i,j}^{\mathrm{CL}}[\boldsymbol{X}]_{i,j}^* + C_{i,j}^{\mathrm{CL}}\right\}, \forall i, j, \\ \qquad\qquad\qquad \textit{for trace-inverse function,} \tag{85b} \end{cases}$$

*where $A_{i,j}^{\mathrm{L}}$'s, $A_{i,j}^{\mathrm{CL}}$'s, $B_{i,j}^{\mathrm{CL}}$'s and $C_{i,j}^{\mathrm{CL}}$'s are all complex scalars and $C_{i,j}^{\mathrm{L}}$'s are real scalars. In particular, the linear element-wise phase derivative in (85a) can be regarded as a simplified case of its conjugate linear counterpart by setting $B_{i,j}^{\mathrm{CL}} = 0$ and $\Im\left\{C_{i,j}^{\mathrm{CL}}\right\} = 0, \forall i, j$ in (85b). Moreover, there are two points $\left\{[\boldsymbol{X}_1]_{i,j}, [\boldsymbol{X}_2]_{i,j}\right\}$ satisfying $g_{i,j}(\boldsymbol{X}) = 0, \forall i, j$, from which the optimal solution can be determined by the following conjugate linear function*

$$[\boldsymbol{X}_{\mathrm{opt}}]_{i,j} = \arg_{\left\{[\boldsymbol{X}_1]_{i,j}, [\boldsymbol{X}_2]_{i,j}\right\}} \tag{86}$$

$$\begin{cases} \Re\left\{A_{i,j}^{\mathrm{CL}}[\boldsymbol{X}]_{i,j} - B_{i,j}^{\mathrm{CL}}[\boldsymbol{X}]_{i,j}^*\right\} \geq 0, \ \textit{for min problem,} \\ \Re\left\{A_{i,j}^{\mathrm{CL}}[\boldsymbol{X}]_{i,j} - B_{i,j}^{\mathrm{CL}}[\boldsymbol{X}]_{i,j}^*\right\} < 0, \ \textit{for max problem,} \end{cases} \forall i, j.$$

*Proof.* The detailed proof is shown in Appendix A.

Based on **proposition 1**, it is seen that the optimal $[\boldsymbol{X}]_{i,j}$'s are obtained by aligning their phase-shifts with the corresponding counterparts jointly determined by $A_{i,j}^{\mathrm{CL}}$'s, $B_{i,j}^{\mathrm{CL}}$'s and $C_{i,j}^{\mathrm{CL}}$'s, which are all related to $[\boldsymbol{X}]_{m,n}$'s, $m \neq i, n \neq j$

and need to be frequently calculated in each iteration of updating $[\boldsymbol{X}]_{i,j}$'s. In addition, the calculations of $A_{i,j}^{\text{CL}}$'s, $B_{i,j}^{\text{CL}}$'s and $C_{i,j}^{\text{CL}}$'s all involve complicated matrix inversion with complexity of $\mathcal{O}\left(N_r^3\right)$ and SVD with complexity of $\mathcal{O}\left(2N_rN_t^2+N_r^3\right)$. It is evident that this complexity will become enormous as $N_r$ and $N_t$ increases. In order to avoid the frequent calculations of $A_{i,j}^{\text{CL}}$'s, $B_{i,j}^{\text{CL}}$'s and $C_{i,j}^{\text{CL}}$'s, we next derive the optimal solution directly based on the functions $g_{i,j}(\boldsymbol{X})$'s associated with the original element-wise phase derivatives, which is shown in **Proposition 2**.

**Proposition 2.** *Define arbitrary five feasible solutions satisfying constant modulus constraints, i.e., $\widehat{\boldsymbol{X}}_m$, $m=1,\cdots,5$ and calculate their corresponding $g_{i,j}(\widehat{\boldsymbol{X}}_m)$'s, we have*

$$\boldsymbol{w}_{i,j}=\boldsymbol{R}_{i,j}^{-1}\boldsymbol{t}_{i,j} \text{ with } \boldsymbol{R}_{i,j}\in\mathbb{C}^{5\times5}, \boldsymbol{t}_{i,j}\in\mathbb{C}^{5\times1}, \ \forall i,j, \quad (87)$$

*where*

$$[\boldsymbol{R}_{i,j}]_{m,:} = \Big[ \Im\{[\widehat{\boldsymbol{X}}_m]_{i,j}\}-[\boldsymbol{t}_{i,j}]_m\Re\{[\widehat{\boldsymbol{X}}_m]_{i,j}\}, \Re\{[\widehat{\boldsymbol{X}}_m]_{i,j}\}$$
$$+[\boldsymbol{t}_{i,j}]_m\Im\{[\widehat{\boldsymbol{X}}_m]_{i,j}\}, -\Im\{[\widehat{\boldsymbol{X}}_m]_{i,j}\}-[\boldsymbol{t}_{i,j}]_m\Re\{[\widehat{\boldsymbol{X}}_m]_{i,j}\},$$
$$\Re\{[\widehat{\boldsymbol{X}}_m]_{i,j}\}-[\boldsymbol{t}_{i,j}]_m\Im\{[\widehat{\boldsymbol{X}}_m]_{i,j}\}, 1 \Big],$$

$$[\boldsymbol{t}_{i,j}]_m = \tan(\angle g_{i,j}(\widehat{\boldsymbol{X}}_m)), \ \forall i,j,m. \quad (88)$$

*Then, we obtain the optimal solutions of optimization problems with conjugate linear element-wise phase derivatives as*

$$[\boldsymbol{\Theta}]_{i,j}=\begin{cases} -\arctan\left(\frac{\widehat{z}_{2,i,j}}{\widehat{z}_{1,i,j}}\right)+\arcsin\left(-\frac{[\boldsymbol{w}_{i,j}]_5}{\sqrt{\widehat{z}_{1,i,j}^2+\widehat{z}_{2,i,j}^2}}\right) \\ or \ \pi-\arctan\left(\frac{\widehat{z}_{2,i,j}}{\widehat{z}_{1,i,j}}\right)+\arcsin\left(-\frac{[\boldsymbol{w}_{i,j}]_5}{\sqrt{\widehat{z}_{1,i,j}^2+\widehat{z}_{2,i,j}^2}}\right), \\ \qquad\qquad\qquad if \ \widehat{z}_{1,i,j}\geq 0, \ \forall i,j, \\ -\pi+\arctan\left(\frac{\widehat{z}_{2,i,j}}{\widehat{z}_{1,i,j}}\right)+\arcsin\left(-\frac{[\boldsymbol{w}_{i,j}]_5}{\sqrt{\widehat{z}_{1,i,j}^2+\widehat{z}_{2,i,j}^2}}\right) \\ or \ \arctan\left(\frac{\widehat{z}_{2,i,j}}{\widehat{z}_{1,i,j}}\right)+\arcsin\left(-\frac{[\boldsymbol{w}_{i,j}]_5}{\sqrt{\widehat{z}_{1,i,j}^2+\widehat{z}_{2,i,j}^2}}\right), \\ \qquad\qquad\qquad if \ \widehat{z}_{1,i,j}<0, \ \forall i,j, \end{cases} \quad (89)$$

*where $\widehat{z}_{1,i,j} = [\boldsymbol{w}_{i,j}]_1 - [\boldsymbol{w}_{i,j}]_3$ and $\widehat{z}_{2,i,j} = [\boldsymbol{w}_{i,j}]_2 + [\boldsymbol{w}_{i,j}]_4, \forall i,j$. In particular, for the special case of the linear element-wise phase derivatives, we have $\boldsymbol{w}_{i,j} = [[\boldsymbol{w}_{i,j}]_1, [\boldsymbol{w}_{i,j}]_2, 0, 0, 0]^{\text{T}}$ and $\widehat{z}_{1,i,j}$ does not affect the optimal $[\boldsymbol{\Theta}]_{i,j}$.*

*Proof.* The detailed proof is shown in Appendix B.

Based on **Proposition 2**, we next aim to develop a novel AO algorithm with the aid of five arbitrary feasible solutions to determine the optimal solutions of all above optimization problems under constant modulus constraints, which is summarized in **Algorithm 1**.

*Remark 3:* In fact, the element-wise phase derivatives for different optimization problems with constant modulus constraints can be roughly summarized as a general (conjugate) linear form, and the corresponding optimal solutions can be determined from two potential solutions based on a concise conjugate linear function. Based on this, a novel AO algorithm with the advantages of the low complexity and guaranteed performance is developed. Moreover, for the optimization problems that have other nonconvex constraints such as the

**TABLE III**
COMPUTATIONAL COMPLEXITIES OF ALL STUDIED ALGORITHMS

| Algorithm | Computational complexity |
|---|---|
| Proposed novel AO algorithm | $\mathcal{O}\left(KN_r^3\right)$ |
| Element-wise BCD algorithm | $\mathcal{O}\left(2KN_r^2N_t+3KN_r^3\right)$ |
| MM-based algorithm | $\mathcal{O}\left(T_{\text{MM}}K^2\right)$ |
| RCG-based algorithm | $\mathcal{O}\left(K^2+T_{\text{RCG}}\left(K^3\right)\right)$ |

fractional functions (e.g., signal-to-interference-plus-noise ratio) [34] or the difference of convex functions (e.g., secrecy rate) [35], the proposed novel AO algorithm can also be used to derive the closed-form solutions. Specifically, we simultaneously employ the successive convex approximation technique to handle the nonconvex constraints and the Lagrange multiplier method to reformulate the objective function. The corresponding Lagrangian belongs to one of the four types of the considered functions in Sec. II-A and thus can be efficiently solved by the proposed algorithms. However, dual variables are introduced and need to be further determined by the subgradient method.

### D. Convergence and Complexity Analysis of the Novel AO Algorithm

In this subsection, we firstly demonstrate the convergence of the proposed novel AO algorithm to a locally optimal solution. By leveraging the element-wise phase derivative, the proposed algorithm firstly decouples the original nonconvex problem under constant modulus constraints into multiple unconstrained subproblems associated with each matrix element, and then derives their optimal closed-form solutions since it satisfies the first-order KKT optimality condition in each subproblem. Thus, the objective function is monotonically non-increasing over the iterations. Moreover, the objective function is lower bounded since the feasible region of the original problem is closed. Therefore, the proposed algorithm is guaranteed to converge to a stationary solution for the original problem. That is to say, its local optimality is thus ensured [36].

Then, taking the capacity maximization problem in the fully-passive IRS-aided MIMO system as an example, we analyze and compare the computational complexities of the proposed novel AO algorithm and the advanced algorithms for constant modulus constrained problem, which is summarized in Table III. It is observed that the computational complexity of the proposed novel AO algorithm is dominated by the matrix inversion for calculating $g_{i,j}(\widehat{\boldsymbol{X}}_m), \forall m,i,j$, which is computed as $\mathcal{O}\left(KN_r^3\right)$.

Currently, three effective algorithms are widely utilized to optimize the constant modulus constrained problems, i.e., the element-wise BCD algorithm [18], the MM-based algorithm [21] and the Riemannian conjugate gradient (RCG)-based algorithm [38]. Similarly, the computational complexities of these advanced algorithms can be elaborated as below. Specifically, the element-wise BCD algorithm has a high complexity of $\mathcal{O}\left(2KN_r^2N_t+3KN_r^3\right)$, which mainly comes from operations of matrix inversion and SVD. Moreover, the computational complexity of the MM-based algorithm is given by $\mathcal{O}\left(T_{\text{MM}}K^2\right)$ with $T_{\text{MM}}$ denoting the number of iterations. The RCG-based algorithm primarily comprises three components, i.e., the computation of Riemannian gradient, the retraction operator and the Armijo backtracking line search,
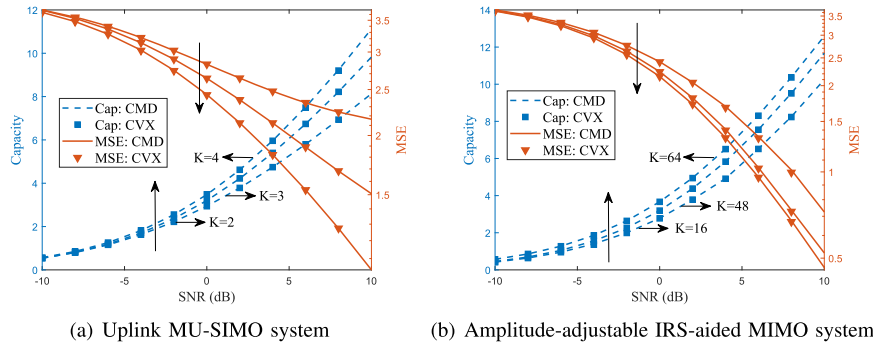
(a) Uplink MU-SIMO system        (b) Amplitude-adjustable IRS-aided MIMO system

Fig. 3. The capacity and MSE performance comparison in the two systems considered in Sec. II.

whose complexity is calculated as $\mathcal{O}\left(K^2 + T_{\mathrm{RCG}}\left(K^3\right)\right)$ with $T_{\mathrm{RCG}}$ being the line search times. In conclusion, since $K \gg \{N_r, N_t\}$, the proposed novel AO algorithm achieves a significant reduction in the computational complexity as compared to the advanced algorithms. Therefore, it is more suitable for practical applications.

## IV. SIMULATIONS AND DISCUSSIONS

In this section, numerical simulation results are provided to evaluate the performance of the derived optimal closed-form solutions based on complex matrix derivatives in Sec. II-B (also referred to as CMD-based algorithm), and the novel AO algorithm in **Algorithm 1**, which are respectively proposed for tackling the optimization problems under diagonal structure constraints and constant modulus constraints.

### A. Diagonal Structure Constraints

We firstly consider the uplink MU-SIMO system, where $K = 2, 3$ and $4$ single-antenna users transmit signals to the BS equipped with $N_t = 6$ antennas, respectively. Moreover, the maximum transmit power of each user is assumed to be $P_k = 28$ dBm, $26$ dBm and $25$ dBm for $K = 2, 3$ and $4$ users, respectively, and the maximum sum power is set as $P = 30$ dBm. Under the assumption that the channel follows the circularly symmetric complex Gaussian distribution with unit noise variance, i.e. $\boldsymbol{H} \sim \mathcal{CN}\left(\boldsymbol{0}, \boldsymbol{I}_{N_t K}\right)$, the SNR is defined as SNR $= 10 \log_{10}(\frac{P}{\sigma^2})$, where noise power $\sigma^2$ varies with SNR. All the results are obtained by averaging over 100 channel realizations. Firstly, Fig. 3(a) compares the capacity and MSE performance achieved by the CMD-based algorithm and the CVX toolbox [37] versus SNR under different numbers of users. For each considered number of users, it is clearly observed that the CMD-based algorithm achieves almost the same capacity and MSE performance as the numerical CVX optimization with a lower complexity. Specifically, the complexity of the CMD-based algorithm is dominated by matrix inversion with the complexity of $\mathcal{O}\left(K^3\right)$, while the numerical CVX optimization adopting the interior-point method has a complexity of $\mathcal{O}\left((K + N_t K)^{3.5} + K^3\right)$. Moreover, as expected that with the expansion of the number of users, the capacity and MSE performance are further enhanced due to the increasing system degrees of freedom. Thus, the global optimality and low complexity of the CMD-based algorithm are demonstrated.

Then, the amplitude-adjustable IRS-aided MIMO system is taken into account, where the BS equipped with $N_t = 6$ antennas and $N_{rf} = 4$ RF chains communicates with the user equipped with $N_r = 4$ antennas, while $K = 16, 48$ and $64$-element IRSs are deployed to enhance the point-to-point communication, respectively. The path loss setting is the same as that in [17] and other parameter settings are the same as those in the uplink MU-SIMO system. Fig. 3(b) illustrates that the capacity and MSE performance attained by the CMD-based algorithm and the numerical CVX optimization as the function of SNR. We find that the CMD-based algorithm attains almost the same optimal performance as the numerical CVX toolbox. Also, the increase in the number of IRS elements results in the improved capacity and MSE performance.

### B. Constant Modulus Constraints

In the point-to-point hybrid analog-digital MIMO system, the BS equipped with $N_t = 6$ antennas and $N_{rf} = 4$ RF chains serves a user equipped with $N_r = 4$ antennas. Other parameter settings are the same as those of the uplink MU-SIMO system. Moreover, we compare the proposed novel AO algorithm with the following benchmark schemes: **BCD**: The authors of [20] propose an element-wise BCD algorithm for optimizing the analog beamforming matrix. **MM**: The authors of [21] equivalently reformulate the original problem into the WMMSE minimization problem, and then optimize the analog beamforming matrix using the MM-based algorithm. **RCG**: The RCG-based algorithm is applied to effectively solve the corresponding optimization problem [38].

Taking the capacity maximization problem as an example, we firstly depict Fig. 4(a) to show the convergence behavior of all studied algorithms, where SNR$= -2$ dB. As can be seen, the proposed novel AO algorithm obviously outperforms the MM-based and RCG-based algorithms in terms of the convergence speed and capacity performance. Moreover, both the proposed algorithm and the element-wise BCD algorithm converge within 20 iterations, while the proposed algorithm shows a little higher capacity and a faster speed than the element-wise BCD algorithm.

Next, we demonstrate the effectiveness of the proposed novel AO algorithm in Fig. 4(b), where three sets of different feasible solutions satisfying constant modulus constraints are generated randomly to calculate the optimal solutions of the capacity maximization and MSE minimization problems for different numbers of BS antennas, respectively. The same capacity and MSE performance are attained for different feasible solutions under each antenna number setting and these two performance are enhanced with the increasing number of BS antennas, implying the stability and effectiveness of the proposed novel AO algorithm.

Fig. 4(c) compares the capacity and MSE performance of the proposed novel AO algorithm and benchmark schemes
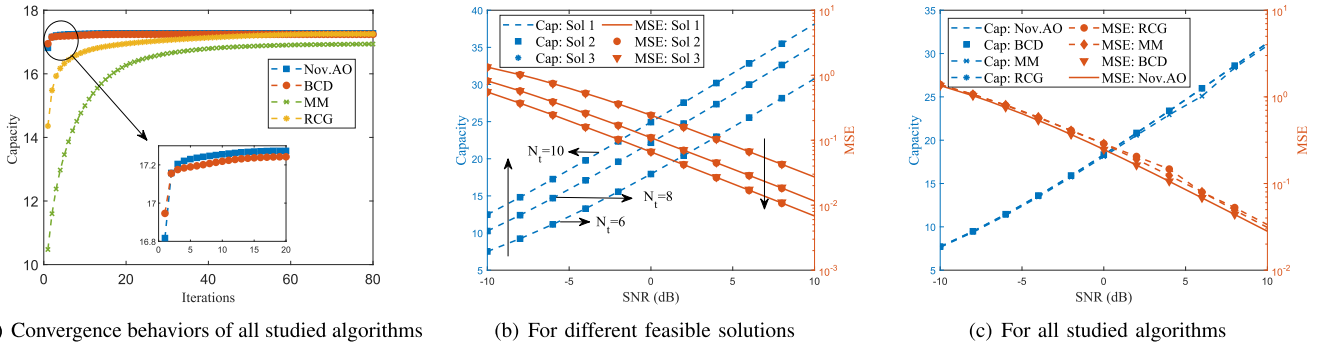
(a) Convergence behaviors of all studied algorithms

(b) For different feasible solutions

(c) For all studied algorithms

Fig. 4. The capacity and MSE performance comparison in the hybrid analog-digital MIMO system.



(a) For all studied algorithms

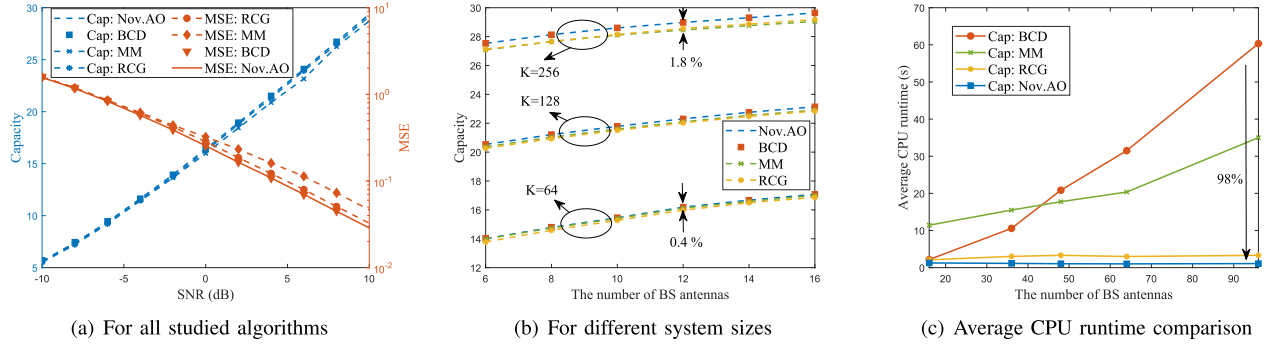(b) For different system sizes

(c) Average CPU runtime comparison

Fig. 5. The capacity and MSE performance comparison in the fully-passive IRS-aided MIMO system.

versus SNR. Obviously, the proposed novel AO algorithm consistently performs better than the MM-based and RCG-based algorithms, since both benchmark schemes do not directly act on the original objective function. To be specific, the MM-based and RCG-based algorithms respectively adopt the surrogate function and the gradient projection to tackle constant modulus constraints, respectively, whose performance depends on the choice of the initial point and the descending step size. Moreover, the proposed novel AO algorithm attains nearly identical performance to the element-wise BCD algorithm, since both the two algorithms derive the optimal closed-form solution associated with each matrix element.

Furthermore, we consider the fully-passive IRS-aided MIMO system with the same parameter settings as the amplitude-adjustable IRS-aided MIMO system. In Fig. 5(a), the capacity and MSE performance of the proposed novel AO algorithm are respectively compared with that of benchmark schemes in the hybrid analog-digital MIMO system. Notice that the above three benchmark schemes, i.e., the element-wise BCD algorithm [18], [19], the MM-based algorithm and the RCG-based algorithm, have also been widely adopted in most existing IRS related works and thus are still considered here. The observed trend in Fig. 5(a) indicates that the proposed novel AO algorithm is able to achieve better performance than the MM-based and RCG-based algorithms, and almost the same performance as the element-wise BCD algorithm while maintaining a lower complexity.

Then, taking the capability maximization problem as an example, we compare the proposed novel AO algorithm and benchmark schemes versus the numbers of BS antennas and IRS elements in Fig. 5(b), where SNR= $-2$ dB. It is naturally seen that the capacity performance becomes better with the increasing sizes of BS antennas and IRS elements. Moreover, the performance gap between the proposed novel AO algorithm and benchmark schemes becomes more significant with the increase of the system size. For example, for $N_t = 12$, when the number of IRS elements increases from $K = 64$ to $K = 256$, there is $1.4\%$ capacity increment achieved by the proposed algorithm relative to the MM-based algorithm. Moreover, the proposed algorithm always achieves almost the same performance as the element-wise BCD algorithm while maintaining a lower complexity. Therefore, the proposed algorithm's scalability and optimality are ensured.

Finally, in order to demonstrate the low-complexity advantage of the proposed AO algorithm, we depict Fig. 5(c) to intuitively compare its average CPU runtime with that of benchmark schemes for solving the capacity maximization problem, where SNR= 5 dB and $N_r = 8$. We firstly observe that the proposed novel AO algorithm has the lowest average CPU runtime, since it avoids the complicated matrix operation. Moreover, it shows a sharp decrease of the average CPU runtime relative to the element-wise BCD algorithm. For example, at $N_t = 96$, there is $98\%$ CPU runtime decrement achieved by the proposed algorithm implying its low-complexity advantage. Whereas, the element-wise BCD algorithm suffers from the highest time overhead which also increases as the number of BS antennas increases, since the dimension of the involved matrix inversion is the same as the number of BS antennas.

## V. Conclusions

In this paper, we investigated complex matrix derivatives for two special matrices, i.e., diagonal structured matrices and constant modulus structured matrices. Under the diagonal structure constraints, the optimal closed-form solutions of the capacity maximization problem, the MSE minimization problem and their variants can be obtained using complex matrix derivatives. Whereas for constant modulus constraints, the optimal solutions of these classical optimization problems are derived utilizing element-wise phase derivatives. Further, in

$$\tan(\angle g_{i,j}(\widehat{\boldsymbol{X}}_m)) = \frac{\Im\{g_{i,j}(\widehat{\boldsymbol{X}}_m)\}}{\Re\{g_{i,j}(\widehat{\boldsymbol{X}}_m)\}}$$

$$=\frac{\Re\{A_{i,j}^{\mathrm{CL}}\}\Im\{\widehat{\boldsymbol{X}}_m\} + \Im\{A_{i,j}^{\mathrm{CL}}\}\Re\{\widehat{\boldsymbol{X}}_m\} + \Im\{B_{i,j}^{\mathrm{CL}}\}\Re\{\widehat{\boldsymbol{X}}_m\} - \Re\{B_{i,j}^{\mathrm{CL}}\}\Im\{\widehat{\boldsymbol{X}}_m\} + \Im\{C_{i,j}^{\mathrm{CL}}\}}{\Re\{A_{i,j}^{\mathrm{CL}}\}\Re\{\widehat{\boldsymbol{X}}_m\} - \Im\{A_{i,j}^{\mathrm{CL}}\}\Im\{\widehat{\boldsymbol{X}}_m\} + \Re\{B_{i,j}^{\mathrm{CL}}\}\Re\{\widehat{\boldsymbol{X}}_m\} + \Im\{B_{i,j}^{\mathrm{CL}}\}\Im\{\widehat{\boldsymbol{X}}_m\} + \Re\{C_{i,j}^{\mathrm{CL}}\}}$$

$$=\frac{[\boldsymbol{w}_{i,j}]_1\Im\{\widehat{\boldsymbol{X}}_m\} + [\boldsymbol{w}_{i,j}]_2\Re\{\widehat{\boldsymbol{X}}_m\} + [\boldsymbol{w}_{i,j}]_4\Re\{\widehat{\boldsymbol{X}}_m\} - [\boldsymbol{w}_{i,j}]_3\Im\{\widehat{\boldsymbol{X}}_m\} + [\boldsymbol{w}_{i,j}]_5}{[\boldsymbol{w}_{i,j}]_1\Re\{\widehat{\boldsymbol{X}}_m\} - [\boldsymbol{w}_{i,j}]_2\Im\{\widehat{\boldsymbol{X}}_m\} + [\boldsymbol{w}_{i,j}]_3\Re\{\widehat{\boldsymbol{X}}_m\} + [\boldsymbol{w}_{i,j}]_4\Im\{\widehat{\boldsymbol{X}}_m\} + 1}, \ m = 1,\cdots,5. \tag{90}$$

order to avoid the complicated matrix operations, we explore the inherent structure of the element-wise phase derivatives, and develop a novel AO algorithm with the aid of several arbitrary feasible solutions. Finally, numerical simulations demonstrate the global optimality and low complexity of the proposed novel AO algorithm.

## APPENDIX

### A. Proof of Proposition 1

Based on Sec. III-B, we can easily conclude that for both capacity maximization problems (i.e. **Prob.8**, **Prob.11**) with log-determinant functions and WMMSE minimization problems (i.e. **Prob.10**, **Prob.13**) with trace-linear and trace-quadratic functions, the corresponding element-wise phase derivatives have the same linear forms as in (85a). Whereas, for the MSE minimization problems (i.e. **Prob.9**, **Prob.12**) with trace-inverse functions, the element-wise phase derivatives satisfy the conjugate linear forms in (85b).

Furthermore, in order to determine the optimal solution from two zero-derivative points $\left\{[\boldsymbol{X}_1]_{i,j}, [\boldsymbol{X}_2]_{i,j}\right\}$ satisfying $g_{i,j}(\boldsymbol{X}) = 0, \forall i,j$, we resort to the second-order derivative of the corresponding objective function $f(\boldsymbol{X})$ w.r.t. $[\boldsymbol{\Theta}]_{i,j}$'s. Specifically, by taking the conjugate linear element-wise phase derivatives as an example, we have

$$\frac{\partial^2 f(\boldsymbol{X})}{\partial [\boldsymbol{\Theta}]_{i,j}^2} = \frac{\partial g_{i,j}(\boldsymbol{X})}{\partial [\boldsymbol{\Theta}]_{i,j}} = 2j\Re\left\{A_{i,j}^{\mathrm{CL}}[\boldsymbol{X}]_{i,j} - B_{i,j}^{\mathrm{CL}}[\boldsymbol{X}]_{i,j}^*\right\}. \tag{91}$$

According to the optimization theory, it is readily inferred that these two zero-derivative points are local minimum when $\frac{\partial^2 f(\boldsymbol{X})}{\partial [\boldsymbol{\Theta}]_{i,j}^2} \geq 0$ holds. In contrast, they become local maximum when $\frac{\partial^2 f(\boldsymbol{X})}{\partial [\boldsymbol{\Theta}]_{i,j}^2} < 0$ holds [25]. This completes the proof.

### B. Proof of Proposition 2

It follows from **Proposition 1** that the derivation of the optimal $\boldsymbol{X}_{i,j}$'s is based on $A_{i,j}^{\mathrm{CL}}$'s, $B_{i,j}^{\mathrm{CL}}$'s and $C_{i,j}^{\mathrm{CL}}$'s. Thus, we firstly define a 5-dimensional vector as follows.

$$\boldsymbol{w}_{i,j} = [[\boldsymbol{w}_{i,j}]_1, [\boldsymbol{w}_{i,j}]_2, \cdots, [\boldsymbol{w}_{i,j}]_5]^{\mathrm{T}} \tag{92}$$

$$=\left[\frac{\Re\{A_{i,j}^{\mathrm{CL}}\}}{\Re\{C_{i,j}^{\mathrm{CL}}\}}, \frac{\Im\{A_{i,j}^{\mathrm{CL}}\}}{\Re\{C_{i,j}^{\mathrm{CL}}\}}, \frac{\Re\{B_{i,j}^{\mathrm{CL}}\}}{\Re\{C_{i,j}^{\mathrm{CL}}\}}, \frac{\Im\{B_{i,j}^{\mathrm{CL}}\}}{\Re\{C_{i,j}^{\mathrm{CL}}\}}, \frac{\Im\{C_{i,j}^{\mathrm{CL}}\}}{\Re\{C_{i,j}^{\mathrm{CL}}\}}\right]^{\mathrm{T}}, \forall i,j.$$

Moreover, by recalling **Prob.12**, $A_{i,j}^{\mathrm{CL}}$'s, $B_{i,j}^{\mathrm{CL}}$'s and $C_{i,j}^{\mathrm{CL}}$'s are closely related to the functions $g_{i,j}(\boldsymbol{X})$'s associated with the original element-wise phase derivatives as (90), as shown at the top of this page. Note that the equations in (90) can be further rewritten as the following homogeneous linear equations, i.e.,

$$[\boldsymbol{w}_{i,j}]_1\left(\Im\{\widehat{\boldsymbol{X}}_m\} - \tan(\angle g_{i,j}(\widehat{\boldsymbol{X}}_m))\Re\{\widehat{\boldsymbol{X}}_m\}\right)$$

$$+[\boldsymbol{w}_{i,j}]_2\left(\Re\{\widehat{\boldsymbol{X}}_m\} + \tan(\angle g_{i,j}(\widehat{\boldsymbol{X}}_m))\Im\{\widehat{\boldsymbol{X}}_m\}\right)$$

$$-[\boldsymbol{w}_{i,j}]_3\left(\Im\{\widehat{\boldsymbol{X}}_m\} + \tan(\angle g_{i,j}(\widehat{\boldsymbol{X}}_m))\Re\{\widehat{\boldsymbol{X}}_m\}\right)$$

$$+[\boldsymbol{w}_{i,j}]_4\left(\Re\{\widehat{\boldsymbol{X}}_m\} - \tan(\angle g_{i,j}(\widehat{\boldsymbol{X}}_m))\Im\{\widehat{\boldsymbol{X}}_m\}\right)$$

$$+[\boldsymbol{w}_{i,j}]_5 - \tan(\angle g_{i,j}(\widehat{\boldsymbol{X}}_m)) = 0, \ m = 1,\cdots,5. \tag{93}$$

It follows from (93) that the optimal $\boldsymbol{w}_{i,j}$'s can be obtained by jointly solving its involved five homogeneous linear equations, whose closed-form structures are further shown in (87) and the proof of **Proposition 2** is completed.

## REFERENCES

[1] J. Yang, and S. Roy, "On joint transmitter and receiver optimization or multiple-input-multiple-output (MIMO) transmission systems," *IEEE Trans. Commun.*, vol. 42, no. 12, pp. 3221–3231, Dec. 1994.

[2] S. Sugiura, S. Chen, and L. Hanzo, "A universal space-time architecture for multiple-antenna aided systems," *IEEE Commun. Surveys Tuts.*, vol. 14, no. 2, pp. 401–420, Apr.–Jun. 2012.

[3] I. E. Telatar, "Capacity of multi-antenna gaussian channels," *European Trans. Commun.*, vol. 10, no. 6, pp. 585–595, Nov.–Dec. 1999.

[4] D. P. Palomar, J. M. Cioffi, and M. A. Lagunas, "Joint Tx-Rx beamforming design for multicarrier MIMO channels: A unified framework for convex optimization," *IEEE Trans. Signal Process.*, vol. 51, no. 9, pp. 2381–2401, Sep. 2003.

[5] E. Vlachos, G. C. Alexandropoulos, and J. Thompson, "Massive MIMO channel estimation for millimeter wave systems via matrix completion," *IEEE Signal Process. Lett.*, vol. 25, no. 11, pp. 1675–1679, Nov. 2018.

[6] C. Xing, S. Wang, S. Chen, S. Ma, H. V. Poor, and L. Hanzo, "Matrix-monotonic optimization - Part I: Single-variable optimization," *IEEE Trans. Signal Process.*, vol. 69, pp. 738–754, 2021.

[7] Z. Chen, N. Zhao, D. K. C. So, J. Tang, X. Y. Zhang, and K. -K. Wong, "Joint altitude and hybrid beamspace precoding optimization for UAV-enabled multiuser mmWave MIMO system," *IEEE Trans. Veh. Tech.*, vol. 71, no. 2, pp. 1713–1725, Feb. 2022.

[8] H. Vaezy, M. J. Omidi, M. M. Naghsh, and H. Yanikomeroglu, "Energy efficient transceiver design in MIMO interference channels: The selfish, unselfish, worst-case, and robust methods," *IEEE Trans. Commun.*, vol. 67, no. 8, pp. 5377–5389, Aug. 2019.

[9] O. E. Ayach, S. Rajagopal, S. Abu-Surra, Z. Pi, and R. W. Heath, "Spatially sparse precoding in millimeter wave MIMO systems," *IEEE Trans. Wireless Commun.*, vol. 13, no. 3, pp. 1499–1513, Mar. 2014.

[10] W. M. Jang, and W. Wu, "Distributed and centralized multiuser detection with antenna arrays," *IEEE Trans. Wireless Commun.*, vol. 4, no. 3, pp. 855–860, May 2005.

[11] S. Han, C. -l. I, Z. Xu, and C. Rowell, "Large-scale antenna systems with hybrid analog and digital beamforming for millimeter wave 5G," *IEEE Commun. Mag.*, vol. 53, no. 1, pp. 186–194, Jan. 2015.

[12] X. Yu, J.-C. Shen, J. Zhang, and K. B. Letaief, "Alternating minimization algorithms for hybrid precoding in millimeter wave MIMO systems," *IEEE J. Sel. Topics Signal Process.*, vol. 10, no. 3, pp. 485–500, Apr. 2016.

[13] T. Qiao, Y. Cao, J. Tang, N. Zhao, and K. -K. Wong, "IRS-aided uplink security enhancement via energy-harvesting jammer," *IEEE Trans. Commun.*, vol. 70, no. 12, pp. 8286–8297, Dec. 2022.

[14] J. Xu, and Y. Liu, "A novel physics-based channel model for reconfigurable intelligent surface-assisted multi-user communication systems," *IEEE Trans. Wireless Commun.*, vol. 21, no. 2, pp. 1183–1196, Feb. 2022.

[15] X. Jin, P. Zhang, C. Wan, D. Ma, and Y. Yao, "RIS assisted dual-function radar and secure communications based on frequency-shifted chirp spread spectrum index modulation," *China Commun.*, vol. 20, no. 10, pp. 85–99, 2023.

[16] S. Gong, C. Xing, Y. Jing, S. Wang, J. Wang, S. Chen, and L. Hanzo, "A unified MIMO optimization framework relaying on the KKT conditions," *IEEE Trans. Commun.*, vol. 69, no. 11, pp. 7251–7268, Aug. 2021.

[17] C. Xing, S. Xie, S. Gong, X. Yang, S. Chen, and L. Hanzo, "A KKT conditions based transceiver optimization framework for RIS-aided multi-user MIMO networks," *IEEE Trans. Commun.*, vol. 71, no. 5, pp. 2602–2617, May 2023.

[18] S. Zhang, and R. Zhang, "Capacity characterization for intelligent reflecting surface aided MIMO communication," *IEEE J. Sel. Areas Commun.*, vol. 38, no. 8, pp. 1823–1838, Aug. 2020.

[19] X. Zhao, K. Xu, S. Ma, S. Gong, G. Yang, and C. Xing, "Joint transceiver optimization for IRS-aided MIMO communications," *IEEE Trans. Commun.*, vol. 70, no. 5, pp. 3467–3482, Mar. 2022.

[20] F. Sohrabi, and W. Yu, "Hybrid digital and analog beamforming design for large-scale antenna arrays," *IEEE J. Sel. Areas Commun.*, vol. 10, no. 3, pp. 501–513, Apr. 2016.

[21] S. Gong, C. Xing, V. K. N. Lau, S. Chen, and L. Hanzo, "Majorization-minimization aided hybrid transceivers for MIMO interference channels," *IEEE Trans. Signal Process.*, vol. 68, pp. 4903–4918, 2020.

[22] C. Xing, Y. Li, S. Gong, J. An, S. Chen, and L. Hanzo, "A general matrix variable optimization framework for MIMO assisted wireless communications," *IEEE Trans. Veh. Tech.*, 2023

[23] H. -T. Wai, W. -K. Ma, and A. M. -C. So, "Cheap semidefinite relaxation MIMO detection using row-by-row block coordinate descent," *IEEE Int. Conf. Acoust., Speech, Signal Process. (ICASSP)*, Prague, Czech Republic, 2011, pp. 3256–3259.

[24] A. Hjørungnes, *Complex-valued matrix derivatives: With applications in signal processing and communications*. Cambridge University Press: Cambridge, UK, 2011.

[25] S. Boyd, and L. Vandenberghe, *Convex optimization*. Cambridge University Press, 2004.

[26] Y. Yu, and Y. Hua, "Power allocation for a MIMO relay system with multiple-antenna users," *IEEE Trans. Signal Process.*, vol. 58, no. 5, pp. 2823–2835, May 2010

[27] S. Abeywickrama, R. Zhang, Q. Wu, and C. Yuen, "Intelligent reflecting surface: Practical phase shift model and beamforming optimization," *IEEE Trans. Commun.*, vol. 68, no. 9, pp. 5849–5863, Sept. 2020.

[28] Q. Shi, M. Razaviyayn, Z.-Q. Luo, and C. He, "An iteratively weighted MMSE approach to distributed sum-utility maximization for a MIMO interfering broadcast channel," *IEEE Trans. Signal Process.*, vol. 59, no. 9, pp. 4331–4340, Sep. 2011

[29] S. S. Christensen, R. Agarwal, E. d. Carvalho, and J. M. Cioffi, "Weighted sum-rate maximization using weighted MMSE for MIMO-BC beamforming design," *IEEE Trans. Wireless Commun.*, vol. 7, no. 12, pp. 4792–4799, Dec. 2008.

[30] S. Serbetli, and A. Yener, "Transceiver optimization for multisuer MIMO systems," *IEEE Trans. Signal Process.*, vol. 52, no. 1, pp. 214–226, Jan. 2004.

[31] S.-R. Lee, J.-S. Kim, S.-H. Moon, H.-B. Kong, and I. Lee, "Zero-forcing beamforming in multiuser MISO downlink systems under per-antenna power constraint and equal-rate metric," *IEEE Trans. Wireless Commun.*, vol. 12, no. 1, pp. 228–236, Jan. 2013

[32] D. P. Palomar, and Y. Jiang, "MIMO transceiver design via majorization theory," *Found. Trends Commun. Inf. Theory*, vol. 3, no. 4, pp. 331–551, 2006.

[33] C. Xing, X. Zhao, W. Xu, X. Dong, and G. Y. Li, "A framework on hybrid MIMO transceiver design based on matrix-monotonic optimization," *IEEE Trans. Signal Process.*, vol. 67, no. 13, pp. 3531–3546, Jul. 2019.

[34] Q. Shi, M. Razaviyayn, M. Hong, and Z. -Q. Luo, "SINR constrained beamforming for a MIMO multi-user downlink system: Algorithms and convergence analysis," *IEEE Trans. Signal Process.*, vol. 64, no. 11, pp. 2920–2933, Jun. 2016.

[35] K. Cumanan, Z. Ding, B. Sharif, G. Y. Tian, and K. K. Leung, "Secrecy rate optimizations for a MIMO secrecy channel with a multiple-antenna eavesdropper," *IEEE Trans. Veh. Tech.*, vol. 63, no. 4, pp. 1678–1690, May 2014.

[36] M. V. Solodov, "On the convergence of constrained parallel variable distribution algorithms," *SIAM J. Optim.*, vol. 8, no. 1, pp. 187–196, Feb. 1998.

[37] M. Grant, and S. Boyd. (Sep. 2013). *CVX: MATLAB Software for Disciplined Convex Programming, Version 2.0 Beta.* [Online]. Available: http://cvxr.com/cvx

[38] X. Yu, J. Shen, J. Zhang, and K. B. Letaief, "Alternating minimization algorithms for hybrid precoding in millimeter wave MIMO systems," *IEEE J. Sel. Topics Signal Process.*, vol. 10, no. 3, pp. 485–500, Feb. 2016.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60