# Efficient Channel Estimation for RIS-Aided MIMO Communications with Unitary Approximate Message Passing

Yabo Guo, Peng Sun, Zhengdao Yuan, Chongwen Huang, Qinghua Guo, *Senior Member, IEEE*, Zhongyong Wang, and Chau Yuen, *Fellow, IEEE*

*Abstract*—Reconfigurable intelligent surface (RIS) is very promising for wireless networks to achieve high energy efficiency, extended coverage, improved capacity, massive connectivity, etc. To unleash the full potentials of RIS-aided communications, acquiring accurate channel state information is crucial, which however is very challenging. For RIS-aided multiple-input and multiple-output (MIMO) communications, the existing channel estimation methods have computational complexity growing rapidly with the number of RIS units $N$ (e.g., in the order of $N^2$ or $N^3$) and/or have special requirements on the matrices involved (e.g., the matrices need to be sparse for algorithm convergence to achieve satisfactory performance), which hinder their applications. In this work, instead of using the conventional signal model in the literature, we derive a new signal model obtained through proper vectorization and reduction operations. Then, leveraging the unitary approximate message passing (UAMP), we develop a more efficient channel estimator that has complexity linear with $N$ and does not have special requirements on the relevant matrices, thanks to the robustness of UAMP. These facilitate the applications of the proposed algorithm to a general RIS-aided MIMO system with a larger $N$. Moreover, extensive numerical results show that the proposed estimator delivers much better performance and/or requires significantly less number of training symbols, thereby leading to notable reductions in both training overhead and latency.

*Index Terms*—Reconfigurable intelligent surface (RIS), channel estimation, approximate message passing (AMP).

## I. INTRODUCTION

Y. Guo, P. Sun and Z. Wang are with the School of Information Engineering, Zhengzhou University, Zhengzhou 450002, China (e-mail: ieybguo@163.com, iepengsun@zzu.edu.cn, zywangzzu@gmail.com).

Z. Yuan is with the Artificial Intelligence Technology Engineering Research Center, Open University of Henan, and also with the School of Information Engineering, Zhengzhou University, Zhengzhou 450002, China (e-mail: yuan_zhengdao@163.com).

C. Huang is with the College of Information Science and Electronic Engineering, Zhejiang University, Hangzhou 310007, China, and Zhejiang Provincial Key Lab of Information Processing, Communication and Networking (IPCAN), Hangzhou 310007, China, and the International Joint Innovation Center, Zhejiang University, Haining 314400, China (e-mail: chongwenhuang@zju.edu.cn).

Q. Guo is with the School of Electrical, Computer and Telecommunications Engineering, University of Wollongong, Wollongong, NSW 2522, Australia (e-mail: qguo@uow.edu.au).

C. Yuen is with the Engineering Product Development (EPD) Pillar, Singapore University of Technology and Design, Singapore 487372 (e-mail: yuenchau@sutd.edu.sg).

RECONFIGURABLE intelligent surface (RIS) is a reconfigurable planar array with a massive number of passive reflecting elements, which is capable of altering the wireless propagation environment to achieve desired channel responses. RIS has been recognized as a promising technology in future wireless communications. With the aid of RIS, wireless networks are able to achieve high energy efficiency, improve the system capacity and radio coverage, enhance massive connectivity, etc [1]–[16].

The research on RIS has attracted tremendous attention and many works have been conducted to explore the potentials of RIS-aided communications. Under the assumption of perfect channel state information (CSI), energy-efficient designs were studied in [1] with the assist of RIS for wireless communications. The works in [2] and [8] studied RIS-aided secrecy communications. The work in [15] investigated the use of RIS to enhance the energy harvesting and information transmission capabilities of the internet of things systems. Energy efficient transmission with distributed RISs was studied in [6]. The works in [14] and [17] studied RIS-aided non-othogonal multiple access. RIS-aided communications were extended to millimeter-wave and Terahertz communications in [9], [16], [18]. The works in [7] and [19] maximize the weighted sum rate for both single and multi-cells networks. To unleash the potentials of RIS aided communications, efficient accurate CSI acquisition is crucial [5], [20]–[30].

In this work, we focus on channel estimation of RIS-aided multiple-input multiple-output (MIMO) communications, which is challenging especially when the number of RIS unit $N$ is large. In [25], a two-stage algorithm that includes a sparse matrix factorization stage and a matrix completion stage was proposed for the estimation of cascaded (transmitter-RIS and RIS-receiver) channel, where the RIS phase matrix needs to be sparse and the channel between RIS and BS is required to be a low-rank matrix. The work in [26] proposed a trilinear semi-blind cascaded channel estimation problem, in which the receiver estimates the channel coefficients and the transmitted signals jointly, and a message passing algorithm was developed. The algorithm has relatively low complexity, but it also requires that the relevant matrices are sparse for the convergence of the algorithm, thereby achieving good performance. A message-passing algorithm was proposed in [27] to estimate the cascaded channels by exploiting the information on the slow-varying channel components and the hidden channel sparsity, and its complexity growing with $N^2$.

The sparseness of the matrices may be achieved by using a sparse RIS phase matrix, which however is not optimal in terms of channel estimation and transmission, or by restricting the applications to special scenarios with sparse channel matrices. The work in [28] proposed a three-phase pilot-based channel estimation framework for RIS-assisted uplink multiuser communications, and the complexity of the method is cubic in the number of RIS units $N$. In [29], leveraging the parallel factor (PARAFAC) decomposition to represent high dimensional tensors that involve unknown channels in different unfolded forms, channel estimation methods based on alternating least squares (ALS) and vector approximate message passing (VAMP) were proposed to estimate the channels from BS to RIS and from RIS to users alternatively. Although there are no requirements on the sparseness of the involved matrices, its complexity grows with $N^3$, which is a concern for a large or even a moderate $N$. In summary, the existing methods have special requirements on the involved matrices and/or have the scalability issue as their computational complexity rapidly grows with $N$, which hinder their applications.

To overcome the problems of the existing methods, leveraging the unitary approximate message passing (UAMP) algorithm [31]–[33], we design a new RIS channel estimator in this work. As a variant of the AMP algorithm [34], UAMP achieves remarkably improved robustness by using a unitary transformation, which enables it to deal with a linear reverse problem with a general (or tough) system transfer matrix while with low complexity [32], [33]. In this work, instead of using the conventional signal model for RIS channel estimation in the literature, we derive a new signal model through proper vectorization and reduction operations, which reformulates the channel estimation to a structured signal recovery problem. Then a factor graph representation is developed, and a message passing algorithm is derived, where UAMP plays a crucial role. Thanks to the low complexity and robustness of UAMP, the proposed algorithm is very efficient, which has complexity linear with the number of RIS units $N$, and does not have special requirements on the relevant matrices. These enable the applications of the proposed algorithm to a general RIS aided MIMO system with a larger $N$. The Cramér-Rao lower bound (CRLB) for the channel estimation considered is derived to serve as a performance benchmark. Extensive numerical results show that, with much lower complexity, the proposed channel estimator significantly outperforms existing ones in terms of channel estimation performance and/or the training overhead.

The remainder of the paper is organized as follows. In Section II, we introduce the system model and problem formulation of RIS channel estimation. In Section III, the problem is reformulated, a new signal model is obtained and the problem is represented as a factor graph. Then the UAMP based message passing algorithm is developed in Section III-B, and the CRLB for channel estimation is derived in Section V. Numerical results are provided in Section VI, followed by conclusions in Section VII.

*Notations*-Boldface lower-case and upper-case letters denote vectors and matrices, respectively. Superscripts $\boldsymbol{A}^H$ and $\boldsymbol{A}^T$ represent conjugate transpose and transpose, respectively, and $\boldsymbol{A}^*$ represents the conjugate of $\boldsymbol{A}$. A Gaussian distribution of $x$ with mean $\hat{x}$ and variance $\nu_x$ is denoted by $\mathcal{N}(x; \hat{x}, \nu_x)$. Notations $\otimes$ and $\odot$ represent the Kronecker and Khatri-Rao products, respectively. The relation $f(x) = cg(x)$ for some positive constant $c$ is written as $f(x) \propto g(x)$. We use $\boldsymbol{a} \cdot \boldsymbol{b}$ and $\boldsymbol{a} \cdot /\boldsymbol{b}$ to represent the element-wise product and division between vectors $\boldsymbol{a}$ and $\boldsymbol{b}$, respectively. The notation $\boldsymbol{a}^{-1}$ denotes the element-wise inverse operation to vector $\boldsymbol{a}$. We use $|\boldsymbol{A}|^2$ to denote element-wise magnitude squared operation for $\boldsymbol{A}$, $\|\boldsymbol{a}\|$ to denote the $l_2$ norm of $\boldsymbol{a}$, and $\|\boldsymbol{A}\|_F$ to denote the Frobenius norm of $\boldsymbol{A}$. The notation $<\boldsymbol{a}>$ denotes the average operation for $\boldsymbol{a}$, i.e., the sum of the elements of $\boldsymbol{a}$ divided by its length. We use $\boldsymbol{1}$ and $\boldsymbol{0}$ to denote an all-one vector and an all-zero vector with a proper length, respectively. The notation $\mathbf{Diag}(\boldsymbol{a})$ represents a diagonal matrix with $\boldsymbol{a}$ as its diagonal and $\boldsymbol{I}_n$ donates a $n \times n$ identity matrix.

## II. System Model and Problem Formulation



Fig. 1: Illustration of RIS-aided MIMO uplink transmission.

We consider RIS-aided MIMO uplink transmission as shown in Fig. 1, where the BS equipped with $M$ antennas receives signals from $K$ users. A RIS with $N$ passive reflecting elements is equipped between the BS and users, each user equipped with a single antenna (the extension to multiple antennas is straightforward). The RIS is attached to the facade of a building in the vicinity of users. Due to the highly attenuation caused by unfavorable propagation environments such as tall buildings, the direct propagation path between the BS and users is neglected [29].

We denote the channel matrix between the BS and the RIS by $\boldsymbol{G} \in \mathbb{C}^{M \times N}$, and use $\boldsymbol{H} \triangleq [\boldsymbol{h}_1^u, \ldots, \boldsymbol{h}_K^u] \in \mathbb{C}^{N \times K}$ to represent the channel matrix between the RIS and the $K$ users, where $\boldsymbol{h}_k^u$ denote the channel vector from the $k$-th user to the RIS. With the $l$-th ($l = 1, 2, \ldots, L$) RIS phase configuration, the $K$ consecutive received signals $\boldsymbol{Y}_l \in \mathbb{C}^{M \times K}$ is given by

$$\boldsymbol{Y}_l = \boldsymbol{G}\mathbf{Diag}(\boldsymbol{\Phi}_{l,:})\boldsymbol{H}\boldsymbol{X} + \boldsymbol{W}_l, \qquad (1)$$

where $\boldsymbol{\Phi}_{l,:}$ is the $l$-th row of the RIS phase matrix $\boldsymbol{\Phi} \in \mathbb{C}^{L \times N}$ (part of the RIS phase configurations) and $\boldsymbol{W}_l$ models the zero mean complex additive white Gaussian noise (AWGN) with precision $\beta$ (i.e., variance $\beta^{-1}$), and $\boldsymbol{X} \in \mathbb{C}^{K \times K}$ denotes the transmitted orthogonal training matrix from the users, i.e.,

$XX^H = I_K$. Right-multiplying both sides of (1) by $X^H$ leads to

$$\tilde{Y}_l = G\text{Diag}(\Phi_{l,:})H + \tilde{W}_l, \qquad (2)$$

where $\tilde{Y}_l \triangleq Y_l X^H \in \mathbb{C}^{M\times K}$, and $\tilde{W}_l \triangleq W_l X^H \in \mathbb{C}^{M\times K}$. As $X$ is a unitary matrix, entries of $\tilde{W}_l$ are still zero-mean white Gaussian with the same noise precision $\beta$.

During the training process, $L$ configurations of the RIS are used, leading to $L$ matrices $\{\tilde{Y}_l, l = 1, ..., L\}$, based which we aim to estimate the channel matrices. As the number of RIS unit $N$ can be large, low complexity of channel estimators is crucial for practical implementation. The cubic or quadratic complexity (in terms of $N$) of existing channel estimation algorithms can be a concern in their applications. During the training process, the RIS undergoes $L$ configurations. In order to reduce the training overhead and communication latency, a small $L$ is desirable. We address the above challenges by reformulating the RIS channel estimation to a structured signal recovery problem with a new signal model and develop an efficient message passing algorithm. In particular, we do not make any special requirements on the matrices $G$, $H$ or $\Phi$.

## III. PROBLEM REFORMULATION AND FACTOR GRAPH REPRESENTATION

### A. New Model for RIS Channel Estimation

Instead of using model (2) directly, we reformulate it to a new model through vectorization and reduction, which leads to a structured signal recovery problem. Vectorizing (2) leads to

$$\mathbf{vec}(\tilde{Y}_l) = (H^T \otimes G)\mathbf{vec}(\text{Diag}(\Phi_{l,:})) + \mathbf{vec}(\tilde{W}_l), \qquad (3)$$

where the vector $\mathbf{vec}(\text{Diag}(\Phi_{l,:}))$ can be represented as

$$\mathbf{vec}(\text{Diag}(\Phi_{l,:})) = [\phi_{l,1}, \mathbf{0}_N^T, \phi_{l,2}, \mathbf{0}_N^T, \ldots, \mathbf{0}_N^T, \phi_{l,N}]^T, \quad (4)$$

and $\phi_{l,n}$ donates the $n$-th ($n = 1, 2, \ldots, N$) element of $\Phi_{l,:}$. Here we note that, the vector $\mathbf{vec}(\text{Diag}(\Phi_{l,:}))$ is highly sparse, and the non-zeros elements of $\mathbf{vec}(\text{Diag}(\Phi_{l,:}))$ are separated by all-zero vector $\mathbf{0}_N$. This can be exploited to significantly reduce the dimension of (3). We drop the zero elements in $\mathbf{vec}(\text{Diag}(\Phi_{l,:}))$ and the corresponding columns in $H^T \otimes G$, then (3) can be reduced to

$$\tilde{y}_l = \tilde{S}\Phi_{l,:}^T + \tilde{w}_l, \qquad (5)$$

where $\tilde{y}_l \triangleq \mathbf{vec}(\tilde{Y}_l) \in \mathbb{C}^{KM\times 1}$, $\tilde{S} \triangleq H^T \odot G \in \mathbb{C}^{KM\times N}$, $\tilde{S}\Phi_{l,:}^T = (H^T \otimes G)\mathbf{vec}(\text{Diag}(\Phi_{l,:}))$, and $\tilde{w}_l \triangleq \mathbf{vec}(\tilde{W}_l) \in \mathbb{C}^{KM\times 1}$. By stacking $\{\tilde{y}_l, l = 1, ..., L\}$ into a matrix, we have $\tilde{Y} = \tilde{S}\Phi^T + \tilde{W}$, which is rewritten as

$$Y = \Phi S + W. \qquad (6)$$

where $Y = \tilde{Y}^T$, $W = \tilde{W}^T$ and

$$S = \tilde{S}^T = (H^T \odot G)^T \in \mathbb{C}^{N\times KM}. \qquad (7)$$

Now we can see from (6) that, the channel estimation is reformulated as the recovery of the matrix $S$, which admits the structure (7). We can also treat $S$ as an intermediate variable, as our aim is to estimate $H$ and $G$.

### B. Probabilistic and Factor Graph Representation

We consider recovering the signal using the message passing techniques, in particular leveraging UAMP to achieve low complexity while with high robustness. We first represent the problem in a probabilistic form. It is noted that UAMP works with a unitary transform of the linear observation model (6), which is crucial to achieving high robustness. So, to facilitate the use of UAMP later, we first carry out a unitary transformation to (6) based on the singular value decomposition (SVD) $\Phi = U\Lambda V$, leading to

$$R = \Psi S + \overline{W}, \qquad (8)$$

where $R = U^H Y$, $\Psi = U^H\Phi = \Lambda V$ and $\overline{W} = U^H W$. Since $U$ is unitary, entries in $\overline{W}$ are still AWGN with precision $\beta$.

Note that $\tilde{S} = [\tilde{s}_1, \ldots, \tilde{s}_N]$, $H^T = [h_1, \ldots, h_N]$ and $G = [g_1, \ldots, g_N]$. Then, according to $\tilde{S} = H^T \odot G$, we have

$$\tilde{s}_n = h_n \otimes g_n, \qquad (9)$$

where $h_n = [h_{1,n}, \ldots, h_{K,n}]^T$, $g_n = [g_{1,n}, \ldots, g_{M,n}]^T$ and $\tilde{s}_n = [\tilde{s}_{1,1,n}, \ldots, \tilde{s}_{m,k,n}, \ldots, \tilde{s}_{M,K,n}]^T$ with $\tilde{s}_{m,k,n} = h_{k,n}g_{m,n}$. Let $J = KM$, and note that $R = [r_1, \ldots, r_J]$, $S = [s_1, \ldots, s_J]$, and $\overline{W} = [w_1, \ldots, w_J]$. Define an auxiliary variable $Z \triangleq [z_1, \ldots, z_J]$ with $z_j = \Psi s_j$. Then, the joint distribution of $H, G, S, \tilde{S}, Z$ and $\beta$ given $R$ can be factorized as

$$p(H, G, S, \tilde{S}, Z, \beta|R)$$
$$\propto p(R|Z, \beta)p(\tilde{S}|S)p(Z|S)p(\tilde{S}|H, G)p(H)p(G)p(\beta)$$
$$= p(\tilde{S}|S)p(\beta)\prod_n p(\tilde{s}_n|h_n, g_n)p(h_n)p(g_n)\prod_j p(r_j|z_j, \beta)$$
$$\times p(z_j|s_j)$$
$$\triangleq f_S(\tilde{S}, S)f_\beta(\beta)\prod_n f_{\tilde{s}_n}(\tilde{s}_n, h_n, g_n)f_{h_n}(h_n)f_{g_n}(g_n)$$
$$\times \prod_j f_{r_j}(z_j, \beta)f_{z_j}(z_j, s_j). \qquad (10)$$

where the involved distributions are listed in Table I. To facilitate the factor graph representation of the factorization in (10), local functions (factors) are defined, and the correspondence between the distributions and local functions are also shown in Table I. It is noted that $p(H)$ and $p(G)$ represent the priors for the channel matrices $H$ and $G$, respectively. When no priors are available for the channel matrices, the priors can be set to be a non-informative one, e.g., $\rho_h = \rho_g = +\infty$ in Table I. The factor graph representation is depicted in Fig. 2.

TABLE I: Factors and distributions in (10).

| Factor | Distribution | Function |
|---|---|---|
| $f_\beta$ | $p(\beta)$ | $\propto \beta^{-1}$ |
| $f_{h_n}$ | $p(h_n)$ | $\mathcal{CN}(h_n; \mathbf{0}_K, \rho_h I_K)$ |
| $f_{g_n}$ | $p(g_n)$ | $\mathcal{CN}(g_n; \mathbf{0}_M, \rho_g I_M)$ |
| $f_{r_j}$ | $p(r_j|z_j, \beta)$ | $\mathcal{N}(r_j; z_j, \beta^{-1}I_L)$ |
| $f_{z_j}$ | $p(z_j|s_j)$ | $\delta(z_j - \Psi s_j)$ |
| $f_{\tilde{s}_n}$ | $p(\tilde{s}_n|h_n, g_n)$ | $\delta(\tilde{s}_n - h_n \otimes g_n)$ |
| $f_S$ | $p(\tilde{S}|S)$ | $\delta(\tilde{S}^T - S)$ |

Fig. 2: Factor graph representation of (10).

Our aim is to find the posteriori distributions $p(\boldsymbol{H}|\boldsymbol{R})$ and $p(\boldsymbol{G}|\boldsymbol{R})$ and their estimates in terms of a posteriori means, i.e., $\hat{\boldsymbol{H}} = \mathbb{E}\{\boldsymbol{H}|\boldsymbol{R}\}$ and $\hat{\boldsymbol{G}} = \mathbb{E}\{\boldsymbol{G}|\boldsymbol{R}\}$. We note that it is difficult to find the exact a posteriori distributions, and approximate inference has to be resorted. In next section, we will develop a low complexity message passing algorithm to find their approximations. To facilitate the message passing algorithm design, a scalar factor graph representation is shown in Fig. 3.

## IV. UAMP Based Message Passing Algorithm For RIS Channel Estimation

In this section, we develop an efficient message passing algorithm for approximate inference, where UAMP is incorporated to deal with the most computationally intensive part of message passing, which is crucial to achieving low complexity while with high robustness. The message passing algorithm carries out an iterative process, where each iteration involves a forward message passing process and backward message passing process in the graph shown in Fig. 2 or Fig. 3. We use $m_{A \to B}(x)$ to denote a message passed from node $A$ to node $B$, which is a function of $x$. For Gaussian messages, the arrows above its mean and variance indicate the direction of the message passing. In addition, we use $b(x)$ to denote the belief of a variable $x$. Note that, if a forward computation requires backward messages, the relevant messages in the previous iteration is used by default.

### A. Forward Message Passing

In the forward direction, according to the rules of variational message passing [35], we have

$$m_{f_{\boldsymbol{r}_j} \to \boldsymbol{z}_j}(\boldsymbol{z}_j) \propto \exp\left\{ \int_\beta \mathfrak{b}(\beta) \log f_{\boldsymbol{r}_j} \right\}$$
$$\propto \mathcal{N}\left(\boldsymbol{z}_j; \boldsymbol{r}_j, \hat{\beta}^{-1}\right), \quad (11)$$

where

$$\mathfrak{b}(\beta) \propto m_{f_{\boldsymbol{r}_j} \to \beta}(\beta) f_\beta$$
$$\propto \beta^{LJ-1} \exp\left\{ \sum_{j=1} -\beta \left(\|\boldsymbol{r}_j - \hat{\boldsymbol{z}}_j\|^2 + \mathbf{1}^T \boldsymbol{\nu}_{\boldsymbol{z}_j}\right) \right\}, \quad (12)$$

and

$$\hat{\beta} = \int_\beta \beta \mathfrak{b}(\beta) = \frac{LJ}{\sum_{j=1}(\|\boldsymbol{r}_j - \hat{\boldsymbol{z}}_j\|^2 + \mathbf{1}_L^T \boldsymbol{\nu}_{\boldsymbol{z}_j})}. \quad (13)$$

It is noted that in the above equation, $\hat{\boldsymbol{z}}_j$ and the message $m_{f_{\boldsymbol{r}_j} \to \beta}(\beta)$ are required to compute $\hat{\beta}$, which are obtained from the last iteration and their computations are delayed to (71) and (74).

The Gaussian form of the message in (11) suggests the following model

$$\boldsymbol{r}_j = \boldsymbol{z}_j + \boldsymbol{w}_j, \quad (14)$$

where the noise $\boldsymbol{w}_j$ is Gaussian with mean zero and precision $\hat{\beta}$. This allows seamless integration with the forward recursion of UAMP. Here, assume that the mean and variance of $\boldsymbol{s}_j$ are available, i.e., $\hat{\boldsymbol{s}}_j$ and $\nu_{\boldsymbol{s}_j}$. Specially, we assume $\boldsymbol{s}_j$ have a common variance $\nu_{\boldsymbol{s}_j}$, and the corresponding computation will be detailed later. Following UAMP, we define a vector $\boldsymbol{\psi} \in \mathbb{C}^{L \times 1}$ as

$$\boldsymbol{\psi} = |\boldsymbol{\Psi}|^2 \mathbf{1}_N. \quad (15)$$

Then we calculate two vectors $\boldsymbol{\nu}_{\boldsymbol{p}_j}$ and $\boldsymbol{p}_j$ as

$$\boldsymbol{\nu}_{\boldsymbol{p}_j} = \boldsymbol{\psi} \nu_{\boldsymbol{s}_j}, \quad (16)$$

$$\boldsymbol{p}_j = \boldsymbol{\Psi} \hat{\boldsymbol{s}}_j - \boldsymbol{\nu}_{\boldsymbol{p}_j} \cdot \boldsymbol{\mu}_j, \quad (17)$$

where $\boldsymbol{\mu}_j$ is a vector that is computed in last iteration. According to UAMP, we update the intermediate vectors $\boldsymbol{\nu}_{\boldsymbol{\mu}_j}$ and $\boldsymbol{\mu}_j$ by

$$\boldsymbol{\nu}_{\boldsymbol{\mu}_j} = \mathbf{1}_L . / (\boldsymbol{\nu}_{\boldsymbol{p}_j} + \hat{\beta}^{-1} \mathbf{1}_L), \quad (18)$$

$$\boldsymbol{\mu}_j = \boldsymbol{\nu}_{\boldsymbol{\mu}_j} \cdot (\boldsymbol{r}_j - \boldsymbol{p}_j), \quad (19)$$

Then compute vectors $\boldsymbol{\nu}_{\boldsymbol{q}_j}$ and $\boldsymbol{q}_j$ with

$$\boldsymbol{\nu}_{\boldsymbol{q}_j} = \mathbf{1}_N . / |\boldsymbol{\Psi}^H|^2 \boldsymbol{\nu}_{\boldsymbol{\mu}_j}, \quad (20)$$

$$\boldsymbol{q}_j = \hat{\boldsymbol{s}}_j + \boldsymbol{\nu}_{\boldsymbol{q}_j} . \boldsymbol{\Psi}^H \boldsymbol{\mu}_j, \quad (21)$$

The message $\boldsymbol{q}_j$ and $\boldsymbol{\nu}_{\boldsymbol{q}_j}$ are the mean and variance of $\boldsymbol{s}_j$. According to the belief propagation derivation of (U)AMP, we have

$$m_{\boldsymbol{s}_j \to f_{\boldsymbol{s}}}(\boldsymbol{s}_j) = \mathcal{N}(\boldsymbol{s}_j; \boldsymbol{q}_j, \mathbf{Diag}(\boldsymbol{\nu}_{\boldsymbol{q}_j})). \quad (22)$$

Stack $\boldsymbol{q}_j$ and $\boldsymbol{\nu}_{\boldsymbol{q}_j}$ into matrices as

$$\boldsymbol{Q} = [\boldsymbol{q}_1, \ldots, \boldsymbol{q}_J], \quad (23)$$

$$\boldsymbol{\nu}_Q = [\boldsymbol{\nu}_{\boldsymbol{q}_1}, \ldots, \boldsymbol{\nu}_{\boldsymbol{q}_J}], \quad (24)$$

Due to the deterministic relationship between $\boldsymbol{S}$ and $\tilde{\boldsymbol{S}}$, i.e., $\tilde{\boldsymbol{S}}^T = \boldsymbol{S}$, as shown in Table I, we have

$$\tilde{\boldsymbol{Q}} = \boldsymbol{Q}^T = [\tilde{\boldsymbol{q}}_1, \ldots, \tilde{\boldsymbol{q}}_N], \quad (25)$$

$$\boldsymbol{\nu}_{\tilde{Q}} = \boldsymbol{\nu}_Q^T = [\boldsymbol{\nu}_{\tilde{\boldsymbol{q}}_1}, \ldots, \boldsymbol{\nu}_{\tilde{\boldsymbol{q}}_N}], \quad (26)$$

and the vectors $\tilde{\boldsymbol{q}}_n \in \mathbb{C}^{J \times 1}$ and $\boldsymbol{\nu}_{\tilde{\boldsymbol{q}}_n} \in \mathbb{C}^{J \times 1}$ can be divided into $K$ length-$M$ vectors, i.e.,

$$\tilde{\boldsymbol{q}}_n = [\tilde{\boldsymbol{q}}_{1,n}^T, \ldots, \tilde{\boldsymbol{q}}_{K,n}^T]^T, \quad (27)$$

Fig. 3: Scalar factor graph representation of (10).

$$\boldsymbol{\nu}_{\tilde{\boldsymbol{q}}_n} = [\boldsymbol{\nu}_{\tilde{\boldsymbol{q}}_{1,n}}^T, \dots, \boldsymbol{\nu}_{\tilde{\boldsymbol{q}}_{K,n}}^T]^T. \tag{28}$$

Let $\nu_{\tilde{\boldsymbol{q}}_{k,n}} = \langle \boldsymbol{\nu}_{\tilde{\boldsymbol{q}}_{k,n}} \rangle$, and the message $m_{\tilde{\boldsymbol{s}}_n \to f_{\tilde{\boldsymbol{s}}_n}}(\tilde{\boldsymbol{s}}_n)$ can be expressed as

$$m_{\tilde{\boldsymbol{s}}_n \to f_{\tilde{\boldsymbol{s}}_n}}(\tilde{\boldsymbol{s}}_n) = \mathcal{N}(\tilde{\boldsymbol{s}}_n; \tilde{\boldsymbol{q}}_n, \mathbf{Diag}(\boldsymbol{\nu}'_{\tilde{\boldsymbol{q}}_n})), \tag{29}$$

where

$$\boldsymbol{\nu}'_{\tilde{\boldsymbol{q}}_n} = [\nu_{\tilde{\boldsymbol{q}}_{1,n}}, \dots, \nu_{\tilde{\boldsymbol{q}}_{K,n}}]^T \otimes \mathbf{1}_M. \tag{30}$$

Next, we compute the message at $f_{\tilde{\boldsymbol{s}}_n}$, $\boldsymbol{h}_n$ and $\boldsymbol{g}_n$. It is noted that $\tilde{\boldsymbol{s}}_n \in \mathbb{C}^{J \times 1}$ also can be divided into $K$ length-$M$ vectors and expressed as

$$\tilde{\boldsymbol{s}}_n = [\tilde{\boldsymbol{s}}_{1,n}^T, \dots, \tilde{\boldsymbol{s}}_{K,n}^T]^T, \tag{31}$$

where $\tilde{\boldsymbol{s}}_{k,n} \triangleq [\tilde{s}_{1,k,n}, \dots, \tilde{s}_{M,k,n}]^T$. We further factorize the function $f_{\tilde{\boldsymbol{s}}_n}(\tilde{\boldsymbol{s}}_n, \boldsymbol{h}_n, \boldsymbol{g}_n)$ as

$$f_{\tilde{\boldsymbol{s}}_n}(\tilde{\boldsymbol{s}}_n, \boldsymbol{h}_n, \boldsymbol{g}_n) = \prod_{m,k} f_{\tilde{s}_{m,k,n}}(h_{k,n}, g_{m,n}), \tag{32}$$

and the factor $f_{\tilde{s}_{m,k,n}}(h_{k,n}, g_{m,n})$ is shown in Fig. 4.



Fig. 4: Factor graph representation of $f_{\tilde{s}_{m,k,n}}$.

With the definition $\tilde{\boldsymbol{q}}_{k,n} \triangleq [\tilde{q}_{1,k,n}, \dots, \tilde{q}_{M,k,n}]^T$, (29) implies that

$$m_{\tilde{s}_{m,k,n} \to f_{\tilde{s}_{m,k,n}}}(\tilde{s}_{m,k,n}) = \mathcal{N}(\tilde{s}_{m,k,n}; \tilde{q}_{m,k,n}, \nu_{\tilde{\boldsymbol{q}}_{k,n}}), \tag{33}$$

and we note that the factor $f_{\tilde{s}_{m,k,n}} = \delta(\tilde{s}_{m,k,n} - h_{k,n}g_{m,n})$ as shown in Table I. To compute the message $m_{f_{\tilde{s}_{m,k,n}} \to g_{m,n}}(g_{m,n})$ with belief propagation at factor node $f_{\tilde{s}_{m,k,n}}$, we need to integrate out $s_{m,k,n}$ and $h_{k,n}$.

However, due to the multiplication of $g_{m,n}$ and $h_{k,n}$, the message will be intractable even the incoming message $m_{h_{k,n} \to f_{\tilde{s}_{m,k,n}}}(h_{k,n})$ is Gaussian. To solve this, we first apply belief propagation and eliminate the variable $\tilde{s}_{m,k,n}$ to get an intermediate function node $\tilde{f}_{\tilde{s}_{m,k,n}}(h_{k,n}, g_{m,n})$, i.e.,

$$\tilde{f}_{\tilde{s}_{m,k,n}}(h_{k,n}, g_{m,n}) = \int_{\tilde{s}_{m,k,n}} m_{\tilde{s}_{m,k,n} \to f_{\tilde{s}_{m,k,n}}}(\tilde{s}_{m,k,n}) \cdot f_{\tilde{s}_{m,k,n}}$$
$$= \mathcal{N}(\tilde{s}_{m,k,n}; \tilde{q}_{m,k,n}, \nu_{\tilde{\boldsymbol{q}}_{k,n}}). \tag{34}$$

This turns the function node $f_{\tilde{s}_{m,k,n}}$ with the hard constraint $\delta(\tilde{s}_{m,k,n} - h_{k,n}g_{m,n})$ to a soft function node, enabling the use of variational inference to handle $h_{k,n}$ and $g_{m,n}$. With the intermediate local function $\tilde{f}_{\tilde{s}_{m,k,n}}(h_{k,n}, g_{m,n})$, we can compute the outgoing message $m_{f_{\tilde{s}_{m,k,n}} \to g_{m,n}}(g_{m,n})$ as

$$m_{f_{\tilde{s}_{m,k,n}} \to g_{m,n}}(g_{m,n}) \propto \exp\left\{ \int_{h_{k,n}} \mathfrak{b}(h_{k,n}) \log \tilde{f}_{\tilde{s}_{m,k,n}} \right\}$$
$$\propto \mathcal{N}(g_{m,n}; \vec{g}_{m,k,n}, \vec{\nu}_{g_{m,k,n}}), \tag{35}$$

where

$$\vec{\nu}_{g_{m,k,n}} = \frac{\nu_{\tilde{\boldsymbol{q}}_{k,n}}}{\left|\hat{h}_{k,n}\right|^2 + \nu_{h_{k,n}}}, \tag{36}$$

$$\vec{g}_{m,k,n} = \frac{\tilde{q}_{m,k,n} \hat{h}_{k,n}^*}{\left|\hat{h}_{k,n}\right|^2 + \nu_{h_{k,n}}}, \tag{37}$$

with $\hat{h}_{k,n}$ and $\nu_{h_{k,n}}$ being the approximate a posteriori mean and variance of $h_{k,n}$, which are computed in (52) and (51). With belief propagation and referring to Fig. 4, the message $m_{g_{m,n} \to f_{g_{m,n}}}(g_{m,n})$ can be represented as

$$m_{g_{m,n} \to f_{g_{m,n}}}(g_{m,n}) = \mathcal{N}(g_{m,n}; \vec{g}_{m,n}, \vec{\nu}_{g_{m,n}}), \tag{38}$$

with

$$\vec{\nu}_{g_{m,n}} = 1 / \sum_{k=1}^{K} \frac{1}{\vec{\nu}_{g_{m,k,n}}}, \tag{39}$$

$$\vec{g}_{m,n} = \vec{\nu}_{g_{m,n}} \sum_{k=1}^{K} \frac{\vec{g}_{m,k,n}}{\vec{\nu}_{g_{m,k,n}}}. \tag{40}$$

So, the marginal of $g_{m,n}$ can be expressed as

$$\mathfrak{b}\left(g_{m,n}\right) = m_{g_{m,n}\to f_{g_{m,n}}}\left(g_{m,n}\right) f_{g_{m,n}}$$
$$\propto \mathcal{N}\left(g_{m,n}; \hat{g}_{m,n}, \nu_{g_{m,n}}\right), \tag{41}$$

with

$$\nu_{g_{m,n}} = \frac{\vec{\nu}_{g_{m,n}}\rho_g}{\rho_g + \vec{\nu}_{g_{m,n}}}, \tag{42}$$

$$\hat{g}_{m,n} = \frac{\vec{g}_{m,n}\rho_g}{\rho_g + \vec{\nu}_{g_{m,n}}}. \tag{43}$$

Similarly, we can compute the message $m_{f_{\tilde{s}_{m,k,n}}\to h_{k,n}}$ as

$$m_{f_{\tilde{s}_{m,k,n}}\to h_{k,n}}\left(h_{k,n}\right) \propto \mathcal{N}\left(h_{k,n}; \vec{h}_{m,k,n}, \vec{\nu}_{h_{m,k,n}}\right), \tag{44}$$

where

$$\vec{\nu}_{h_{m,k,n}} = \frac{\nu_{\tilde{q}_{k,n}}}{|\hat{g}_{m,n}|^2 + \nu_{g_{m,n}}}, \tag{45}$$

$$\vec{h}_{m,k,n} = \frac{\tilde{q}_{m,k,n}\hat{g}^*_{m,n}}{|\hat{g}_{m,n}|^2 + \nu_{g_{m,n}}}, \tag{46}$$

with $\hat{g}_{m,n}$ and $\nu_{g_{m,n}}$ being the approximate a posteriori mean and variance of $g_{m,n}$, which are computed in (43) and (42). With belief propagation, the message $m_{h_{k,n}\to f_{h_{k,n}}}\left(h_{k,n}\right)$ can be represented as

$$m_{h_{k,n}\to f_{h_{k,n}}}\left(h_{k,n}\right) = \mathcal{N}\left(h_{k,n}; \vec{h}_{k,n}, \vec{\nu}_{h_{k,n}}\right), \tag{47}$$

with

$$\vec{\nu}_{h_{k,n}} = 1/\sum_{m=1}^{M}\frac{1}{\vec{\nu}_{h_{m,k,n}}}, \tag{48}$$

$$\vec{h}_{k,n} = \vec{\nu}_{h_{k,n}} \sum_{m=1}^{M}\frac{\vec{h}_{m,k,n}}{\vec{\nu}_{h_{m,k,n}}}. \tag{49}$$

The marginal of $h_{k,n}$ can be expressed as

$$\mathfrak{b}\left(h_{k,n}\right) = m_{h_{k,n}\to f_{h_{k,n}}}\left(h_{k,n}\right) f_{h_{k,n}}$$
$$\propto \mathcal{N}\left(h_{k,n}; \hat{h}_{k,n}, \nu_{h_{k,n}}\right), \tag{50}$$

with

$$\nu_{h_{k,n}} = \frac{\vec{\nu}_{h_{k,n}}\rho_h}{\rho_h + \vec{\nu}_{h_{k,n}}}, \tag{51}$$

$$\hat{h}_{k,n} = \frac{\vec{h}_{k,n}\rho_h}{\rho_h + \vec{\nu}_{h_{k,n}}}. \tag{52}$$

This is the end of forward message passing.

## B. Backward message passing

Next, we elaborate the backward message passing. The backward message from $h_{k,n}$ to $f_{\tilde{s}_{m,k,n}}$ can be expressed as

$$m_{h_{k,n}\to f_{\tilde{s}_{m,k,n}}}\left(h_{k,n}\right) = \frac{\mathfrak{b}\left(h_{k,n}\right)}{m_{f_{\tilde{s}_{m,k,n}}\to h_{k,n}}\left(h_{k,n}\right)}, \tag{53}$$

They are represented collectively as $m_{\boldsymbol{h}_n\to f_{\tilde{\boldsymbol{s}}_n}}\left(\boldsymbol{h}_n\right)$, which is Gaussian with mean $\overleftarrow{\boldsymbol{h}}_n$ and variance $\mathbf{Diag}(\overleftarrow{\boldsymbol{\nu}}_{\boldsymbol{h}_n})$. With the factor graph shown in Fig. 4, the mean and variance can be computed as

$$\overleftarrow{\boldsymbol{\nu}}_{\boldsymbol{h}_n} = \left((\mathbf{1}_K./\boldsymbol{\nu}_{\boldsymbol{h}_n})\otimes\mathbf{1}_M - \mathbf{1}_J./\overrightarrow{\boldsymbol{\nu}}_{\boldsymbol{h}_n}\right)^{.-1}, \tag{54}$$

$$\overleftarrow{\boldsymbol{h}}_n = \overleftarrow{\boldsymbol{\nu}}_{\boldsymbol{h}_n} \cdot \left((\hat{\boldsymbol{h}}_n./\boldsymbol{\nu}_{\boldsymbol{h}_n})\otimes\mathbf{1}_M - \overrightarrow{\boldsymbol{h}}_n./\overrightarrow{\boldsymbol{\nu}}_{\boldsymbol{h}_n}\right), \tag{55}$$

where $\boldsymbol{\nu}_{\boldsymbol{h}_n} = [\nu_{h_{1,n}},\ldots,\nu_{h_{K,n}}]^T$, $\hat{\boldsymbol{h}}_n = [\hat{h}_{1,n},\ldots,\hat{h}_{K,n}]^T$, $\left[\overrightarrow{\boldsymbol{\nu}}_{\boldsymbol{h}_n}\right]_{(k-1)M+m} = \overrightarrow{\nu}_{h_{m,k,n}}$ and $\left[\overrightarrow{\boldsymbol{h}}_n\right]_{(k-1)M+m} = \overrightarrow{h}_{m,k,n}$.

Similarly, the message $m_{\boldsymbol{g}_n\to f_{\boldsymbol{s}_n}}\left(\boldsymbol{g}_n\right)$ is also Gaussian with mean $\overleftarrow{\boldsymbol{g}}_n$ and variance $\mathbf{Diag}(\overleftarrow{\boldsymbol{\nu}}_{\boldsymbol{g}_n})$, which can be computed as

$$\overleftarrow{\boldsymbol{\nu}}_{\boldsymbol{g}_n} = \left(\mathbf{1}_K\otimes(\mathbf{1}_M./\boldsymbol{\nu}_{\boldsymbol{g}_n}) - \mathbf{1}_J./\overrightarrow{\boldsymbol{\nu}}_{\boldsymbol{g}_n}\right)^{.-1}, \tag{56}$$

$$\overleftarrow{\boldsymbol{g}}_n = \overleftarrow{\boldsymbol{\nu}}_{\boldsymbol{g}_n} \cdot \left(\mathbf{1}_K\otimes(\hat{\boldsymbol{g}}_n./\boldsymbol{\nu}_{\boldsymbol{g}_n}) - \overrightarrow{\boldsymbol{g}}_n./\overrightarrow{\boldsymbol{\nu}}_{\boldsymbol{g}_n}\right), \tag{57}$$

where $\boldsymbol{\nu}_{\boldsymbol{g}_n} = [\nu_{g_{1,n}},\ldots,\nu_{g_{M,n}}]^T$, $\hat{\boldsymbol{g}}_n = [\hat{g}_{1,n},\ldots,\hat{g}_{M,n}]^T$, $\left[\overrightarrow{\boldsymbol{\nu}}_{\boldsymbol{g}_n}\right]_{(m-1)K+k} = \overrightarrow{\nu}_{g_{m,k,n}}$ and $\left[\overrightarrow{\boldsymbol{g}}_n\right]_{(m-1)K+k} = \overrightarrow{g}_{m,k,n}$. Then, the backward message $m_{f_{\tilde{\boldsymbol{s}}_n}\to\tilde{\boldsymbol{s}}_n}(\tilde{\boldsymbol{s}}_n) = \mathcal{N}(\tilde{\boldsymbol{s}}_n; \overleftarrow{\tilde{\boldsymbol{s}}}_n, \overleftarrow{\boldsymbol{\nu}}_{\tilde{\boldsymbol{s}}_n})$ with

$$\overleftarrow{\tilde{\boldsymbol{s}}}_n = \overleftarrow{\boldsymbol{h}}_n \cdot \overleftarrow{\boldsymbol{g}}_n, \tag{58}$$

$$\overleftarrow{\boldsymbol{\nu}}_{\tilde{\boldsymbol{s}}_n} = |\overleftarrow{\boldsymbol{h}}_n|^2 \cdot \overleftarrow{\boldsymbol{\nu}}_{\boldsymbol{g}_n} + \overleftarrow{\boldsymbol{\nu}}_{\boldsymbol{h}_n} \cdot |\overleftarrow{\boldsymbol{g}}_n|^2 + \overleftarrow{\boldsymbol{\nu}}_{\boldsymbol{h}_n} \cdot \overleftarrow{\boldsymbol{\nu}}_{\boldsymbol{g}_n}, \tag{59}$$

where $\overleftarrow{\tilde{\boldsymbol{s}}}_n = [\overleftarrow{\tilde{\boldsymbol{s}}}_{1,n}^T,\ldots,\overleftarrow{\tilde{\boldsymbol{s}}}_{K,n}^T]^T$ and $\overleftarrow{\boldsymbol{\nu}}_{\tilde{\boldsymbol{s}}_n} = [\overleftarrow{\boldsymbol{\nu}}_{\tilde{\boldsymbol{s}}_{1,n}}^T,\ldots,\overleftarrow{\boldsymbol{\nu}}_{\tilde{\boldsymbol{s}}_{K,n}}^T]^T$. Following (U)AMP, the backward message is combined with message $m_{\tilde{\boldsymbol{s}}_n\to f_{\tilde{\boldsymbol{s}}_n}}\left(\tilde{\boldsymbol{s}}_n\right)$, i.e.,

$$\boldsymbol{\nu}_{\tilde{\boldsymbol{s}}_{k,n}} = \left(1/\nu_{\boldsymbol{q}_{k,n}}\mathbf{1}_M + \mathbf{1}_M./\overleftarrow{\boldsymbol{\nu}}_{\tilde{\boldsymbol{s}}_{k,n}}\right)^{.-1}, \tag{60}$$

$$\hat{\tilde{\boldsymbol{s}}}_{k,n} = \boldsymbol{\nu}_{\tilde{\boldsymbol{s}}_{k,n}} \cdot \left(1/\nu_{\tilde{\boldsymbol{s}}_{k,n}}\hat{\boldsymbol{q}}_{k,n} + \overleftarrow{\tilde{\boldsymbol{s}}}_{k,n}./\overleftarrow{\boldsymbol{\nu}}_{\tilde{\boldsymbol{s}}_{k,n}}\right), \tag{61}$$

Stack $\boldsymbol{\nu}_{\tilde{\boldsymbol{s}}_{k,n}}$ and $\hat{\tilde{\boldsymbol{s}}}_{k,n}$ into

$$\boldsymbol{\nu}_{\tilde{\boldsymbol{S}}} = [\boldsymbol{\nu}_{\tilde{\boldsymbol{s}}_1},\ldots,\boldsymbol{\nu}_{\tilde{\boldsymbol{s}}_N}], \tag{62}$$

$$\hat{\tilde{\boldsymbol{S}}} = [\hat{\tilde{\boldsymbol{s}}}_1,\ldots,\hat{\tilde{\boldsymbol{s}}}_N], \tag{63}$$

where

$$\boldsymbol{\nu}_{\tilde{\boldsymbol{s}}_n} = [\boldsymbol{\nu}_{\tilde{\boldsymbol{s}}_{1,n}}^T,\ldots,\boldsymbol{\nu}_{\tilde{\boldsymbol{s}}_{K,n}}^T]^T, \tag{64}$$

$$\hat{\tilde{\boldsymbol{s}}}_n = [\hat{\tilde{\boldsymbol{s}}}_{1,n}^T,\ldots,\hat{\tilde{\boldsymbol{s}}}_{K,n}^T]^T. \tag{65}$$

Further, we can obtain $\boldsymbol{\nu_S} = \boldsymbol{\nu}_{\tilde{\boldsymbol{S}}}^T$ and $\hat{\boldsymbol{S}} = \hat{\tilde{\boldsymbol{S}}}^T$ with

$$\boldsymbol{\nu_S} = [\boldsymbol{\nu_{s_1}}, \ldots, \boldsymbol{\nu_{s_J}}], \tag{66}$$

$$\nu_{\boldsymbol{s}_j} = \langle \boldsymbol{\nu_{s_j}} \rangle, \tag{67}$$

$$\hat{\boldsymbol{S}} = [\hat{\boldsymbol{s}}_1, \ldots, \hat{\boldsymbol{s}}_J], \tag{68}$$

$$\hat{\boldsymbol{s}}_j = [\hat{s}_{j,1}, \ldots, \hat{s}_{j,N}], \tag{69}$$

According to the belief propagation derivation of (U)AMP, it holds that

$$m_{\boldsymbol{z}_j \to f_{\boldsymbol{r}_j}}(\boldsymbol{z}_j) = m_{f_{\boldsymbol{z}_j} \to \boldsymbol{z}_j}(\boldsymbol{z}_j) = \mathcal{N}(\boldsymbol{z}_j; \boldsymbol{p}_j, \mathbf{Diag}(\boldsymbol{\nu_{p_j}})), \tag{70}$$

where $\boldsymbol{p}_j$ and $\boldsymbol{\nu_{p_j}}$ are computed in the forward direction. It is noted that the factor node $f_{\boldsymbol{r}_j}$ connects the variable node $\beta$. According to the rules of the variational message passing, the message $m_{f_{\boldsymbol{r}_j} \to \beta}(\beta)$ can be expressed as

$$m_{f_{\boldsymbol{r}_j} \to \beta}(\beta) \propto \exp \left\{ \sum_{j=1}^{J} \int_{\boldsymbol{z}_j} \mathfrak{b}(\boldsymbol{z}_j) \log f_{\boldsymbol{r}_j} \right\}, \tag{71}$$

where $\mathfrak{b}(\boldsymbol{z}_j)$ is the approximate marginal of $\boldsymbol{z}_j$, which can be expressed as

$$\begin{aligned} \mathfrak{b}(\boldsymbol{z}_j) &\propto m_{f_{\boldsymbol{r}_j} \to \boldsymbol{z}_j}(\boldsymbol{z}_j) m_{\boldsymbol{z}_j \to f_{\boldsymbol{r}_j}}(\boldsymbol{z}_j) \\ &= \mathcal{N}\left(\boldsymbol{z}_j; \hat{\boldsymbol{z}}_j, \mathbf{Diag}\left(\boldsymbol{\nu_{z_j}}\right)\right), \end{aligned} \tag{72}$$

where

$$\boldsymbol{\nu_{z_j}} = \mathbf{1}_L ./ \left( \mathbf{1}_L ./ \boldsymbol{\nu_{p_j}} + \hat{\beta}\mathbf{1}_L \right), \tag{73}$$

$$\hat{\boldsymbol{z}}_j = \boldsymbol{\nu_{z_j}} \cdot \left( \boldsymbol{p}_j ./ \boldsymbol{\nu_{p_j}} + \hat{\beta}\boldsymbol{r}_j \right), \tag{74}$$

with $\hat{\beta}$ being the approximate a posteriori mean of the noise precision that is obtained with (13). It is noted that in the above derivation, the message $m_{f_{\boldsymbol{r}_j} \to \boldsymbol{z}_j}(\boldsymbol{z}_j)$ is required, which is Gaussian, i.e., $m_{f_{\boldsymbol{r}_j} \to \boldsymbol{z}_j}(\boldsymbol{z}_j) = \mathcal{N}(\boldsymbol{z}_j; \boldsymbol{r}_j, \hat{\beta}^{-1})$, and its derivation is shown in (11). Then, it is not hard to show that the message

$$m_{f_{\boldsymbol{r}_j} \to \beta}(\beta) \propto \beta^{LJ} \exp \left\{ \sum_{j=1}^{J} -\beta \left( \|\boldsymbol{r}_j - \hat{\boldsymbol{z}}_j\|^2 + \mathbf{1}_L^T \boldsymbol{\nu_{z_j}} \right) \right\}. \tag{75}$$

This is the end of backward message passing.

The message passing algorithm is summarized in Algorithm 1 and it can be terminated when it reaches a maximum number of iteration or the difference between the estimates of two consecutive iterations is less than a threshold.

### C. Computational Complexity Analysis

We analyze the computational complexity of the proposed algorithm and compare it with that of sate-of-the-art algorithms. The UAMP-based message passing algorithm needs pre-processing, i.e., performing a single economic SVD for $\boldsymbol{\Phi}$ and unitary transformation, and the complexity is $\mathcal{O}(NL^2)$. It noted that the SVD can be carried out offline and there

---

**Algorithm 1** UAMP-Based Channel Estimation for RIS-Aided MIMO System

---

**Input:** A feasible $\boldsymbol{\Phi}$, $\epsilon > 0$ and the maximum number of iteration $I_{max}$.

**Initialize:** $\hat{h}_{k,n}, \nu_{h_{k,n}} = 1, \hat{\boldsymbol{s}}_j = \mathbf{0}, \nu_{\boldsymbol{s}_j} = 1, \boldsymbol{\mu}_j = \mathbf{0}, \forall k, n, j,$ and $\hat{\beta} = 1$.

**Repeat:**

1: update noise precision $\hat{\beta}$ with (13);
2: $\forall j$: update $\boldsymbol{\nu_{p_j}}$ and $\boldsymbol{p}_j$ with (16) and (17);
3: $\forall j$: update $\boldsymbol{\nu_{\mu_j}}$ and $\boldsymbol{\mu}_j$ with (18) and (19);
4: $\forall j$: update $\boldsymbol{\nu_{q_j}}$ and $\boldsymbol{q}_j$ with (20) and (21);
5: $\forall n$: update $\boldsymbol{\nu}'_{\tilde{\boldsymbol{q}}_n}$ and $\tilde{\boldsymbol{q}}_n$ with (30) and (27);
6: $\forall m, k, n$: update $\vec{\nu}_{g_{m,k,n}}$ and $\vec{g}_{m,k,n}$ with (36) and (37);
7: $\forall m, n$: update $\vec{\nu}_{g_{m,n}}$ and $\vec{g}_{m,n}$ with (39) and (40);
8: $\forall m, n$: update $\nu_{g_{m,n}}$ and $\hat{g}_{m,n}$ with (42) and (43);
9: $\forall m, k, n$: update $\vec{\nu}_{h_{m,k,n}}$ and $\vec{h}_{m,k,n}$ with (45) and (46);
10: $\forall k, n$: update $\vec{\nu}_{h_{k,n}}$ and $\vec{h}_{k,n}$ with (48) and (49);
11: $\forall k, n$: update $\nu_{h_{k,n}}$ and $\hat{h}_{k,n}$ with (51) and (52);
12: $\forall n$: update $\overleftarrow{\nu}_{\boldsymbol{h}_n}$ and $\overleftarrow{\boldsymbol{h}}_n$ with (54) and (55);
13: $\forall n$: update $\overleftarrow{\nu}_{\boldsymbol{g}_n}$ and $\overleftarrow{\boldsymbol{g}}_n$ with (56) and (57);
14: $\forall n$: update $\overleftarrow{\nu}_{\tilde{\boldsymbol{s}}_n}$ and $\overleftarrow{\tilde{\boldsymbol{s}}}_n$ with (59) and (58);
15: $\forall k, n$: update $\nu_{\tilde{s}_{k,n}}$ and $\hat{\tilde{s}}_{k,n}$ with (60) and (61);
16: $\forall j$: update $\nu_{\boldsymbol{s}_j}$ and $\hat{\boldsymbol{s}}_j$ with (67) and (69);
17: $\forall j$: update $\boldsymbol{\nu_{z_j}}$ and $\hat{\boldsymbol{z}}_j$ with (73) and (74);
18: Construct $\hat{\boldsymbol{H}} = [\hat{\boldsymbol{h}}_1, \ldots, \hat{\boldsymbol{h}}_N]^T$ with $\hat{\boldsymbol{h}}_n = [\hat{h}_{1,n}, \ldots, \hat{h}_{K,n}]^T$ and $\hat{\boldsymbol{G}} = [\hat{\boldsymbol{g}}_1, \ldots, \hat{\boldsymbol{g}}_N]$ with $\hat{\boldsymbol{g}}_n = [\hat{g}_{1,n}, \ldots, \hat{g}_{M,n}]^T$.

**Until** $\|\hat{\boldsymbol{H}} - \boldsymbol{H}\|_F^2 \|\boldsymbol{H}\|_F^{-2} < \epsilon$ and $\|\hat{\boldsymbol{G}} - \boldsymbol{G}\|_F^2 \|\boldsymbol{G}\|_F^{-2} < \epsilon$ or the number of iteration is more than $I_{max}$.

**Output:** $\hat{\beta}$, $\hat{\boldsymbol{H}}$ and $\hat{\boldsymbol{G}}$ that are the estimations of $\beta$, $\boldsymbol{H}$ and $\boldsymbol{G}$, respectively.

---

is no matrix inversion involved in Algorithm 1. Also note that the formulated problem is a multiple measurement vector one. The complexity of the proposed algorithm is dominated by the computation of $\boldsymbol{p}$ in step 1, which requires $\mathcal{O}(NLKM + LKM)$, and the computations of $\nu_q$ and $\boldsymbol{q}$ in step 6, which require $\mathcal{O}(NLKM + NKM)$ and $\mathcal{O}(NLKM)$, respectively. It can be shown that the overall complexity of the algorithm is $\mathcal{O}(NLKM) + \mathcal{O}(NKM)$ per iteration. As we consider a general RIS-aided MIMO system without any special requirements on the (channel or RIS phase) matrices, the most relevant algorithms for comparison are the ALS-based algorithm and VAMP-based algorithm in [29]. The complexity of the ALS-based algorithm is $\mathcal{O}(NLKM + L(M+K)N^2) + \mathcal{O}(N^3)$ per iteration, and that of the VAMP-based algorithm is $\mathcal{O}((M+K)N^3 + (M+K)N^2)$ per iteration. It is noted that $N$ is the number of RIS units, which can be much larger than $M$, $K$ and $L$. From the analysis, we can

see that the complexity of the proposed algorithm, which is linear with $N$, is significantly smaller than that of the ALS or VAMP-based algorithm.

### TABLE II: Complexity Comparison

| Algorithm | Complexity |
|---|---|
| UAMP | $\mathcal{O}(NLKM) + \mathcal{O}(NKM)$ |
| ALS | $\mathcal{O}(NLKM + L(M+K)N^2) + \mathcal{O}(N^3)$ |
| VAMP | $\mathcal{O}((M+K)N^3 + (M+K)N^2)$ |

## V. CRAMÉR-RAO LOWER BOUND

In this section, we derive the CRLB for the RIS channel estimation, which is used to serve as another performance benchmark, besides the ALS and VAMP-based algorithms.

We firstly rewrite the system model as

$$\tilde{Y} = (H^T \odot G)\Phi^T + \tilde{W}, \tag{76}$$

and define a complex parameter $\theta \in \mathbb{C}^{2N(M+K)\times 1}$, which includes all of unknown complex parameters in $H$ and $G$ as

$$\theta \triangleq \left[ h_1^T, \ldots, h_N^T, g_1^T, \ldots, g_N^T, h_1^H, \ldots, h_N^H, g_1^H, \ldots, g_N^H \right]^T. \tag{77}$$

The likelihood function of $\tilde{Y}$ can be expressed as

$$p(\tilde{Y};\theta) = (\pi\sigma^2)^{-KML} \exp\left\{ -\sigma^{-2} \sum_{l=1}^{L} \left\| \tilde{y}_l - (H^T \odot G)\Phi_{l,:}^T \right\|^2 \right\}, \tag{78}$$

and the logarithm of likelihood function can be expressed as

$$\ln(p(\tilde{Y};\theta)) = -KML \ln(\pi\sigma^2) - \sigma^{-2} \sum_{l=1}^{L} \left\| \tilde{y}_l - (H^T \odot G)\Phi_{l,:}^T \right\|^2, \tag{79}$$

Then, define $f_\theta \triangleq \ln(p(\tilde{Y};\theta))$ and the Fisher information matrix (FIM) $\mathcal{J}_\theta \in \mathbb{C}^{(2N(M+K))\times(2N(M+K))}$ can be obtained by

$$\mathcal{J}_\theta = \mathbb{E}\left\{ \left( \frac{\partial f_\theta}{\partial \theta} \right) \left( \frac{\partial f_\theta}{\partial \theta} \right)^H \right\}. \tag{80}$$

The partial derivatives of $f_\theta$ with respect to $\theta$ can be expressed as

$$\frac{\partial f_\theta}{\partial h_{k,n}} = \sigma^{-2} \sum_{l=1}^{L} \left\{ [\Phi_{l,:}^T]_n g_n^T \left( ([Y_l^{'}]_{:,k}) \right)^* - ((H^T)_{k,:} \odot G)^* (\Phi_{l,:}^T)^* \right\}, \tag{81}$$

$$\frac{\partial f_\theta}{\partial g_{m,n}} = \sigma^{-2} \sum_{l=1}^{L} \left\{ [\Phi_{l,:}^T]_n h_n^T \left( (Y_l^{'})_{m,:} \right)^H - (H^T \odot G_{m,:})^* (\Phi_{l,:}^T)^* \right\}, \tag{82}$$

$$\frac{\partial f_\theta}{\partial h_{k,n}^*} = \left( \frac{\partial f_\theta}{\partial h_{k,n}} \right)^*, \qquad \frac{\partial f_\theta}{\partial g_{m,n}^*} = \left( \frac{\partial f_\theta}{\partial g_{m,n}} \right)^*. \tag{83}$$

where $[\Phi_{l,:}^T]_n$ is the $n$-th element of $\Phi_{l,:}^T$, $Y_l^{'} \triangleq [\tilde{y}_{1,l}, \ldots, \tilde{y}_{K,l}] \in \mathbb{C}^{M\times K}$ and $\tilde{y}_{k,l} \in \mathbb{C}^{M\times 1}$ represents the sub-vectors of $\tilde{y}_l$ as

$$\tilde{y}_l = [\tilde{y}_{1,l}^T, \ldots, \tilde{y}_{K,l}^T]^T. \tag{84}$$

Hence, the FIM $\mathcal{J}_\theta$ can be expressed as

$$\mathcal{J}_\theta = \left[ \begin{array}{cc} \mathcal{P} & 0 \\ 0 & \mathcal{P}^* \end{array} \right], \tag{85}$$

where $\mathcal{P}$ is shown in (86) at the top of next page and the size of $\mathcal{P}$ and $0$ are $(N(M+K)) \times (N(M+K))$. The inverse of the FIM of $\theta$ gives, under some regularity conditions, a lower bound for the augmented covariance matrix of an unbiased estimator of $\theta$ as

$$\mathcal{J}_\theta^{-1} = \left[ \begin{array}{cc} \mathcal{P}^{-1} & 0 \\ 0 & (\mathcal{P}^{-1})^* \end{array} \right], \tag{87}$$

where $\mathcal{P}^{-1}$ is given by

$$\mathcal{P}^{-1} = \left[ \begin{array}{cc} \Omega_H & \mathcal{I} \\ \mathcal{I}^H & \Omega_G \end{array} \right], \tag{88}$$

and $\Omega_H \in \mathbb{C}^{KN\times KN}$ and $\Omega_G \in \mathbb{C}^{MN\times MN}$ are the CRLB matrices for the estimates of $H$ and $G$, respectively, and $\mathcal{I} \in \mathbb{C}^{KN\times MN}$ represents the remaining sub-matrices.

Furthermore, $\mathcal{P}$ can be divided into four sub-matrices as

$$\mathcal{P} = \left[ \begin{array}{cc} \mathcal{P}_{HH} & \mathcal{P}_{HG} \\ \mathcal{P}_{HG}^H & \mathcal{P}_{GG} \end{array} \right], \tag{89}$$

where $\mathcal{P}_{HH} \in \mathbb{C}^{KN\times KN}$, $\mathcal{P}_{GG} \in \mathbb{C}^{MN\times MN}$ and $\mathcal{P}_{HG} \in \mathbb{C}^{KN\times MN}$, as shown in (86). According to the formula for inverse of a partitioned Hermitian matrix in [36], we can obtain $\Omega_H$ and $\Omega_G$ as following

$$\Omega_H = (\mathcal{P}_{HH} - \mathcal{P}_{HG}\mathcal{P}_{GG}^{-1}\mathcal{P}_{HG}^H)^{-1}, \tag{90}$$

$$\Omega_G = (\mathcal{P}_{GG} - \mathcal{P}_{HG}^H\mathcal{P}_{HH}^{-1}\mathcal{P}_{HG})^{-1}. \tag{91}$$

So the CRLB of $H$ and $G$ can be donated as

$$\text{CRLB}_H = \frac{\text{trace}(\Omega_H)}{KN}, \tag{92}$$

$$\text{CRLB}_G = \frac{\text{trace}(\Omega_G)}{MN}. \tag{93}$$

## VI. SIMULATION RESULTS

In this section, we provide extensive numerical experiments to demonstrate the superior performance of the proposed UAMP-based channel estimation algorithm. For comparison, we also include the ALS-based and VAMP-based channel estimation algorithms in [29]. The threshold $\epsilon = 10^{-3}$ and the maximum number of iterations is set to 30. The entries of $H$ and $G$ are independently drawn from a complex Gaussian distribution with zero mean and unit variance. The scaling ambiguity of the estimation is eliminated in the calculation of the normalized mean square error (NMSE). The SNR (in dB) is defined as

$$\text{SNR} = 10 \log_{10} \left( \frac{N\mathbb{E}\left\{ \|\Phi\|_F^2 \right\}}{L\beta^{-1}} \right). \tag{94}$$

$$\mathcal{P} = \begin{bmatrix} \mathbb{E}\left\{\left(\frac{\partial f}{\partial h_{1,1}}\right)\left(\frac{\partial f}{\partial h_{1,1}^*}\right)\right\} & \cdots & \mathbb{E}\left\{\left(\frac{\partial f}{\partial h_{1,1}}\right)\left(\frac{\partial f}{\partial h_{K,N}^*}\right)\right\} & \mathbb{E}\left\{\left(\frac{\partial f}{\partial h_{1,1}}\right)\left(\frac{\partial f}{\partial g_{1,1}^*}\right)\right\} & \cdots & \mathbb{E}\left\{\left(\frac{\partial f}{\partial h_{1,1}}\right)\left(\frac{\partial f}{\partial g_{M,N}^*}\right)\right\} \\ \vdots & & \vdots & \vdots & & \vdots \\ \mathbb{E}\left\{\left(\frac{\partial f}{\partial h_{K,N}}\right)\left(\frac{\partial f}{\partial h_{1,1}^*}\right)\right\} & \cdots & \mathbb{E}\left\{\left(\frac{\partial f}{\partial h_{K,N}}\right)\left(\frac{\partial f}{\partial h_{K,N}^*}\right)\right\} & \mathbb{E}\left\{\left(\frac{\partial f}{\partial h_{K,N}}\right)\left(\frac{\partial f}{\partial g_{1,1}^*}\right)\right\} & \cdots & \mathbb{E}\left\{\left(\frac{\partial f}{\partial h_{K,N}}\right)\left(\frac{\partial f}{\partial g_{M,N}^*}\right)\right\} \\ \mathbb{E}\left\{\left(\frac{\partial f}{\partial g_{1,1}}\right)\left(\frac{\partial f}{\partial h_{1,1}^*}\right)\right\} & \cdots & \mathbb{E}\left\{\left(\frac{\partial f}{\partial g_{1,1}}\right)\left(\frac{\partial f}{\partial h_{K,N}^*}\right)\right\} & \mathbb{E}\left\{\left(\frac{\partial f}{\partial g_{1,1}}\right)\left(\frac{\partial f}{\partial g_{1,1}^*}\right)\right\} & \cdots & \mathbb{E}\left\{\left(\frac{\partial f}{\partial g_{1,1}}\right)\left(\frac{\partial f}{\partial g_{M,N}^*}\right)\right\} \\ \vdots & & \vdots & \vdots & & \vdots \\ \mathbb{E}\left\{\left(\frac{\partial f}{\partial g_{M,N}}\right)\left(\frac{\partial f}{\partial h_{1,1}^*}\right)\right\} & \cdots & \mathbb{E}\left\{\left(\frac{\partial f}{\partial g_{M,N}}\right)\left(\frac{\partial f}{\partial h_{K,N}^*}\right)\right\} & \mathbb{E}\left\{\left(\frac{\partial f}{\partial g_{M,N}}\right)\left(\frac{\partial f}{\partial g_{1,1}^*}\right)\right\} & \cdots & \mathbb{E}\left\{\left(\frac{\partial f}{\partial g_{M,N}}\right)\left(\frac{\partial f}{\partial g_{M,N}^*}\right)\right\} \end{bmatrix}. \tag{86}$$

In order to estimate the RIS channels during the training phase, the RIS goes through $L$ configurations, which leads to the RIS phase matrix $\boldsymbol{\Phi}$ in (6). In the simulations, we consider two methods. One is that the RIS units are turned on or off randomly, leading to a matrix $\boldsymbol{\Phi}$ with entries 1 or 0 [37], which is called a binary matrix hereafter. In the simulations we assume that each entry in $\boldsymbol{\Phi}$ takes 1 or 0 with the same probability. The other one is that the phases of the RIS units are set to some discrete values, and in particular, the phase matrix $\boldsymbol{\Phi}$ is part of the DFT matrix (called partial DFT matrix), as in [22].

We evaluate the performance of estimators in terms of the NMSE of estimated channel matrices $\hat{\boldsymbol{H}}$ and $\hat{\boldsymbol{G}}$. The NMSE performance of various estimators versus SNR with $L = N = K = M = 64$ is shown in Fig. 5, where Fig. 5 (a) and (b) are for partial DFT matrix and binary phase matrix, respectively. It is observed that the proposed UAMP-based algorithm significantly outperforms the ALS-based and VAMP-based channel estimators for both $\boldsymbol{H}$ and $\boldsymbol{G}$, especially in the case that $\boldsymbol{\Phi}$ is a binary matrix. Note that the NMSEs of $\boldsymbol{H}$ and $\boldsymbol{G}$ are very similar. To keep the figures clear, we only show the NMSE performance of $\boldsymbol{H}$ in the subsequent simulation results.

In Fig. 6, with $N = K = M = 64$, we compare the NMSE performance versus SNR of the estimators with different values of $L$. According to the results, the performance of the UAMP-based method is significantly better than that of the ALS and VAMP-based methods, especially when $L$ is relatively small. As smaller $L$ (the number of RIS phase configurations needed for channel estimation) is highly desirable to reduce the training overhead and latency, next we vary the value of $L$ and examine the performance of the estimators. The results are shown in Fig. 7, where the SNR is set to 20dB, and $\boldsymbol{\Phi}$ is partial DFT matrix in (a) and binary matrix in (b). It can be seen that, with the increase of $L$, the performance of all estimators improves, as expected. However, the performance improvement of the ALS and VAMP-based estimators is very slow with $L$. We can also see from Fig. 7 (a) that, the UAMP-based estimator with $L = 16$ achieves the same performance of ALS and VAMP-based estimators with $L = 32$. According to Fig. 7 (b), the UAMP-based estimator with $L = 16$ even outperforms the ALS and VAMP-based estimators with $L = 32$. The results demonstrate that the use of the proposed algorithm can lead to a huge reduction in training overheads.



Fig. 5: NMSE performance of the estimators versus SNR, where $L = N = K = M = 64$. (a) Partial DFT matrix; (b) Binary matrix.

We next examine the impact of the number of RIS units $N$ on the performance of channel estimation, where we set $L = K = M = 32$. As a smaller $N$ leads to a less number of channel coefficients to be estimated, for a fixed $L$ the performance of the estimators improves with the decrease of $N$. In Fig. 8, we can find that, in the case of binary matrix, the performance of the UAMP-based estimator with $N = 128$ is even better than that of the ALS and VAMP-based estimators with $N = 32$, which again demonstrates the superior performance of the proposed one. The performance of

Fig. 6: NMSE performance of the estimators versus SNR for $N = K = M = 64$ and different $L$ (a) Partial DFT matrix; (b) Binary matrix.



Fig. 7: NMSE performance of the estimators versus $L$, where $N = K = M = 64$ and SNR=20dB. (a) Partial DFT matrix; (b) Binary matrix.

the UAMP-based algorithm, the ALS algorithm and the CRLB are shown in Fig. 9, where $L = K = M = N = 16$. We can see that the performance of UAMP-based algorithm is almost the same as the CRLB, which is significantly better than that of the ALS algorithm.

From the above results, we find that the performance of the proposed UAMP-based algorithm consistently show good performance for both partial DFT matrix and binary matrix. In contrast, the ALS and VAMP-based algorithms exhibit significantly worse performance in the case of binary matrix, compared to the partial DFT matrix.

We also investigate the performance of the estimation with iteration number for different $\Phi$ when the SNR is 20dB. The results are shown in Fig. 10. It can be seen that the proposed algorithm converges fast in different cases, especially when $\Phi$ is a partial DFT matrix. It is worth mentioning that the estimation of the noise precision (reciprocal of the variance) is incorporated in the UAMP-based channel estimator, so no separate noise power estimator is needed. In Fig. 11, we compare the estimated noise power and its true value, where $\Phi$ is a partial DFT matrix. In Fig. 11 (a), the estimate of the noise variance with the iteration number is shown, where the SNR is 20dB. We can see that the convergence is fast.

The results in Fig. 11 (b) show that the proposed algorithm provides accurate noise variance estimates for a wide range of SNRs.

## VII. Conclusions

In this paper, we have addressed the issue of channel estimation in RIS-aided MIMO communications. Through vectorization and reduction, we obtain a new signal model for channel estimation, based on which a message passing based algorithm is developed, leveraging UAMP. Compared to the state-of-the-art algorithms, the proposed algorithm does not have any special requirements on the matrices involved, and it shows significant advantages in computational complexity, estimation performance and training overhead. Extensive numerical results demonstrate the merits of the proposed algorithm.

## References

[1] C. Huang, A. Zappone, G. C. Alexandropoulos, M. Debbah, and C. Yuen, "Reconfigurable intelligent surfaces for energy efficiency in wireless communication," *IEEE Transactions on Wireless Communications*, vol. 18, no. 8, pp. 4157–4170, 2019.

Fig. 8: NMSE performance of $\hat{H}$ versus SNR with different $N$, where $L = K = M = 32$. (a) Partial DFT matrix; (b) Binary matrix.



Fig. 9: Performance of the estimators and CRLB, where $L = K = M = N = 16$. (a) Partial DFT matrix; (b) Binary matrix.

[2] X. Guan, Q. Wu, and R. Zhang, "Intelligent reflecting surface assisted secrecy communication via joint beamforming and jamming," *CoRR*, vol. abs/1907.12839, 2019. [Online]. Available: http://arxiv.org/abs/1907.12839

[3] Q. Nadeem, A. Kammoun, A. Chaaban, M. Debbah, and M. Alouini, "Asymptotic analysis of large intelligent surface assisted MIMO communication," *CoRR*, vol. abs/1903.08127, 2019. [Online]. Available: http://arxiv.org/abs/1903.08127

[4] Y. Liang, R. Long, Q. Zhang, J. Chen, H. V. Cheng, and H. Guo, "Large intelligent surface/antennas (LISA): making reflective radios smart," *CoRR*, vol. abs/1906.06578, 2019. [Online]. Available: http://arxiv.org/abs/1906.06578

[5] Q. Wu and R. Zhang, "Towards smart and reconfigurable environment: Intelligent reflecting surface aided wireless network," *IEEE Communications Magazine*, vol. 58, no. 1, pp. 106–112, 2020.

[6] Z. Yang, M. Chen, W. Saad, W. Xu, M. Shikh-Bahaei, H. V. Poor, and S. Cui, "Energy-efficient wireless communications with distributed reconfigurable intelligent surfaces," *IEEE Transactions on Wireless Communications*, pp. 1–1, 2021.

[7] H. Guo, Y.-C. Liang, J. Chen, and E. G. Larsson, "Weighted sum-rate maximization for intelligent reflecting surface enhanced wireless networks," in *2019 IEEE Global Communications Conference (GLOBECOM)*, 2019, pp. 1–6.

[8] Y. Liang, J. Yang, W. Xie, M. Hasna, and M. Renzo, "Secrecy performance analysis of RIS-aided wireless communication systems," *IEEE Transactions on Vehicular Technology*, vol. PP, no. 99, pp. 1–1, 2020.

[9] P. Wang, J. Fang, X. Yuan, Z. Chen, and H. Li, "Intelligent reflecting surface-assisted millimeter wave communications: Joint active and passive precoding design," *IEEE Transactions on Vehicular Technology*, vol. 69, no. 12, pp. 14 960–14 973, 2020.

[10] E. Basar, "Reconfigurable intelligent surface-based index modulation: A new beyond mimo paradigm for 6G," *IEEE Transactions on Communications*, vol. 68, no. 5, pp. 3187–3196, 2020.

[11] C. Huang, S. Hu, G. C. Alexandropoulos, A. Zappone, C. Yuen, R. Zhang, M. D. Renzo, and M. Debbah, "Holographic MIMO surfaces for 6G wireless networks: Opportunities, challenges, and trends," *IEEE Wireless Communications*, vol. 27, no. 5, pp. 118–125, 2020.

[12] A. Zappone, M. Di Renzo, F. Shams, X. Qian, and M. Debbah, "Overhead-aware design of reconfigurable intelligent surfaces in smart radio environments," *IEEE Transactions on Wireless Communications*, vol. 20, no. 1, pp. 126–141, 2021.

[13] L. Yang, Y. Yang, M. O. Hasna, and M.-S. Alouini, "Coverage, probability of SNR gain, and DOR analysis of RIS-aided communication systems," *IEEE Wireless Communications Letters*, vol. 9, no. 8, pp. 1268–1272, 2020.

[14] Y. Li, M. Jiang, Q. Zhang, and J. Qin, "Joint beamforming design in multi-cluster MISO NOMA reconfigurable intelligent surface-aided downlink communication networks," *IEEE Transactions on Communications*, vol. 69, no. 1, pp. 664–674, 2021.

[15] Z. Zhu, Z. Li, Z. Chu, G. Sun, W. Hao, P. Liu, and I. Lee, "Resource allocation for intelligent reflecting surface assisted wireless powered iot systems with power splitting," *IEEE Transactions on Wireless Communications*, pp. 1–1, 2021.

[16] C. Huang, Z. Yang, G. C. Alexandropoulos, K. Xiong, L. Wei, C. Yuen, Z. Zhang, and M. Debbah, "Multi-hop ris-empowered terahertz communications: A drl-based hybrid beamforming design," *IEEE Journal on Selected Areas in Communications*, vol. 39, no. 6, pp. 1663–1677, 2021.

Fig. 10: NMSE performance of $\hat{H}$ and $\hat{G}$ versus the iteration number, where $L = 20$ and $K = M = N = 32$.



Fig. 11: Noise variance estimation of the proposed algorithm.

[17] T. Hou, Y. Liu, Z. Song, X. Sun, Y. Chen, and L. Hanzo, "Reconfigurable intelligent surface aided noma networks," *IEEE Journal on Selected Areas in Communications*, vol. 38, no. 11, pp. 2575–2588, 2020.

[18] I. F. Akyildiz, C. Han, and S. Nie, "Combating the distance problem in the millimeter wave and terahertz frequency bands," *IEEE Communications Magazine*, vol. 56, no. 6, pp. 102–108, 2018.

[19] C. Pan, H. Ren, K. Wang, W. Xu, M. Elkashlan, A. Nallanathan, and L. Hanzo, "Multicell mimo communications relying on intelligent reflecting surfaces," *IEEE Transactions on Wireless Communications*, vol. 19, no. 8, pp. 5218–5233, 2020.

[20] D. Mishra and H. Johansson, "Channel estimation and low-complexity beamforming design for passive intelligent surface assisted MISO wireless energy transfer," in *ICASSP 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2019, pp. 4659–4663.

[21] B. Zheng and R. Zhang, "Intelligent reflecting surface-enhanced ofdm: Channel estimation and reflection optimization," *IEEE Wireless Communications Letters*, vol. 9, no. 4, pp. 518–522, 2020.

[22] T. L. Jensen and E. De Carvalho, "An optimal channel estimation scheme for intelligent reflecting surfaces based on a minimum variance unbiased estimator," in *ICASSP 2020 - 2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2020, pp. 5000–5004.

[23] Q.-U.-A. Nadeem, H. Alwazani, A. Kammoun, A. Chaaban, M. Debbah, and M.-S. Alouini, "Intelligent reflecting surface-assisted multi-user MISO communication: Channel estimation and beamforming design," *IEEE Open Journal of the Communications Society*, vol. 1, pp. 661–680, 2020.

[24] A. Taha, M. Alrabeiah, and A. Alkhateeb, "Enabling large intelligent surfaces with compressive sensing and deep learning," *IEEE Access*, vol. 9, pp. 44 304–44 321, 2021.

[25] Z.-Q. He and X. Yuan, "Cascaded channel estimation for large intelligent metasurface assisted massive MIMO," *IEEE Wireless Communications Letters*, vol. 9, no. 2, pp. 210–214, 2020.

[26] Z.-Q. He, H. Liu, X. Yuan, Y.-J. A. Zhang, and Y.-C. Liang, "Semi-blind cascaded channel estimation for reconfigurable intelligent surface aided massive mimo," 2021.

[27] H. Liu, X. Yuan, and Y.-J. A. Zhang, "Matrix-calibration-based cascaded channel estimation for reconfigurable intelligent surface assisted multiuser MIMO," *IEEE Journal on Selected Areas in Communications*, vol. 38, no. 11, pp. 2621–2636, 2020.

[28] Z. Wang, L. Liu, and S. Cui, "Channel estimation for intelligent reflecting surface assisted multiuser communications: Framework, algorithms, and analysis," *IEEE Transactions on Wireless Communications*, vol. 19, no. 10, pp. 6607–6620, 2020.

[29] L. Wei, C. Huang, G. C. Alexandropoulos, C. Yuen, Z. Zhang, and M. Debbah, "Channel estimation for RIS-empowered multi-user MISO wireless communications," *IEEE Transactions on Communications*, vol. 69, no. 6, pp. 4144–4157, 2021.

[30] L. Wei, C. Huang, Q. Guo, Z. Zhang, M. Debbah, and C. Yuen, "Bidirectional approximate message passing for ris-assisted multi-user miso communications," in *2021 IEEE 94th Vehicular Technology Conference (VTC2021-Fall)*, 2021, pp. 1–6.

[31] Q. Guo and J. Xi, "Approximate message passing with unitary transformation," *CoRR*, vol. abs/1504.04799, 2015. [Online]. Available: http://arxiv.org/abs/1504.04799

[32] Z. Yuan, Q. Guo, and M. Luo, "Approximate message passing with unitary transformation for robust bilinear recovery," *IEEE Transactions on Signal Processing*, vol. 69, pp. 617–630, 2021.

[33] M. Luo, Q. Guo, M. Jin, Y. C. Eldar, D. Huang, and X. Meng, "Unitary approximate message passing for sparse bayesian learning," *IEEE Transactions on Signal Processing*, vol. 69, pp. 6023–6039, 2021.

[34] D. L. Donoho, A. Maleki, and A. Montanari, "Message passing algorithms for compressed sensing: I. motivation and construction," in *2010 IEEE Information Theory Workshop on Information Theory (ITW 2010, Cairo)*, 2010, pp. 1–5.

[35] J. Dauwels, "On variational message passing on factor graphs," in *2007 IEEE International Symposium on Information Theory*, 2007, pp. 2546–2550.

[36] X. Liu and N. Sidiropoulos, "Cramer-rao lower bounds for low-rank decomposition of multidimensional arrays," *IEEE Transactions on Signal Processing*, vol. 49, no. 9, pp. 2074–2086, 2001.

[37] D. Mishra and H. Johansson, "Channel estimation and low-complexity beamforming design for passive intelligent surface assisted MISO wireless energy transfer," in *ICASSP 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2019, pp. 4659–4663.