

# Deep Transfer Learning Based Downlink Channel Prediction for FDD Massive MIMO Systems

Yuwen Yang, Feifei Gao, Zhimeng Zhong, Bo Ai and Ahmed Alkhateeb

**Abstract**—Artificial intelligence (AI) based downlink channel state information (CSI) prediction for frequency division duplexing (FDD) massive multiple-input multiple-output (MIMO) systems has attracted growing attention recently. However, existing works focus on the downlink CSI prediction for the users under a given environment and is hard to adapt to users in new environment especially when labeled data is limited. To address this issue, we formulate the downlink channel prediction as a deep transfer learning (DTL) problem, where each learning task aims to predict the downlink CSI from the uplink CSI for one single environment. Specifically, we develop the direct-transfer algorithm based on the fully-connected neural network architecture, where the network is trained on the data from all previous environments in the manner of classical deep learning and is then fine-tuned for new environments. To further improve the transfer efficiency, we propose the meta-learning algorithm that trains the network by alternating inner-task and across-task updates and then adapts to a new environment with a small number of labeled data. Simulation results show that the direct-transfer algorithm achieves better performance than the deep learning algorithm, which implies that the transfer learning benefits the downlink channel prediction in new environments. Moreover, the meta-learning algorithm significantly outperforms the direct-transfer algorithm in terms of both prediction accuracy and stability, especially when the number of samples is very small, which validates its effectiveness and superiority.

**Index Terms**—Deep transfer learning (DTL), meta-learning, few-sample learning, downlink CSI prediction, FDD, massive MIMO

## I. INTRODUCTION

The acquisition of downlink channel state information (CSI) is a very challenging task for frequency division duplexing (FDD) massive multiple-input multiple-output (MIMO) systems due to the prohibitively high overheads associated with downlink training and uplink feedback [1]–[3]. By exploiting the angular and delay reciprocities between the uplink and the downlink [4]–[6], conventional methods proposed to reduce the downlink training overhead by extracting frequency-independent information from the uplink CSI, or to reduce the uplink feedback overhead by using compressive sensing based algorithms [7]–[10]. Nevertheless, the conventional methods

either assume that the propagation paths are distinguishable and limited or highly rely on the sparsity of channels.

Recently, artificial intelligence (AI, including machine learning and deep learning, etc.) has been recognized as a potential solution to deal with the high complexity and overheads of wireless communication system [11]–[25]. Great success has been achieved in various applications such as channel estimation [12], [13], data detection [14], [15], CSI feedback [16], beamforming [17], [18], and hybrid precoding [19], [20], etc. Furthermore, AI based downlink CSI prediction for FDD systems has also been studied in many works [22]–[25]. In [22], a convolutional neural network is proposed to predict the downlink CSI from the uplink CSI for single-antenna FDD systems. Based on the CSI correlations between different base station (BS) antennas, linear and support vector based regression methods have been proposed to use the downlink CSI of partial BS antennas to predict the whole downlink CSI, which can reduce the overheads of both downlink pilots and uplink feedback [23]. In fact, [24] develops the channel mapping in space and frequent concept which proves that deep neural networks (DNNs) can learn not only the correlation between closely-located BS antennas but also between the base station arrays that are positioned at different locations or different frequency bands in the same environment. By exploiting the mapping relation between the uplink and downlink CSI in FDD massive MIMO systems, an efficient complex-valued based DNN is trained to predict the downlink CSI only from the uplink CSI, i.e., no downlink pilot is needed at all [25].

Compared with conventional methods (e.g. [7], [8]), AI aided downlink CSI prediction benefits from the excellent learning capability of DNNs and does not require accurate channel models or high computational operations. However, existing AI aided downlink CSI prediction methods [22]–[25] all focus on training models for users in a certain environment. Since users may experience different wireless transmission environments, data collection and training the models from scratch are required for the users in new environments [22]–[25]. Typically thousands of samples and training epochs are required for training one DNN, and thus training a new DNN for each user would face unacceptable time and data cost. Therefore, it is highly desirable to design a method that can adapt to new environments with a small amount of data.

Inspired by human's capability to transfer knowledge from previous experience, transfer learning, which aims to improve the performance of target tasks by exploiting the knowledge from source tasks, becomes a promising technology in machine learning area to solve similar tasks with limited labeled data. Studies on transfer learning date back from 1995 in

Y. Yang and F. Gao are with Institute for Artificial Intelligence Tsinghua University (THUAI), State Key Lab of Intelligent Technologies and Systems, Beijing National Research Center for Information Science and Technology (BNRist), Department of Automation, Tsinghua University, Beijing, 100084, P. R. China (email: yyw18@mails.tsinghua.edu.cn, feifeigao@ieee.org).

Z. Zhong is with the Huawei Technologies Ltd., Shanghai 210206, China (e-mail: zhongzhimeng@huawei.com).

B. Ai is with the State Key Laboratory of Rail Traffic Control and Safety, Beijing Jiaotong University, Beijing 100044, China (email: boai@bjtu.edu.cn).

A. Alkhateeb is with the School of Electrical, Computer and Energy Engineering at Arizona State University, Tempe, AZ 85287, USA (e-mail: alkhateeb@asu.edu).

different names, such as incremental/cumulative learning, life-long learning, multi-task learning, and learning to learn, etc [26]. Generally, transfer learning algorithms acquire knowledge from source tasks by pre-training a model on a large-scale source dataset and then fine-tune the pre-trained model on a small-scale target dataset<sup>1</sup>. With the rapid development of deep learning, the concept of deep transfer learning (DTL) is proposed by combining transfer learning with deep learning [28]. The methodology in transfer learning can be easily generalized to DTL, including instance-transfer, feature-representation-transfer, and parameter-transfer, etc [26]. In particular, the model-agnostic meta-learning algorithm proposed in [29], a typical parameter-transfer learning algorithm, is known for its universal applicability and the state-of-the-art performance. Unlike prior meta-learning algorithms [30]–[32] that learn an update function and impose restrictions on the model architecture, the model-agnostic meta-learning algorithm learns a model initialization that can effectively adapt to a new task with a small amount of labeled data.

In this paper, we formulate the downlink channel prediction for FDD massive MIMO systems as a DTL problem, where each learning task aims to predict the downlink CSI from the uplink CSI for users in a certain environment. Based on the fully-connected neural network (FNN) architecture, we develop the *direct-transfer algorithm*, where the network is trained on the data from all previous environments in the manner of classical deep learning and is then fine-tuned with limited data from a new environment. To achieve more effective transfer, we further proposed the *meta-learning algorithm* that learns a model initialization by an alternating training procedure, consisting of the inner-task and the across-task updates. In the inner-task updates, the meta-learning algorithm learns task-specific parameters by performing gradient descents on task-specific loss functions. In the across-task updates, the meta-learning algorithm update the network parameters by optimizing the loss function associated with multiple tasks. The proposed meta-learning algorithm can adapt to any new environment with a few number of labeled data. We also propose the *no-transfer algorithm* as a baseline algorithm, where the network is trained in the manner of classical deep learning and is directly tested on the data from a new environment without parameter adaption. Simulation results show that the direct-transfer algorithm outperforms the no-transfer algorithm, which validates that transfer learning can effectively improve the performance of downlink channel prediction in new environments. Moreover, the meta-learning algorithm significantly outperforms the direct-transfer algorithm, especially when the number of samples is very small, which demonstrates its superiority over the direct-transfer algorithm.

The remainder of this paper is organized as follows. The system model for FDD massive MIMO systems is introduced in Section II. The DTL problem is formulated in Section III. The no-transfer and direct-transfer algorithms are presented

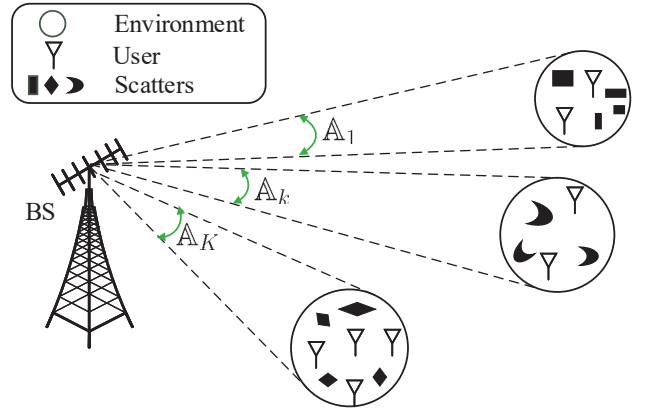


Fig. 1. Downlink CSI prediction for FDD Massive MIMO systems.

in Section IV. The meta-learning algorithm is developed in Section V. Numerical results are given in Section VI. Our main conclusions are given in Section VII.

**Notations:** The bold and lowercase letters denote vectors while the bold and capital letters denote matrices. The notations  $[z]_p$  and  $\text{len}(z)$  denote the  $p$ -th entry and the length of the vector  $\mathbf{x}$ , respectively. The notations  $\Re[\cdot]$  and  $\Im[\cdot]$ , respectively, denote the real and imaginary parts of matrices, vectors or scales. The notation  $\|\mathbf{x}\|_1$  denotes the  $L_1$  norm of  $\mathbf{x}$ . The notation  $(\cdot)^T$  denotes the transpose of a matrix or a vector. The notation  $\mathbb{R}^{2M}$  represents the  $2M$ -dimensional real vector space. The notation  $\circ$  represents the composite mapping operation. The notation  $\mathcal{N}_C(\mathbf{0}, \mathbf{I})$  represents the standard complex Gaussian distribution. The notation  $E[\cdot]$  represents the expectation with respect to all random variables within the brackets. The notation  $\leftarrow$  represents the assignment operation while the notation  $\rightarrow$  represents the direction of the trajectory.

## II. SYSTEM MODEL

We consider an FDD massive MIMO system, where BS is equipped with  $M \gg 1$  antennas in the form of uniform linear array (ULA)<sup>2</sup> and serves multiple single-antenna users, as shown in Fig. 1. The users are randomly distributed in  $K$  regions and users in the same region share the same propagation environment. Denote  $\mathbb{B}_k$  as the set of users in the  $k$ -th region. Then, the channel between the  $u$ -th ( $u \in \mathbb{B}_k$ ) user and BS can be expressed as [8]:

$$\begin{aligned} \mathbf{h}_u(f) &= \int_{\theta \in \mathbb{A}_k} \alpha_u(\theta) \mathbf{a}(\theta) d\theta \\ &= \int_{\bar{\vartheta}_k^-}^{\bar{\vartheta}_k^+} |\alpha_u(\theta)| e^{-j2\pi f \tau_u(\theta) + j\phi_u(\theta)} \mathbf{a}(\theta) d\theta, \end{aligned} \quad (1)$$

where  $f$  is the carrier frequency of the  $u$ -th user, while  $|\alpha_u(\theta)|$ ,  $\phi_u(\theta)$  and  $\tau_u(\theta)$  are the attenuation amplitude, the phase shift and the delay of the incident signal ray coming from direction of arrival (DOA)  $\theta$ , respectively. Meanwhile,

<sup>1</sup>In unsupervised learning [27], labeled data in target tasks are not required. Unsupervised learning algorithms require strong assumptions of data distributions and are hard to be generalized to different problems, and therefore will not be discussed here.

<sup>2</sup>We adopt the ULA model here for simpler illustration, nevertheless, the proposed approach does not restrict to the specific array shape, and therefore is applicable for array with arbitrary geometry.

$\mathbb{A}_k \triangleq [\bar{\vartheta}_k^-, \bar{\vartheta}_k^+]$  is the incident angle spread (AS) of the users in the  $k$ -th region with  $\bar{\vartheta}_k^-$  and  $\bar{\vartheta}_k^+$  being the corresponding lower and upper bounds of AS. The incident AS is assumed to be limited in a narrow region due to the limited local scattering effects at the BS side [4]–[6]. Moreover,  $\mathbf{a}(\theta)$  is the the array manifold vector defined as

$$\mathbf{a}(\theta) = \left[1, e^{-j\varpi \sin \theta}, \dots, e^{-j\varpi(M-1) \sin \theta}\right]^T, \quad (2)$$

where  $\varpi = 2\pi df/c$ ,  $d$  is the antenna spacing, and  $c$  represents the speed of light.

### III. FORMULATION OF DTL PROBLEM

In this section, we first prove the feasibility of applying deep learning to predict the downlink CSI from uplink CSI for a single user. Then, we formulate the downlink channel prediction as a DTL problem and analyze the related transfer learning algorithms.

#### A. Deep Learning for Uplink-to-Downlink Mapping

Since the downlink and the uplink of a certain user share common physical paths and similar spatial propagation characteristics, there exists a deterministic mapping between the downlink and the uplink channels when the position-to-channel mapping is bijective (see more details referring to [24], [25]). Denote  $f_{U,u}$  and  $f_{D,u}$  as the uplink and the downlink frequencies of the  $u$ -th user, respectively. Then, the uplink-to-downlink mapping function can be written as

$$\Psi_{f_{U,u} \rightarrow f_{D,u}} : \{\mathbf{h}_u(f_{U,u})\} \rightarrow \{\mathbf{h}_u(f_{D,u})\}. \quad (3)$$

Although the mapping function  $\Psi_{f_{U,u} \rightarrow f_{D,u}}$  cannot be described by known mathematical tools, it can be approximated by deep neural networks as shown in the following.

Let us consider the simplest three-layers FNN with only one input layer, one hidden layer with  $N$  neurons, and one output layer. Denote  $\mathbf{x}$  and  $\Omega$  as the input data and the network parameter of FNN, respectively. Then the corresponding output can be expressed as  $\text{NET}_N(\mathbf{x}, \Omega)$ . Since deep learning algorithms work in real domain, we introduce the isomorphism between the complex and real domains as

$$\xi : \mathbf{z} \rightarrow (\Re(\mathbf{z}^T), \Im(\mathbf{z}^T))^T. \quad (4)$$

Denote the inverse mapping of  $\xi$  as  $\xi^{-1}$  that can be given as

$$\xi^{-1} : (\Re(\mathbf{z}^T), \Im(\mathbf{z}^T))^T \rightarrow \mathbf{z}. \quad (5)$$

Based on the universal approximation theorem [33], we obtain the following theorem:

*Theorem 1:* For any given error  $\varepsilon > 0$ , there exists a positive constant  $N$  large enough such that

$$\sup_{\mathbf{x} \in \mathbb{H}} \|\text{NET}_N(\mathbf{x}, \Omega) - \Psi_{\xi, u, U \rightarrow D}(\mathbf{x})\| \leq \varepsilon, \quad \mathbb{H} = \{\xi \circ \mathbf{h}_u(f_{U,u})\} \subseteq \mathbb{R}^{2M}, \quad (6)$$

with  $\Psi_{\xi, u, U \rightarrow D}(\mathbf{x}) = \xi \circ \Psi_{f_{U,u} \rightarrow f_{D,u}} \circ \xi^{-1}(\mathbf{x})$ .

*Proof.* (i) Since  $\forall \mathbf{x} \in \mathbb{H}$ ,  $|\mathbf{x}| = |\mathbf{h}_u(f_{U,u})| \leq \sqrt{M}\vartheta_u |\alpha_u(\theta)|$ , we know  $\mathbf{x}$  is bounded. Therefore,  $\mathbb{H}$  is a compact set;

(ii) Since  $\Psi_{f_{U,u} \rightarrow f_{D,u}}$  and  $\xi$  are continuous mapping and composition of continuous mappings is still a continuous mapping, we know  $\Psi_{\xi, u, U \rightarrow D}(\mathbf{x})$  is a continuous function for  $\forall \mathbf{x} \in \mathbb{H}$ . Therefore, the  $i$ -th element of  $\Psi_{\xi, u, U \rightarrow D}(\mathbf{x})$ , denoted by  $[\Psi_{\xi, u, U \rightarrow D}(\mathbf{x})]_i$ , is also a continuous function of  $\mathbf{x}$ . Based on (i), (ii) and universal approximation theorem [33, Theorem 1], we know for any  $\varepsilon_i > 0$ , there exists a positive constant  $N_i$  such that

$$\sup_{\mathbf{x} \in \mathbb{H}} |\text{NET}_{N_i}(\mathbf{x}, \Omega_i) - [\Psi_{\xi, u, U \rightarrow D}(\mathbf{x})]_i| \leq \varepsilon_i, \quad (7)$$

where  $\text{NET}_{N_i}(\mathbf{x}, \Omega_i)$  is the output of the network with only one neuron in the output layer, and  $\Omega_i$  denotes the corresponding network parameters. By stacking  $2M$  of the above networks together, we can construct a larger network  $\text{NET}_N(\mathbf{x}, \Omega)$  with  $N = \sum_{i=1}^{2M} N_i$ , where  $\Omega$  denotes the corresponding network parameters. By choosing  $\varepsilon_i = \varepsilon/\sqrt{2M}$ , Eq. (6) can be proved.  $\square$

*Theorem 1* reveals that the uplink-to-downlink mapping function can be approximated arbitrarily well by an FNN with a single hidden layer. It should be mentioned that although the network proposed in this work cannot obtain arbitrarily high prediction accuracy due to inadequate learning or insufficient number of hidden units [33], *Theorem 1* still provides theoretical foundation for the application of deep neural networks.

#### B. Definitions of DTL

For the users in the  $k$ -th environment (region), let  $\mathcal{X}_k$  and  $\mathcal{Y}_k$  denote the space<sup>3</sup> of the uplink and the downlink channels, respectively. Then, the definitions of the “domain” and the “task” are given in the following two definitions:

**Definition 1:** The “domain”  $\mathcal{D}(k)$  is composed of the feature space  $\mathcal{X}_k$  and the marginal probability distribution  $P(\mathbf{h}_u(f_{U,u}))|_{u \in \mathbb{B}_k}$ , i.e.,  $\mathcal{D}(k) = \{\mathcal{X}_k, P(\mathbf{h}_u(f_{U,u}))|_{u \in \mathbb{B}_k}\}$ .

**Definition 2:** Define the “task”  $\mathcal{T}(k)$  as the prediction of the downlink channels from the uplink channels for the users in the  $k$ -th environment. Given the specific domain  $\mathcal{D}(k)$ , the “task”  $\mathcal{T}(k)$  is composed of the label space  $\mathcal{Y}_k$  and the prediction function  $\mathcal{F}_k$ , i.e.,  $\mathcal{T}(k) = \{\mathcal{Y}_k, \mathcal{F}_k\}$ .

*Remark 1:* The prediction function  $\mathcal{F}_k$  can be learned from the training data of the  $k$ -th environment and then be used to predict the downlink channels for the users in the  $k$ -th environment. Based on the analysis in Section III-A, the prediction function  $\mathcal{F}_k$  can be interpreted as the fitting function for all the uplink-to-downlink mapping functions defined in the  $k$ -th environment, i.e.,  $\{\Psi_{f_{U,u} \rightarrow f_{D,u}}\}_{u \in \mathbb{B}_k}$ . The prediction function  $\mathcal{F}_k$  can also be interpreted as the conditional probability distribution  $P(\mathbf{h}_u(f_{D,u})|\mathbf{h}_u(f_{U,u}))|_{u \in \mathbb{B}_k}$  from probabilistic view.

Classical transfer learning consists of two aspects, namely, the source domain transfer and the target domain adaption. Based on [26], the definition of transfer learning can be given as following:

<sup>3</sup> The space of the uplink channels is  $\mathcal{X}$  means that any possible uplink channel vector belongs to  $\mathcal{X}$ , i.e.  $\mathbf{h}_u(f_{U,u}) \in \mathcal{X}$  holds for any possible  $u$  or  $f_{U,u}$ .

**Definition 3:** Given the source task  $\mathcal{T}_S$ , the source domain  $\mathcal{D}_S$ , the target task  $\mathcal{T}_T$ , and the target domain  $\mathcal{D}_T$ , the aim of transfer learning is to improve the performance of the target task  $\mathcal{T}_T$  by using the knowledge from  $\mathcal{T}_S$  and  $\mathcal{D}_T$ , where  $\mathcal{D}_T \neq \mathcal{D}_S$  or  $\mathcal{T}_T \neq \mathcal{T}_S$ .

Here we extend the single-source domain transfer to the multi-source domain transfer. Then, a more generalized definition of transfer learning can be provided as follows:

**Definition 4:** Given the number of source tasks  $K_s$ , the source tasks  $\{\mathcal{T}_S(k)\}_{k=1}^{K_s}$ , the source domains  $\{\mathcal{D}_S(k)\}_{k=1}^{K_s}$ , the target task  $\mathcal{T}_T$ , and the target domain  $\mathcal{D}_T$ , the aim of transfer learning is to improve the performance of the target task  $\mathcal{T}_T$  by using the knowledge from  $\{\mathcal{T}_S(k)\}_{k=1}^{K_s}$  and  $\{\mathcal{D}_S(k)\}_{k=1}^{K_s}$ , where  $\mathcal{D}_T \neq \mathcal{D}_S(k)$  or  $\mathcal{T}_T \neq \mathcal{T}_S(k)$  holds for  $k = 1, \dots, K$ .

In Definition 4, the condition  $\mathcal{D}_T \neq \mathcal{D}_S(k)$  means that either the corresponding feature space  $\mathcal{X}_T \neq \mathcal{X}_S(k)$  holds or the corresponding marginal probability distribution  $P_T(\mathbf{h}_u(f_{U,u}))|_{u \in \mathbb{B}_T} \neq P_{S(k)}(\mathbf{h}_u(f_{U,u}))|_{u \in \mathbb{B}_k}$  holds, where  $\mathbb{B}_T$  is the set of users in the target environment. The condition  $\mathcal{T}_T \neq \mathcal{T}_S(k)$  means that either the corresponding label space  $\mathcal{Y}_T \neq \mathcal{Y}_S(k)$  holds or the corresponding conditional probability distribution  $P_T(\mathbf{h}_u(f_{D,u})|\mathbf{h}_u(f_{U,u}))|_{u \in \mathbb{B}_T} \neq P_{S(k)}(\mathbf{h}_u(f_{D,u})|\mathbf{h}_u(f_{U,u}))|_{u \in \mathbb{B}_k}$  holds. Since the conditional probability distributions for different prediction tasks are different, the condition  $\mathcal{T}_T \neq \mathcal{T}_S(k)$  is satisfied. Therefore, the downlink channel prediction for massive MIMO systems can be formulated as a transfer learning problem.

DTL is to transfer knowledge by deep neural networks, which is defined as follows [28]:

**Definition 5:** Given a transfer learning task characterized by  $\langle \{\mathcal{T}_S(k)\}_{k=1}^{K_s}, \{\mathcal{D}_S(k)\}_{k=1}^{K_s}, \mathcal{T}_T, \mathcal{D}_T \rangle$ , it is a DTL task when the prediction function  $\mathcal{F}_T$  of  $\mathcal{T}_T$  is a non-linear function that is approximated by a deep neural network.

Based on Definition 4 and Definition 5, the downlink channel prediction for massive MIMO systems can be formulated as a typical DTL problem, where the  $k$ -th learning task is to predict the downlink channel  $\mathbf{h}_u(f_{D,u})$  from the uplink channel  $\mathbf{h}_u(f_{U,u})$  for the users in the  $k$ -th environment.

### C. Motivation of Meta-learning

Before resorting to DTL methods, we should first consider the question “whether to transfer”. In fact, if the source and target tasks are highly related, a deep neural network would perform well, without the need for fine-tuning the neural network based on the target environment, thanks to its remarkable generalization capability.

The second question is “how to transfer”. One natural solution is to directly use all the labeled data in the source tasks to train the network, and then fine-tune the network with the labeled data in the target task. However, the direct transfer method tends to overfit when the number of labeled data in the target task is small, as will be verified by simulations in Section VI. To overcome this challenge, and motivated by the state-of-the-art performance of the model-agnostic meta-learning algorithm in few-sample learning [29], we will propose the meta-learning algorithm based downlink channel prediction in Section V.

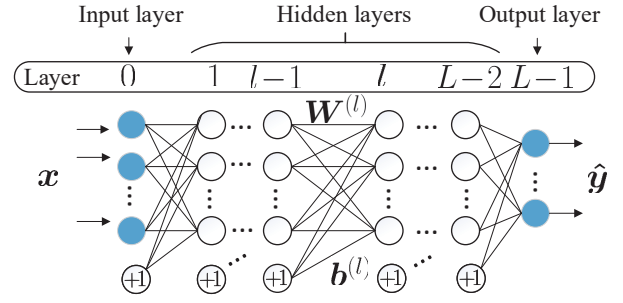


Fig. 2. The FNN architecture.

## IV. NO-TRANSFER AND DIRECT-TRANSFER ALGORITHMS

Based on the analysis in Section III-C, we propose the no-transfer algorithm based on classical deep learning algorithms to investigate the necessity of transfer learning. Then, the direct-transfer algorithm is proposed and used as the benchmark of the meta-learning algorithm. Both the no-transfer and the direct-transfer algorithms adopt the FNN architecture as described in the following subsection.

### A. Network Architecture

The FNN architecture has  $L$  layers, including one input layer,  $L-2$  hidden layers, and one output layer, as shown in Fig. 2. The input of FNN is  $\mathbf{x} = \boldsymbol{\xi} \circ \mathbf{h}_u(f_{U,u})$ . The output of the network is a cascade of nonlinear transformation of  $\mathbf{x}$ , i.e.,

$$\hat{\mathbf{y}} = \text{Net}(\mathbf{x}, \boldsymbol{\Omega}) = \mathbf{f}^{(L-1)} \circ \mathbf{f}^{(L-2)} \circ \dots \circ \mathbf{f}^{(1)}(\mathbf{x}), \quad (8)$$

where  $\boldsymbol{\Omega} \triangleq \{\mathbf{W}^{(l)}, \mathbf{b}^{(l)}\}_{l=1}^{L-1}$  is the network parameters to be trained. Moreover,  $\mathbf{f}^{(l)}$  is the nonlinear transformation function of the  $l$ -th layer and can be written as,

$$\mathbf{f}^{(l)}(\mathbf{x}) = \mathbf{r}^{(l)}(\mathbf{W}^{(l)}\mathbf{x} + \mathbf{b}^{(l)}), \quad 1 \leq l \leq L-1, \quad (9)$$

where  $\mathbf{W}^{(l)}$  is the weight matrix associated with the  $(l-1)$ -th and the  $l$ -th layers, while  $\mathbf{b}^{(l)}$  and  $\mathbf{r}^{(l)}$  are the bias vector and the activation function of the  $l$ -th layer, respectively. The activation function for the hidden layers is selected as the rectified linear unit (ReLU) function  $[\mathbf{r}_{\text{re}}(\mathbf{z})]_p = \max\{0, [\mathbf{z}]_p\}$  with  $p = 1, 2, \dots, \text{len}(\mathbf{z})$ . The activation function for the output layer is the linear function, i.e.,  $\mathbf{r}_{\text{li}}(\mathbf{z}) = \mathbf{z}$ . The loss function can be written as follow:

$$\text{Loss}_{\mathbb{D}}(\boldsymbol{\Omega}) = \frac{1}{V} \sum_{v=0}^{V-1} \left\| \hat{\mathbf{y}}^{(v)} - \mathbf{y}^{(v)} \right\|_2^2, \quad (10)$$

where  $\mathbf{y} = \boldsymbol{\xi} \circ \mathbf{h}_u(f_{D,u})$  is the supervise label,  $V$  is the batch size<sup>4</sup>, the superscript  $(v)$  denotes the index of the  $v$ -th training sample,  $\mathbb{D} = \{(\mathbf{x}^{(v)}, \mathbf{y}^{(v)})\}_{v=1}^V$  is the training dataset, and  $\|\cdot\|_2$  denotes the  $L_2$  norm.

<sup>4</sup>Batch size is the number of samples in one training batch.

### B. Definitions and Generation of Datasets

As mentioned in *Definition 2*, each task represents the prediction of downlink channel from the uplink channel for the users in a certain environment. Denote  $\Delta f$  as the frequency difference between the uplink and the downlink. Then, each uplink frequency corresponds to a certain downlink frequency, i.e.,  $f_{D,u} = f_{U,u} + \Delta f$ . Denote  $U$  as the number of users in the  $k$ -th environment that are involved in the generation of datasets. Denote the training dataset of the  $k$ -th source task as  $\mathbb{D}_{Tr}(k)$ . To generate  $\mathbb{D}_{Tr}(k)$ ,  $N_{Tr}$  sample pairs  $\{(\xi \circ \mathbf{h}_u(f_{D,u}), \xi \circ \mathbf{h}_u(f_{U,u}))\}$  are collected by randomly selecting  $N_{Tr}/U$  uplink frequencies for each user in the  $k$ -th environment.

Denote the adaption and testing datasets of the  $k$ -th target task as  $\mathbb{D}_{Ad}(k)$  and  $\mathbb{D}_{Te}(k)$ , respectively. Similarly, the datasets  $\mathbb{D}_{Ad}(k)$  and  $\mathbb{D}_{Te}(k)$  can be obtained by separately collecting  $N_{Ad}$  and  $N_{Te}$  sample pairs for the  $k$ -th target task. Note that  $\mathbb{D}_{Ad}(k) \cap \mathbb{D}_{Ad}(k) = \emptyset$  should be satisfied to ensure the testing dataset not known to the networks.

### C. No-Transfer Algorithm

The no-transfer algorithm regards the DTL problem as one classical deep learning problem, including the training and the testing stages. The datasets  $\{\mathbb{D}_{Tr}(k)\}_{k=1}^{K_S}$  for  $K_S$  source tasks are generated and are used to train the network. To avoid overfitting, we randomly shuffle all the samples in datasets  $\{\mathbb{D}_{Tr}(k)\}_{k=1}^{K_S}$  and obtain the shuffled training dataset  $\mathbb{D}_{STr}$ .

In each time step,  $V$  training samples are randomly selected as the training batch  $\mathbb{D}_{STrB}$ . Then, we employ the adaptive moment estimation (ADAM) algorithm [34] to minimize the loss function on  $\mathbb{D}_{STrB}$ , i.e.,  $\text{Loss}_{\mathbb{D}_{STrB}}(\Omega)$ . Instead of using the gradient of the current step to guide the updates, ADAM adopts moments of the gradient to determine the direction of the optimization, which can accelerate the training, especially in the case of high curvature or noisy gradients [35, chapter 8]. The  $q$ -th moment of a random variable is defined as the expected value of the  $q$ -th power of the variable. To estimate the 1<sup>st</sup> and the 2<sup>nd</sup> moment vectors of the gradient, i.e.,  $\nabla_{\Omega} \text{Loss}_{\mathbb{D}_{STrB}}(\Omega)$ , ADAM computes the exponentially decaying moving averages of past gradients and elementwise squared gradients, i.e.,

$$\begin{cases} \mu_1 \leftarrow \rho_1 \mu_1 + (1 - \rho_1) \nabla_{\Omega} \text{Loss}_{\mathbb{D}_{STrB}}(\Omega), \\ \nu_1 \leftarrow \rho_2 \nu_1 + (1 - \rho_2) (\nabla_{\Omega} \text{Loss}_{\mathbb{D}_{STrB}}(\Omega))^2, \end{cases} \quad (11)$$

where  $(\nabla_{\Omega} \text{Loss}_{\mathbb{D}_{STrB}}(\Omega))^2$  represents the elementwise square of the gradient, and the hyper-parameters  $\rho_1, \rho_2 \in [0, 1]$  are the exponential decay rates of the moving averages. The vectors  $\mu_1$  and  $\nu_1$  are initialized with zeros and updated with moving averages, which results in biased moment estimates of the gradient. Following [34], unbiased 1<sup>st</sup> and 2<sup>nd</sup> moment vectors can be obtained by

$$\begin{cases} \bar{\mu}_1 \leftarrow \mu_1 / (1 - \rho_1^t), \\ \bar{\nu}_1 \leftarrow \nu_1 / (1 - \rho_2^t), \end{cases} \quad (12)$$

where  $t$  represents the time step. Then, the unbiased 1<sup>st</sup> and 2<sup>nd</sup> moment vectors, i.e.,  $\bar{\mu}_1$  and  $\bar{\nu}_1$ , are used to scale the

learning rate individually for all elements in  $\Omega$ , which can be expressed as

$$\Omega \leftarrow \Omega - \eta \bar{\mu}_1 / (\sqrt{\bar{\nu}_1} + \varepsilon), \quad (13)$$

where  $\eta$  is the learning rate,  $\varepsilon$  is the disturbance factor, and the notation “/” is the elementwise division operation.

By iteratively executing the ADAM algorithm on the dataset  $\mathbb{D}_{STr}$  until  $\text{Loss}_{\mathbb{D}_{STrB}}(\Omega)$  converges, we can obtain the trained network parameter  $\Omega_{Nt}$ . In the testing stage, the network parameter  $\Omega_{Nt}$  is fixed. The testing datasets  $\{\mathbb{D}_{Te}(k)\}_{k=1}^{K_T}$  of  $K_T$  target tasks are generated and are used to test the performance of the no-transfer algorithm. Denote  $\hat{\mathbf{h}}_D = \xi^{-1} \circ \hat{\mathbf{y}}$  and  $\mathbf{h}_D = \xi^{-1} \circ \mathbf{y}$  as the estimated and the true downlink channel, respectively. Normalized mean-squared-error (NMSE) is used to measure the prediction accuracy, which is defined as

$$\text{NMSE} = E \left[ \left\| \mathbf{h}_D - \hat{\mathbf{h}}_D \right\|_2^2 / \left\| \mathbf{h}_D \right\|_2^2 \right]. \quad (14)$$

Denote the prediction NMSE of the no-transfer algorithm as  $\text{NMSE}_{Nt}$  that can be calculated by averaging the NMSEs of the no-transfer algorithm evaluated on  $K_T$  target environments. i.e.,  $\text{NMSE}_{Nt} = \sum_{k=1}^{K_T} \text{NMSE}_{Nt}(k) / K_T$ , where  $\text{NMSE}_{Nt}(k)$  is the NMSE evaluated on the  $k$ -th target environment. The concrete steps of the no-transfer algorithm are given in the **Algorithm 1**.

### D. Direct Transfer Algorithm

In the training stage, the direct-transfer algorithm trains the network via the ADAM algorithm on the training dataset  $\mathbb{D}_{STr}$  until  $\text{Loss}_{\mathbb{D}_{STrB}}(\Omega)$  converges. During the training, the network tries to learn a generalized parameter to predict the downlink CSI for users in different environments. After the training finished, the network has learned the parameter  $\Omega_{Nt}$  from the  $K_S$  source tasks. Then, the datasets  $\{\mathbb{D}_{Ad}(k)\}_{k=1}^{K_T}$  and  $\{\mathbb{D}_{Te}(k)\}_{k=1}^{K_T}$  of  $K_T$  target tasks are generated following Section IV-B. For the  $k$ -th target task, the direct-transfer algorithm initializes the target-task-specific parameter  $\Omega_{T,k}$  as the trained network parameter  $\Omega_{Nt}$ , then fine-tunes  $\Omega_{T,k}$  on the adaption dataset  $\mathbb{D}_{Ad}(k)$  via  $G_{Ad}$  steps of ADAM updates, and then tests the network on the testing dataset  $\mathbb{D}_{Te}(k)$ . Specifically, ADAM aims to minimize  $\text{Loss}_{\mathbb{D}_{Ad}(k)}(\Omega_{T,k})$  by updating the parameter  $\Omega_{T,k}$ , and unbiased moment vectors of the gradient  $\nabla_{\Omega_{T,k}} \text{Loss}_{\mathbb{D}_{Ad}(k)}(\Omega_{T,k})$  are used to determine the direction of the update. Based on the derivations in Section IV-C, biased 1<sup>st</sup> and 2<sup>nd</sup> moment vectors are initialized with zeros and are then updated in each step following

$$\begin{cases} \mu_2 \leftarrow \rho_1 \mu_2 + (1 - \rho_1) \nabla_{\Omega_{T,k}} \text{Loss}_{\mathbb{D}_{Ad}(k)}(\Omega_{T,k}), \\ \nu_2 \leftarrow \rho_2 \nu_2 + (1 - \rho_2) (\nabla_{\Omega_{T,k}} \text{Loss}_{\mathbb{D}_{Ad}(k)}(\Omega_{T,k}))^2. \end{cases} \quad (15)$$

Next, we obtain the unbiased 1<sup>st</sup> and 2<sup>nd</sup> moment vectors by

$$\begin{cases} \bar{\mu}_2 \leftarrow \mu_2 / (1 - \rho_1^g), \\ \bar{\nu}_2 \leftarrow \nu_2 / (1 - \rho_2^g), \end{cases} \quad (16)$$

where  $g$  represents the step index. The update of  $\Omega_{T,k}$  can be written as

$$\Omega_{T,k} \leftarrow \Omega_{T,k} - \eta \bar{\mu}_2 / (\sqrt{\bar{\nu}_2} + \varepsilon). \quad (17)$$

**Algorithm 1:** No-transfer algorithm for downlink CSI prediction

**Input:** Source tasks:  $\{\mathcal{T}_S(k)\}_{k=1}^{K_S}$ , Target tasks:  $\{\mathcal{T}_T(k)\}_{k=1}^{K_T}$ , learning rate:  $\eta$ , exponential decay rates:  $\{\rho_1, \rho_2\}$ , disturbance factor:  $\varepsilon$ , batch size:  $V$

**Output:** Trained network parameter:  $\Omega_{Nt} \leftarrow \Omega$ , predicted downlink CSI based on  $\{\mathbb{D}_{Te}(k)\}_{k=1}^{K_T}$ , NMSE of the no-transfer algorithm:  $\text{NMSE}_{Nt}$

```

1 Training stage
2 Randomly initialize the network parameters  $\Omega$ 
3 Initialize the time step:  $t \leftarrow 0$ 
4 Initialize the 1st and 2nd moment vectors:  $\mu_1, \nu_1 \leftarrow \mathbf{0}$ 
5 Generate the training dataset  $\mathbb{D}_{Tr}(k)$  for each source task
6 Randomly shuffle training datasets  $\{\mathbb{D}_{Tr}(k)\}_{k=1}^{K_S}$  and
   obtain the shuffled training datasets  $\mathbb{D}_{STr}$ 
7 while not done do
8   Randomly select  $V$  training samples from  $\mathbb{D}_{STr}$  as the
   training batch  $\mathbb{D}_{STrB}$ 
9    $t \leftarrow t + 1$ 
10  Update biased 1st and 2nd moment vectors using
   Eq. (11)
11  Compute unbiased 1st and 2nd moment vectors using
   Eq. (12)
12  Update the parameter  $\Omega$  using unbiased 1st and 2nd
   moment vectors:  $\Omega \leftarrow \Omega - \eta \bar{\mu}_1 / (\sqrt{\bar{\nu}_1} + \varepsilon)$ 
13 end
14 Testing stage
15 Initialize NMSE:  $\text{NMSE}_{Nt} \leftarrow 0$ 
16 for  $k = 1, \dots, K_T$  do
17   Generate the testing dataset  $\mathbb{D}_{Te}(k)$  for  $\mathcal{T}_T(k)$ 
18   Predict the downlink CSI base on  $\mathbb{D}_{Te}(k)$  and  $\Omega$ 
   using Eq. (8)
19   Calculate  $\text{NMSE}_{Nt}(k)$  using Eq. (14)
20    $\text{NMSE}_{Nt} \leftarrow \text{NMSE}_{Nt} + \text{NMSE}_{Nt}(k) / K_T$ 
21 end

```

After the fine-tuning finished, the target-task-specific parameter  $\Omega_{T,k}$  will be fixed. Let  $\text{NMSE}_{Dt}(k)$  represent the NMSE of the direct-transfer algorithm evaluated on the  $k$ -th target environment that can be obtained by testing the network on the dataset  $\mathbb{D}_{Te}(k)$  using Eqs. (8) and (14). Then, the prediction NMSE of the direct-transfer algorithm can be obtained by  $\text{NMSE}_{Dt} = \sum_{k=1}^{K_T} \text{NMSE}_{Dt}(k) / K_T$ . The concrete steps of the direct-transfer algorithm are given in the **Algorithm 2**.

Notice that the direct-transfer algorithm utilizes the trained network parameter  $\Omega_{Nt}$  as the initialization of the adaption stage for every target task. By contrast, another way is to utilize previous target-task-specific parameter  $\Omega_{T,k-1}$  as the initialization of the adaption stage for the  $k$ -th target task since  $\Omega_{T,k-1}$  appears to contain the information in all the source environments and  $k-1$  target environments, and therefore could enhance the performance of the algorithm. However, the real situation is that a good initialization in DTL is not to train the network parameter with as much data as possible, but to find such a parameter that can easily adapt to any target

**Algorithm 2:** Direct-transfer algorithm for downlink CSI prediction

**Input:** Source tasks:  $\{\mathcal{T}_S(k)\}_{k=1}^{K_S}$ , Target tasks:  $\{\mathcal{T}_T(k)\}_{k=1}^{K_T}$ , learning rate:  $\eta$ , exponential decay rates:  $\{\rho_1, \rho_2\}$ , disturbance factor:  $\varepsilon$ , batch size:  $V$ , number of gradsteps for adaption:  $G_{Ad}$

**Output:** Trained network parameter:  $\Omega_{Nt}$ , predicted downlink CSI based on  $\{\mathbb{D}_{Te}(k)\}_{k=1}^{K_T}$ , NMSE of the direct-transfer algorithm:  $\text{NMSE}_{Dt}$

```

1 Training stage
2 Implement the Algorithm 1 and obtain the trained
   network parameter  $\Omega_{Nt}$ 
3 Adaption and Testing for target tasks
4 Initialize NMSE:  $\text{NMSE}_{Dt} \leftarrow 0$ 
5 for  $k = 1, \dots, K_T$  do
6   Generate the datasets  $\mathbb{D}_{Ad}(k)$  and  $\mathbb{D}_{Te}(k)$  for  $\mathcal{T}_T(k)$ 
7   Adaption stage
8   Load the network parameter  $\Omega_{T,k} \leftarrow \Omega_{Nt}$ 
9   Initialize the 1st and 2nd moment vectors:  $\mu_2, \nu_2 \leftarrow \mathbf{0}$ 
10  for  $g = 1, \dots, G_{Ad}$  do
11    Update biased 1st and 2nd moment vectors using
    Eq. (15)
12    Compute unbiased 1st and 2nd moment vectors
    using Eq. (16)
13    Update the parameter  $\Omega_{T,k}$  using unbiased
    moment vectors:
     $\Omega_{T,k} \leftarrow \Omega_{T,k} - \eta \bar{\mu}_2 / (\sqrt{\bar{\nu}_2} + \varepsilon)$ 
14  end
15  Testing stage
16  Predict the downlink CSI base on  $\mathbb{D}_{Te}(k)$  and  $\Omega_{T,k}$ 
   using Eq. (8)
17  Calculate  $\text{NMSE}_{Dt}(k)$  using Eq. (14)
18   $\text{NMSE}_{Dt} \leftarrow \text{NMSE}_{Dt} + \text{NMSE}_{Dt}(k) / K_T$ 
19 end

```

task. In a way, a good initial parameter can be interpreted as a parameter point in the parameter space that is near (or easy to get) to the optimal parameter point for most target tasks. Fig. 3 illustrates the trajectory of parameter updating during the training and adaption stages of the direct-transfer algorithm. The black solid line represents the updating trajectory of the network parameter  $\Omega$  during the training stage. The blue dashed lines represent the updating trajectories of the target-task-specific parameters  $\{\Omega_{T,k}\}_{k=1}^3$  during the adaption stage following **Algorithm 2**. The red dash-dotted lines represent the parameter adaption process with previous target-task-specific parameter as the initial parameter. As shown in Fig. 3, the trajectory of red dash-dotted lines, i.e.,  $\Omega_{Nt} \rightarrow \Omega_{T,1} \rightarrow \Omega_{T,2} \rightarrow \Omega_{T,3}$ , is longer than the trajectory of blue dashed lines, i.e.,  $\{\Omega_{Nt} \rightarrow \Omega_{T,k}\}_{k=1}^3$ , which indicates  $\Omega_{Nt}$  is a better initialization than  $\Omega_{T,k-1}$ . The intrinsic reason is that  $\Omega_{Nt}$  is obtained by training the network on a large-scale and shuffled dataset, i.e.,  $\mathbb{D}_{STr}$ , which has more stronger representativeness of different tasks than the small-scale adaption dataset  $\mathbb{D}_{Ad}(k)$ . In fact, the fine-tuning on the adaption dataset  $\mathbb{D}_{Ad}(k)$  renders the network parameter away from the generalized parameter



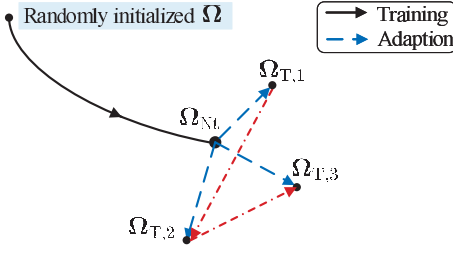


Fig. 3. Illustration of the direct-transfer algorithm, where  $K_T = 3$ .

$\Omega_{Nt}$  and gets closed to the task-specific parameter  $\Omega_{T,k}$ .

*Remark 2:* It should be emphasized that parameter initialization is one of the most crucial tasks in deep learning, especially in DTL. A high-quality initialization can effectively prevent the training from getting stuck in a low-performance local minimum and thus accelerates the convergence [36], [37].

## V. META-LEARNING ALGORITHM

The proposed meta-learning algorithm also adopts the FNN architecture in Section IV-A and has three stages, namely, the meta-training, the meta-adaption and the testing stages.

### A. Definitions and Generation of Datasets

For the  $k$ -th source task, two datasets should be obtained for the meta-training stage, i.e., the training support dataset  $\mathbb{D}_{TrSup}(k)$  and the training query dataset  $\mathbb{D}_{TrQue}(k)$ . Following Section IV-B, we can collect  $N_{Tr}$  sample pairs for each source task and then randomly divide the  $N_{Tr}$  sample pairs, i.e.,  $\{(\xi \circ h_u(f_{D,u}), \xi \circ h_u(f_{U,u}))\}$ , into the training support dataset  $\mathbb{D}_{TrSup}(k)$  and the training query dataset  $\mathbb{D}_{TrQue}(k)$ . Note that  $\mathbb{D}_{TrSup}(k) \cap \mathbb{D}_{TrQue}(k) = \emptyset$  should be satisfied to improve the generalization capability of the network. For the  $k$ -th target task, the adaption dataset  $\mathbb{D}_{Ad}(k)$  and the testing dataset  $\mathbb{D}_{Te}(k)$  can be obtained following Section IV-B.

### B. Meta-training Stage

In the meta-training stage, the meta-learning algorithm aims to learn a network initialization that can effectively adapt to a new task. The parameter of FNN, i.e.,  $\Omega$ , is randomly initialized and is then updated by two iterative processes, i.e., the inner-task and the across-task updates.

1) *Inner-task update:* Denote the source-task-specific parameter as  $\Omega_{S,k}$ . The goal of the inner-task updates is to minimize the loss function on the training support dataset  $\mathbb{D}_{TrSup}(k)$ , i.e.,  $\text{Loss}_{\mathbb{D}_{TrSup}(k)}(\Omega_{S,k})$ , by iteratively updating the parameter  $\Omega_{S,k}$ . Specifically,  $\Omega_{S,k}$  is initialized as the current network parameter  $\Omega$ , and is then updated with  $G_{Tr}$  steps of gradient descents. Since an overlarge gradient leads to the instability of gradient descents, we limit the values of the source-task-specific gradient  $\nabla_{\Omega_{S,k}} \text{Loss}_{\mathbb{D}_{TrSup}(k)}(\Omega_{S,k})$  into a certain range and obtain the truncated source-task-specific gradient  $\mathbf{v}_{S,k}$  as [38]

$$[\mathbf{v}_{S,k}]_p = \min \left\{ \Upsilon, \left[ \nabla_{\Omega_{S,k}} \text{Loss}_{\mathbb{D}_{TrSup}(k)}(\Omega_{S,k}) \right]_p \right\}, \\ p = 1, \dots, \text{len}(\mathbf{v}_{S,k}), \quad (18)$$

where  $\Upsilon$  is the upper threshold of the gradient. Generally, an overlarge  $\Upsilon$  may lead to large fluctuations of the loss function, while a too small  $\Upsilon$  distorts the direction of the update, resulting in too early convergence. By trials and adjustments, a proper  $\Upsilon$  is selected such that the inner-task updates converge stably and efficiently. Notice that the gradient truncation is not required in ADAM since ADAM normalizes the size of update steps for each parameter by using the square root of the 2<sup>nd</sup> moment vector to divide the 1<sup>st</sup> moment vector, as shown in Eq. (17). Then, the parameter  $\Omega_{S,k}$  is updated based on the truncated gradient  $\mathbf{v}_{S,k}$ , i.e.,

$$\Omega_{S,k} \leftarrow \Omega_{S,k} - \beta \mathbf{v}_{S,k}, \quad (19)$$

where  $\beta$  is the inner-task learning rate. It should be mentioned that in [29], the inner-task update only involves one step of gradient descent. Here we extend it to  $G_{Tr}$  steps for better performance.

2) *Across-task update:* In the  $t$ -th time step, we randomly select  $K_B$  source tasks out of  $K_S$  source tasks<sup>5</sup> and generate corresponding support datasets  $\{\mathbb{D}_{TrSup}(k)\}_{k=1}^{K_B}$  and query datasets  $\{\mathbb{D}_{TrQue}(k)\}_{k=1}^{K_B}$  following Section V-A. The optimization objective of the  $t$ -th update is the loss function associated with the  $K_B$  training query datasets  $\{\mathbb{D}_{TrQue}(k)\}_{k=1}^{K_B}$  and the meta-trained source-task-specific parameter  $\Omega_{S,k}$ , i.e.,

$$\sum_{k=1}^{K_B} \text{Loss}_{\mathbb{D}_{TrQue}(k)}(\Omega_{S,k}).$$

The network parameter  $\Omega$  is updated via the ADAM algorithm. Specifically, the biased 1<sup>st</sup> and 2<sup>nd</sup> moment vectors are initialized with zeros and are then updated in the  $t$ -th time step following

$$\begin{cases} \mu_3 \leftarrow \rho_1 \mu_3 + (1 - \rho_1) \nabla_{\Omega} \sum_{k=1}^{K_B} \text{Loss}_{\mathbb{D}_{TrQue}(k)}(\Omega_{S,k}), \\ \nu_3 \leftarrow \rho_2 \nu_3 + (1 - \rho_2) \left( \nabla_{\Omega} \sum_{k=1}^{K_B} \text{Loss}_{\mathbb{D}_{TrQue}(k)}(\Omega_{S,k}) \right)^2. \end{cases} \quad (20)$$

Next, the unbiased 1<sup>st</sup> and 2<sup>nd</sup> moment vectors can be obtained by

$$\begin{cases} \bar{\mu}_3 \leftarrow \mu_3 / (1 - \rho_1^t), \\ \bar{\nu}_3 \leftarrow \nu_3 / (1 - \rho_2^t). \end{cases} \quad (21)$$

The update of network parameter  $\Omega$  can be written as follows:

$$\Omega \leftarrow \Omega - \eta \bar{\mu}_3 / (\sqrt{\bar{\nu}_3} + \varepsilon). \quad (22)$$

After the across-task updates finished, the meta-trained network parameter will be obtained by  $\Omega_{Mt} \leftarrow \Omega$ .

There are two important differences between the inner-task and across-task updates. (i) Each inner-task update is performed on the support dataset of one source task while each across-task update is performed on the query datasets of  $K_B$  source tasks; (ii) The aim of inner-task updates is to optimize the task-specific parameters  $\{\Omega_{T,k}\}_{k=1}^{K_B}$  while the aim of across-task updates is to optimize the overall network parameters  $\Omega$ .

<sup>5</sup>Typically, there is  $K_B \ll K_S$ .

Besides, the across-task learning rate  $\gamma$  is typically set to be smaller than the inner-task learning rate  $\beta$ . This allows the network to acquire task-specific knowledge rapidly and learn across-task knowledge slowly. By alternately implementing the inner-task and across-task updates until the loss converges, the network learns model initialization that can adapt to a new task using only a small number of samples.

It should be mentioned that in [29], the across-task update is preformed by stochastic gradient descent (SGD) like the inner-task update. Here we adopt the ADAM algorithm instead of SGD since the ADAM benefits from the moment technique and the parameter-specific learning rates, and thus significantly outperforms SGD in most situations [39].

### C. Meta-adaption and Testing Stages

Following Section V-A, we generate the adaption datasets  $\{\mathbb{D}_{\text{Ad}}(k)\}_{k=1}^{K_T}$  and the testing datasets  $\{\mathbb{D}_{\text{Te}}(k)\}_{k=1}^{K_T}$  for the meta-adaption and testing stages, respectively. In the meta-adaption stage, the meta-learning algorithm aims to fine-tune the network on the adaption dataset  $\mathbb{D}_{\text{Ad}}(k)$  so that the network can predict the downlink CSI on  $\mathbb{D}_{\text{Te}}(k)$  as accurate as possible. For the  $k$ -th target task, the direct-transfer algorithm initializes the target-task-specific parameter  $\Omega_{T,k}$  as the meta-trained network parameter  $\Omega_{\text{Mt}}$ , and then fine-tunes  $\Omega_{T,k}$  on the adaption dataset  $\mathbb{D}_{\text{Ad}}(k)$  via  $G_{\text{Ad}}$  steps of gradient descents. The optimization objective is the loss function on the  $\mathbb{D}_{\text{Ad}}(k)$ , i.e.,  $\text{Loss}_{\mathbb{D}_{\text{Ad}}(k)}(\Omega_{T,k})$ . Then, the truncated target-task-specific gradient  $\mathbf{v}_{T,k}$  can be obtained by

$$[\mathbf{v}_{T,k}]_p = \min \left\{ \Upsilon, [\nabla_{\Omega_{T,k}} \text{Loss}_{\mathbb{D}_{\text{Ad}}(k)}(\Omega_{T,k})]_p \right\},$$

$$p = 1, \dots, \text{len}(\mathbf{v}_{T,k}). \quad (23)$$

Then, the parameter  $\Omega_{T,k}$  is updated based on the truncated gradient  $\mathbf{v}_{T,k}$ , i.e.,

$$\Omega_{T,k} \leftarrow \Omega_{T,k} - \beta \mathbf{v}_{T,k}. \quad (24)$$

After the fine-tuning finished, the target-task-specific parameter  $\Omega_{T,k}$  will be fixed. Let  $\text{NMSE}_{\text{Mt}}(k)$  represent the NMSE of the meta-learning algorithm evaluated on the  $k$ -th target environment that can be obtained by testing the network on the dataset  $\mathbb{D}_{\text{Te}}(k)$  using Eqs. (8) and (14). Then, the prediction NMSE of the meta-learning algorithm can be obtained by  $\text{NMSE}_{\text{Mt}} = \sum_{k=1}^{K_T} \text{NMSE}_{\text{Mt}}(k) / K_T$ . The concrete steps of the meta-learning algorithm are given in **Algorithm 3**.

Fig. 4 illustrates the trajectories of parameter updating during the meta-training and meta-adaption stages of the meta-learning algorithm. The black solid line represents the updating trajectory of the network parameter  $\Omega$  during the meta-training stage. The blue dashed lines represent the updating trajectories of the target-task-specific parameters  $\{\Omega_{T,k}\}_{k=1}^3$  during the meta-adaption stage. The green dotted lines represent the trajectories of the source-task-specific parameters  $\{\Omega_{S,k}\}_{k=1}^6$  during the inner-task updates. As shown in Fig. 4,  $K_B$  red dash-dotted lines point to the next position of parameter  $\Omega$ , which implies that  $K_B$  trained source-task-specific parameters are involved in the across-task update to determine the updating direction of  $\Omega$ . As depicted in Fig. 3 and Fig. 4, the

---

### Algorithm 3: Meta-learning algorithm for downlink CSI prediction

---

**Input:** Source tasks:  $\{\mathcal{T}_S(k)\}_{k=1}^{K_S}$ , Target tasks:  $\{\mathcal{T}_T(k)\}_{k=1}^{K_T}$ , inner-task learning rate:  $\beta$ , across-task learning rate:  $\gamma$ , exponential decay rates for across-task learning:  $\{\rho_1, \rho_2\}$ , disturbance factor:  $\varepsilon$ , number of gradsteps for inner-task training:  $T_{\text{Tr}}$ , number of tasks in each time index:  $K_B$

**Output:** Meta-trained network parameter:  $\Omega_{\text{Mt}}$ , Predicted downlink CSI based on  $\{\mathbb{D}_{\text{Te}}(k)\}_{k=1}^{K_T}$ , NMSE of the meta-learning algorithm:  $\text{NMSE}_{\text{Mt}}$

```

1 Meta-training stage
2 Randomly initialize the network parameters  $\Omega$ 
3 Initialize the time step:  $t \leftarrow 0$ 
4 Initialize the 1st and 2nd moment vectors:  $\mu, \nu \leftarrow \mathbf{0}$ 
5 while not done do
6   Randomly select  $K_B$  tasks from  $\{\mathcal{T}_S(k)\}_{k=1}^{K_S}$ 
7   Generate corresponding support dataset
      $\{\mathbb{D}_{\text{TrSup}}(k)\}_{k=1}^{K_B}$  and query dataset  $\{\mathbb{D}_{\text{TrQue}}(k)\}_{k=1}^{K_B}$ 
8   for  $k = 1, \dots, K_B$  do
9     Initialize the parameter  $\Omega_k \leftarrow \Omega$ 
10    for  $g = 1, \dots, G_{\text{Tr}}$  do
11      Obtain truncated gradient  $\mathbf{v}_{S,k}$  using Eq. (18)
12      Update the parameter  $\Omega_{S,k}$  with gradient
        descent:  $\Omega_{S,k} \leftarrow \Omega_{S,k} - \beta \mathbf{v}_{S,k}$ 
13    end
14  end
15   $t \leftarrow t + 1$ 
16  Update biased 1st and 2nd moment vectors using
    Eq. (20)
17  Compute unbiased 1st and 2nd moment vectors using
    Eq. (21)
18  Update the parameter  $\Omega$  unbiased 1st and 2nd moment
    vectors:  $\Omega \leftarrow \Omega - \eta \bar{\mu}_3 / (\sqrt{\bar{\nu}_3} + \varepsilon)$ 
19 end
20 Meta-trained network parameter:  $\Omega_{\text{Mt}} \leftarrow \Omega$ 
21 Meta-adaption and Testing
22 Initialize NMSE:  $\text{NMSE}_{\text{Mt}} \leftarrow 0$ 
23 for  $k = 1, \dots, K_T$  do
24   Generate the datasets  $\mathbb{D}_{\text{Ad}}(k)$  and  $\mathbb{D}_{\text{Te}}(k)$  for  $\mathcal{T}_T(k)$ 
25   Meta-adaption stage
26   Load the network parameter  $\Omega_{T,k} \leftarrow \Omega_{\text{Mt}}$ 
27   for  $g = 1, \dots, G_{\text{Ad}}$  do
28     Obtain truncated gradient  $\mathbf{v}_{T,k}$  using Eq. (23)
29     Update the parameter  $\Omega_{T,k}$  with gradient descent:
        $\Omega_{T,k} \leftarrow \Omega_{T,k} - \beta \mathbf{v}_{T,k}$ 
30   end
31   Testing stage
32   Predict the downlink CSI base on  $\mathbb{D}_{\text{Te}}(k)$  and  $\Omega_{T,k}$ 
    using Eq. (8)
33   Calculate  $\text{NMSE}_{\text{Mt}}(k)$  using Eq. (14)
34    $\text{NMSE}_{\text{Mt}} \leftarrow \text{NMSE}_{\text{Mt}} + \text{NMSE}_{\text{Mt}}(k) / K_T$ 
35 end
```

---



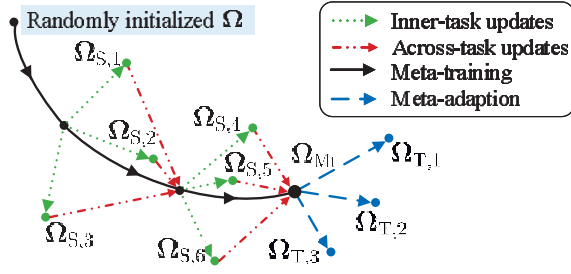


Fig. 4. Illustration of the meta-learning algorithm, where  $K_B = K_T = 3$ .

meta-learning algorithm takes the unique features of different tasks into account and exploits the joint guidance of multiple source tasks while the direct-transfer algorithm simply regards all the source tasks as one source task and ignores the structure information in different tasks.

## VI. SIMULATION RESULTS

Unless otherwise specified, the system parameters are set as follows: The BS is equipped with 64 antennas. To generate the training dataset, the uplink frequency is randomly selected from  $[1, 3]$  GHz. The frequency difference between the uplink and the downlink is 120 MHz. We approximate the channel in Eq. (1) by 200 discrete propagation paths. The incident AS is  $20^\circ$ . For different environments, the power of attenuations and AS are different. For different users, the attenuation of each path follows Rayleigh distribution. The phase and delay of each path follow uniform distribution over  $[-\pi, \pi)$  and  $[0, 10^{-4}]s$ . For the same user, the delays and DOAs of the uplink and downlink channels are the same.

The meta-learning, the no-transfer and the direct-transfer algorithms are all implemented on one computer with one Nvidia GeForce GTX 1080 Ti GPU. TensorFlow 1.4.0 is employed as the deep learning framework. The parameters of the no-transfer, the direct-transfer, and the meta-learning algorithms are given in Tab. I. The numbers of neurons in the input and the output layers are consistent with the lengths of input and output data vectors, respectively. The trainable parameters of FNN are randomly initialized as truncated normal variables<sup>6</sup> with normalized variance<sup>7</sup>. To be fair, the number of training samples in  $\mathbb{D}_{STr}$  are equal to the number of training samples used in the meta-learning algorithm, i.e., the total number of samples in the training set  $K_S \times N_{Tr} = 8000$ .

Fig. 5 depicts the NMSE performance of the no-transfer, the direct-transfer, and the meta-learning algorithms (i.e.,  $NMSE_{Nt}$ ,  $NMSE_{Dt}$  and  $NMSE_{Mt}$ ) versus the number of gradsteps  $G_{Ad}$ . The NMSE performance is the average accuracy of channel estimation from 100 target tasks, i.e., obtained by repetitively (fine-tuning and then) testing 100 different target environments. Since the network parameters of the no-transfer algorithm do not fine-tune based on the target environment, the accuracy curve of the no-transfer algorithm is a horizontal line. As shown in Fig. 5, the prediction accuracy

TABLE I  
DEFAULT PARAMETERS FOR THE NO-TRANSFER, THE DIRECT-TRANSFER, AND THE META-LEARNING ALGORITHMS

Parameter	Value
Number of neurons in hidden layers	(128, 128)
Learning rates: $(\gamma, \beta, \eta)$ ,	(5e-4, 1e-3, 1e-3)
Exponential decay rates for ADAM: $(\rho_1, \rho_2)$	(0.9, 0.999)
Disturbance factor for ADAM: $\varepsilon$	1e-08
Number of source tasks: $K_S$	400
Number of tasks in each training index: $K_B$	10
Number of samples in each source task: $N_{Tr}$	20
Number of samples in the target task: $(N_{Ad}, N_{Te})$	(20, 20)
Number of users in each task: $U$	2
Number of gradsteps: $(T_{Tr}, G_{Ad})$	(10, 5e3)
Batch size: $V$	128

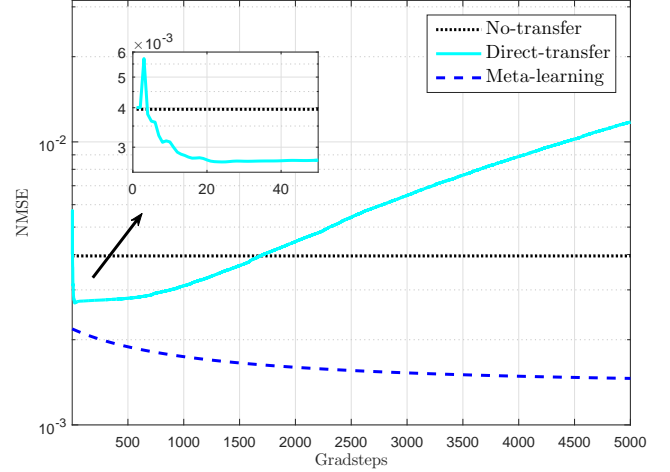


Fig. 5. The NMSE performance of the no-transfer, the direct-transfer, and the meta-learning algorithms versus the number of gradsteps  $G_{Ad}$ , where the number of samples in  $\mathbb{D}_{AdSup}$  is 20.

of the direct-transfer algorithm first improves as  $G_{Ad}$  increases, then degrades rapidly as  $G_{Ad}$  increases, and is even worse than the no-transfer algorithm when  $G_{Ad}$  is larger than 1700. This implies that the direct-transfer algorithm severely overfits as  $G_{Ad}$  increases. As shown in the partial enlarged picture, the optimal performance of the direct-transfer algorithm is achieved when  $G_{Ad}$  is equal to 25. To circumvent the problem of overfitting, the direct-transfer algorithm can stop optimization when the prediction accuracy no longer improves. The optimal prediction accuracy of the direct-transfer algorithm is achieved with different values of  $G_{Ad}$  in different simulation conditions. In the following simulations, we use the optimal performance of the direct-transfer algorithm as the benchmark. Furthermore, the optimal performance of the direct-transfer algorithm is significantly better than the no-transfer algorithm, which indicates that additional samples for the target task and appropriate fine-tuning can improve the prediction accuracy, i.e., transfer learning is more suitable than classical deep learning in downlink channel prediction. It also can be noticed that the performance of the meta-learning algorithm without meta-adaption stage is significantly better than the no-transfer learning algorithm, which demonstrates that the meta-learning algorithm can find a better initialization for fast adaption than the no-transfer learning algorithm. Moreover, the performance

<sup>6</sup>The truncated normal distribution is a normal distribution bounded by two standard deviations from the mean.

<sup>7</sup>The weights of neurons in the  $l$ -th layer are initialized as truncated normal variables with variance  $1/n_l$ .

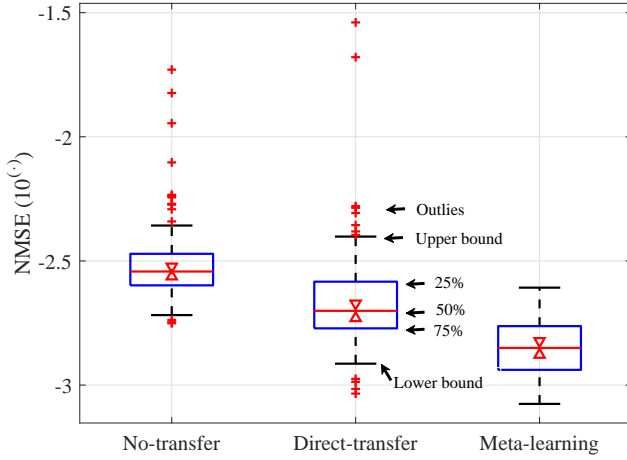


Fig. 6. The stability comparisons of the no-transfer, the direct-transfer, and the meta-learning algorithms.

of the meta-learning algorithm tends to improve as  $G_{Ad}$  increases while the direct-transfer algorithm tends to overfit as  $G_{Ad}$  increases, which also indicates that a better initialization can avoid overfitting effectively.

Fig. 6 compares the accuracy stabilities of the no-transfer, the direct-transfer, and the meta-learning algorithms, where the robustness over different environments are considered. By repetitively (fine-tuning and then) testing 100 new environments, we can obtain the optimal prediction accuracies of the meta-learning and direct-transfer algorithms for each new environment. As shown in Fig. 6, the boxplot displays the distribution of prediction accuracies with respect to 100 environments with six indexes, i.e., “outliers”, “upper bound”, “first quartile”, “median”, “third quartile”, and “lower bound”. The prediction accuracies of 50 environments range between the “lower bound” and “upper bound”, i.e., within the range of the blue box. As shown in Fig. 6, the red “+” denotes the “outliers”. The algorithm with more outliers exhibits more unstable performance. In the no-transfer algorithm, the prediction NMSEs of users in 11 environments is higher than the corresponding upper bound NMSE. In the direct-transfer algorithm, the prediction NMSEs of users in 8 environments is higher than the corresponding upper bound NMSE. There is no outliers in the meta-learning algorithm, which demonstrate its remarkable robustness and stability over different environments.

Fig. 7 depicts the NMSE performance of the no-transfer, the direct-transfer, and the meta-learning algorithms versus the number of samples in  $\mathbb{D}_{AdSup}$ . For each sample size, the NMSE performance is the average accuracy of channel estimation from 100 environments. Since the network parameters of the no-transfer algorithm do not use the dataset  $\mathbb{D}_{AdSup}$ , the accuracy curve of the no-transfer algorithm is a horizontal line. It can be seen from Fig. 7 that the prediction accuracies of the direct-transfer, and the meta-learning algorithms both improves as the sample size increases. Compared with the no-transfer algorithm, the improvement of direct-transfer algorithm is negligible when the sample size  $\mathbb{D}_{AdSup}$  is 10. However, when the sample size is larger than 110, the perfor-

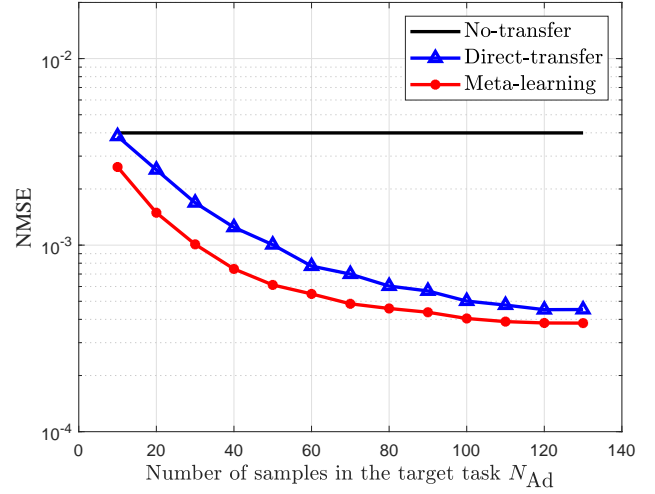


Fig. 7. The NMSE performance of the no-transfer, the direct-transfer, and the meta-learning algorithms versus the number of samples in  $\mathbb{D}_{AdSup}$ .

mance of the direct-transfer, and the meta-learning algorithms both become saturated, which indicates that the sample size is large enough to fine-tune the network. Moreover, the meta-learning algorithm outperforms the direct-transfer algorithm, especially when the sample size is small, which implies its superiority in few-sample learning.

Fig. 8 shows the NMSE performance of the no-transfer, the direct-transfer, and the meta-learning algorithms versus the frequency difference  $\Delta f$ . For each frequency difference point, the NMSE performance is the average accuracy of channel estimation from 100 environments. The prediction accuracies of the three algorithms all degrades as the frequency difference increases. This is because the channel correlation between the uplink and the downlink tends to vanish as the frequency difference increases. The meta-learning algorithm outperforms the no-transfer and the direct-transfer algorithms, and the performance gaps becomes narrower as the frequency difference increases. When the frequency difference is small than 400 MHz, the prediction accuracies of both the meta-learning and the direct-transfer algorithms are higher than  $10^{-2}$ , which validates the effectiveness of the meta-learning and the direct-transfer algorithms.

## VII. CONCLUSION

In this paper, we formulated the downlink channel prediction for FDD massive MIMO systems as a deep transfer learning problem, where each learning task represents the downlink CSI prediction from the uplink CSI for a certain environment. Then, we proposed the no-transfer, direct-transfer and meta-learning algorithms based on the fully-connected neural network architecture. The no-transfer algorithm trains the network in the classical deep learning manner. The direct-transfer algorithm fine-tunes the network based on the initialization of the no-transfer algorithm. The meta-learning algorithm learns a model initialization that can effectively adapt to a new environment with a small amount of labeled data. Simulation results have shown that the meta-learning

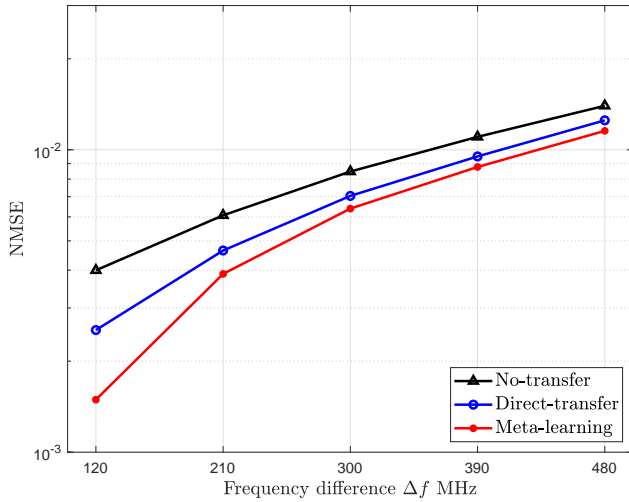


Fig. 8. The NMSE performance of the no-transfer, the direct-transfer, and the meta-learning algorithms versus the frequency difference  $\Delta f$ .

algorithm significantly outperforms the direct-transfer and the no-transfer algorithms in terms of both accuracy and stability, especially when the number of samples are very small, which demonstrates its effectiveness and superiority.

## REFERENCES

- [1] B. Wang, F. Gao, S. Jin, H. Lin, and G. Y. Li, "Spatial- and frequency-wideband effects in millimeter-wave massive MIMO systems," *IEEE Trans. Signal Process.*, vol. 66, no. 13, pp. 3393–3406, Jul. 2018.
- [2] H. Xie, F. Gao, S. Zhang, and S. Jin, "A unified transmission strategy for TDD/FDD massive MIMO systems with spatial basis expansion model," *IEEE Trans. Vehicular Technol.*, vol. 66, no. 4, pp. 3170–3184, Apr. 2017.
- [3] B. Wang, F. Gao, S. Jin, H. Lin, G. Y. Li, S. Sun, and T. S. Rappaport, "Spatial-wideband effect in massive MIMO with application in mmwave systems," *IEEE Commun. Mag.*, vol. 56, no. 12, pp. 134–141, Dec. 2018.
- [4] J. Hoydis, C. Hoek, T. Wild, and S. ten Brink, "Channel measurements for large antenna arrays," in *Proc. Int. Symp. Wireless Commun. Systems (ISWCS)*, Paris, France, Aug. 2012, pp. 811–815.
- [5] Y. Zhou, M. Herdin, A. Sayeed, and E. Bonek, "Experimental study of MIMO channel statistics and capacity via the virtual channel representation," *Univ. Wisconsin-Madison, Madison, WI, USA, Tech. Rep.*, vol. 5, pp. 10–15, 2007.
- [6] K. Hugl, K. Kalliola, and J. Laurila, "Spatial reciprocity of uplink and downlink radio channels in FDD systems," in *Proc. COST*. Citeseer, 2002, vol. 273, p. 066.
- [7] Y. Han, Q. Liu, C. Wen, S. Jin, and K. Wong, "FDD massive MIMO based on efficient downlink channel reconstruction," *IEEE Trans. Commun.*, pp. 1–1, 2019.
- [8] H. Xie, F. Gao, S. Jin, J. Fang, and Y. Liang, "Channel estimation for TDD/FDD massive MIMO systems with channel covariance computing," *IEEE Trans. Wireless Commun.*, vol. 17, no. 6, pp. 4206–4218, Jun. 2018.
- [9] X. Rao and V. K. N. Lau, "Distributed compressive CSIT estimation and feedback for FDD multi-user massive MIMO systems," *IEEE Trans. Signal Process.*, vol. 62, no. 12, pp. 3261–3271, Jun. 2014.
- [10] M. E. Eltayeb, T. Y. Al-Naffouri, and H. R. Bahrani, "Compressive sensing for feedback reduction in MIMO broadcast channels," *IEEE Trans. Commun.*, vol. 62, no. 9, pp. 3209–3222, Sep. 2014.
- [11] Z. Qin, H. Ye, G. Y. Li, and B. F. Juang, "Deep learning in physical layer communications," *IEEE Wireless Commun.*, vol. 26, no. 2, pp. 93–99, Apr. 2019.
- [12] Y. Yang, F. Gao, X. Ma, and S. Zhang, "Deep learning-based channel estimation for doubly selective fading channels," *IEEE Access*, vol. 7, pp. 36579–36589, Mar. 2019.
- [13] H. He, C. Wen, S. Jin, and G. Y. Li, "Deep learning-based channel estimation for beamspace mmwave massive MIMO systems," *IEEE Wireless Commun. Lett.*, vol. 7, no. 5, pp. 852–855, Oct. 2018.
- [14] Z. Jia, W. Cheng, and H. Zhang, "A partial learning based detection scheme for massive MIMO," *IEEE Wireless Commun. Lett.*, pp. 1–1, Apr. 2019.
- [15] H. Ye, G. Y. Li, and B. Juang, "Power of deep learning for channel estimation and signal detection in OFDM systems," *IEEE Wireless Commun. Lett.*, vol. 7, no. 1, pp. 114–117, Feb. 2018.
- [16] C. Wen, W. Shih, and S. Jin, "Deep learning for massive MIMO CSI feedback," *IEEE Wireless Commun. Lett.*, vol. 7, no. 5, pp. 748–751, Oct. 2018.
- [17] S. Jin W. Xu, F. Gao and A. Alkhateeb, "3D scene based beam selection for mmwave communications," *arXiv preprint arXiv:1911.08409*, 2019.
- [18] M. Alrabeiah, A. Hredzak, Z. Liu, and A. Alkhateeb, "ViWi: A deep learning dataset framework for vision-aided wireless communications," *arXiv preprint arXiv:1911.06257*, 2019.
- [19] A. Alkhateeb, O. El Ayach, G. Leus, and R. W. Heath, "Channel estimation and hybrid precoding for millimeter wave cellular systems," *IEEE J. Sel. Topics Signal Process.*, vol. 8, no. 5, pp. 831–846, Oct. 2014.
- [20] X. Li and A. Alkhateeb, "Deep learning for direct hybrid precoding in millimeter wave massive MIMO systems," in *Proc. Asilomar*, *arXiv e-prints*, p. arXiv:1905.13212, May 2019.
- [21] Y. Yang, F. Gao, C. Qian, and G. Liao, "Model-aided deep neural network for source number detection," *arXiv preprint arXiv:1909.13273*, 2019.
- [22] M. Safari and V., "Deep UL2DL: Channel knowledge transfer from uplink to downlink," *arXiv preprint arXiv:1812.07518*, 2018.
- [23] P. Dong, H. Zhang, and G. Y. Li, "Machine learning prediction based CSI acquisition for FDD massive MIMO downlink," in *Proc. IEEE Global Commun. Conf. (GLOBECOM)*, Abu Dhabi, United Arab Emirates, Dec. 2018, pp. 1–6.
- [24] M. Alrabeiah and A. Alkhateeb, "Deep learning for TDD and FDD massive MIMO: Mapping channels in space and frequency," *arXiv preprint arXiv:1905.03761*, 2019.
- [25] Y. Yang, F. Gao, G. Y. Li, and M. Jian, "Deep learning based downlink channel prediction for FDD massive MIMO system," *IEEE Commun. Lett.*, pp. 1–1, 2019.
- [26] S. J. Pan and Q. Yang, "A survey on transfer learning," *IEEE Trans. Knowledge and Data Engineering*, vol. 22, no. 10, pp. 1345–1359, Oct. 2010.
- [27] J. Lin, Y. Wang, Y. Xia, T. He, and Z. Chen, "Learning to transfer: Un-supervised meta domain translation," *arXiv preprint arXiv:1906.00181*, 2019.
- [28] C. Tan, F. Sun, T. Kong, W. Zhang, C. Yang, and C. Liu, "A survey on deep transfer learning," in *Proc. Int. Conf. Artificial Neural Networks*. Springer, 2018, pp. 270–279.
- [29] C. Finn, P. Abbeel, and S. Levine, "Model-agnostic meta-learning for fast adaptation of deep networks," in *Proc. 34th Int. Conf. Machine Learning*. JMLR. org, 2017, vol. 70, pp. 1126–1135.
- [30] M. Andrychowicz, M. Denil, S. Colmenarejo, M. W. Hoffman, D. Pfau, T. Schaul, B. Shillingford, and N. de Freitas, "Learning to learn by gradient descent by gradient descent," in *Proc. 30th Int. Conf. Neural Inf. Process. Systems*, Barcelona, Spain, 2016, pp. 3988–3996.
- [31] A. Santoro, S. Bartunov, M. Botvinick, D. Wierstra, and T. Lillicrap, "Meta-learning with memory-augmented neural networks," in *Proc. Int. Conf. machine learning (ICML)*, 2016, pp. 1842–1850.
- [32] S. Ravi and H. Larochelle, "Optimization as a model for few-shot learning," in *Proc. Int. Conf. Learning Representations (ICLR)*, Toulon, France, 2017.
- [33] K. Hornik, M. Stinchcombe, and H. White, "Multilayer feedforward networks are universal approximators," *Neural netw.*, vol. 2, no. 5, pp. 359–366, 1989.
- [34] D. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.
- [35] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*, MIT Press, 2016, <http://www.deeplearningbook.org>.
- [36] S. Liang, R. Sun, Y. Li, and R. Srikant, "Understanding the loss surface of neural networks for binary classification," *arXiv preprint arXiv:1803.00909*, 2018.
- [37] I. Sutskever, J. Martens, G. Dahl, and G. Hinton, "On the importance of initialization and momentum in deep learning," in *Proc. Int. Conf. Machine Learning (ICML)*, Atlanta, USA, 2013, pp. 1139–1147.
- [38] J. Langford, L. Li, and T. Zhang, "Sparse online learning via truncated gradient," *J. Machine Learning Research*, vol. 10, no. Mar., pp. 777–801, 2009.
- [39] J. Chen and A. Kyrillidis, "Decaying momentum helps neural network training," *arXiv preprint arXiv:1910.04952*, 2019.