Learning to Continuously Optimize Wireless Resource In Episodically Dynamic Environment

Haoran Sun, Wenqiang Pu, Minghe Zhu, Xiao Fu, Tsung-Hui Chang, and Mingyi Hong * Submitted to ICASSP 2021 on October 21, 2020

Abstract

There has been a growing interest in developing data-driven and in particular deep neural network (DNN) based methods for modern communication tasks. For a few popular tasks such as power control, beamforming, and MIMO detection, these methods achieve state-of-the-art performance while requiring less computational efforts, less channel state information (CSI), etc. However, it is often challenging for these approaches to learn in a dynamic environment where parameters such as CSIs keep changing.

This work develops a methodology that enables data-driven methods to continuously learn and optimize in a dynamic environment. Specifically, we consider an "episodically dynamic" setting where the environment changes in "episodes", and in each episode the environment is stationary. We propose to build the notion of continual learning (CL) into the modeling process of learning wireless systems, so that the learning model can incrementally adapt to the new episodes, without forgetting knowledge learned from the previous episodes. Our design is based on a novel min-max formulation which ensures certain "fairness" across different data samples. We demonstrate the effectiveness of the CL approach by customizing it to two popular DNN based models (one for power control and one for beamforming), and testing using both synthetic and real data sets. These numerical results show that the proposed CL approach is not only able to adapt to the new scenarios quickly and seamlessly, but importantly, it maintains high performance over the previously encountered scenarios as well.

1 Introduction

Deep learning has been successful in many applications such as computer vision [2], natural language processing [3], and recommender systems [4]; see [5] for an overview. Recent works have also demonstrated that deep learning can be applied in communication systems, either by replacing an individual component in the system (such as signal detection [6,7], channel decoding [8], channel estimation [9,10], power allocation [11–15], beamforming [16,17] and wireless scheduling [18]), or by jointly optimizing the entire system [19,20], for achieving state-of-the-art performance. Specifically, deep learning is a data-driven method in which a large amount of training data is used to train a deep neural network (DNN) for a specific task (such as power

Processing (ICASSP 2021) on October 21, 2020.

^{*}H. Sun and M. Hong are with Department of Electrical and Computer Engineering, University of Minnesota, Minneapolis, MN 55455, USA. W. Pu is with the Shenzhen Research Institute of Big Data, Shenzhen, China. M. Zhu and T-H. Chang are with Shenzhen Research Institute of Big Data, The Chinese University of Hong Kong, Shenzhen, China. X. Fu is with the School of Electrical Engineering and Computer Science, Oregon State University, Corvallis, OR 97331, USA. Fu is supported by NSF/Intel MLWiNS: CNS-2003082 and M. Hong is supported by NSF/Intel MLWiNS: CNS-2003033.

A short version of this paper [1] has been submitted to the 2021 IEEE International Conference on Acoustics, Speech, & Signal

control). Once trained, such a DNN model will replace conventional algorithms to process data in real time. Existing works have shown that when the real-time data follows similar distribution as the training data, then such an approach can generate high-quality solutions for non-trivial wireless tasks [6–18], while significantly reducing real-time computation, and/or requiring only a subset of channel state information (CSI).

Dynamic environment. However, it is often challenging to use these DNN based algorithms when the environment (such as CSI and user locations) keeps changing. We list a few reasons below.

- 1) It is well-known that these methods typically suffer from severe performance deterioration when the environment changes, that is, when the real-time data follows a different distribution than those used in the training stage [15].
- 2) One can adopt the transfer learning and/or online learning paradigm, by updating the DNN model according to data generated from the new environment [15]. However, these approaches usually degrade or even overwrite the previously learned models [21,22]. Therefore they are sensitive to outliers because once adapted to a transient environment/task, its performance on the existing environment/task can degrade significantly [23]. Such kind of behavior is particularly undesirable for wireless resource allocation tasks, because the randomness of a typical wireless environment can result in an unstable learning model.
- 3) If one chooses to periodically retrain the entire DNN using all the data seen so far [23], then the training can be both time and memory consuming since the number of data needed keeps growing.

Due to these challenges, it is unclear how state-of-the-art DNN based communication algorithms could properly adapt to new environments quickly without experiencing significant performance loss over previously encountered environments. This is one of the main obstacles preventing these data-driven methods from being implemented in real communication systems.

Ideally, one would like to design data-driven models that can adapt to the new environment *efficiently* (i.e., by using as little resource as possible), *seamlessly* (i.e., without knowing when the environment has been changed), *quickly* (i.e., adapt well using only a small amount of data), and *continuously* (i.e., without forgetting the previously learned models). Is achieving all the above goals simultaneously possible? If so, how to do it?

Continual Learning. In the machine learning community, continual learning (CL) has recently been proposed to address the "catastrophic forgetting phenomenon", that is, the tendency of abruptly losing the previously learned models when the current environment information is incorporated [22]. Specifically, consider the setting where different "tasks" (e.g., different CSI distributions) are revealed sequentially. Then CL aims to retain the knowledge learned from the early tasks through one of the following mechanisms: 1) regularize the most important parameters [23–25]; 2) design dynamic neural network architectures and associate neurons with tasks [26, 27]; or 3) introduce a small set of memory for later training rehearsal [28–30]. However, most of the above mentioned methods require the knowledge of the task boundaries, that is, the time stamp where an old task terminates and a new task begins. Unfortunately, such a setting does not suit wireless communication problems well, since the wireless environment usually changes gracefully and swiftly. So it is difficult for wireless systems to identify a precise changing point. Only limited recent CL works have focused on boundary-free environments [31–33], but they both propose general-purpose tools without considering any problem-specific structures in wireless communications. Further, the effectiveness of these approaches have only been tested on typical machine learning tasks such as image classification. It is unclear whether they will be effective in wireless communication tasks.

Contributions. The main contribution of this paper is that we introduce the notion of CL to datadriven wireless system design, and develop a CL formulation together with a training algorithm tailored for core tasks in wireless communications. Specifically, we consider an "episodically dynamic" setting

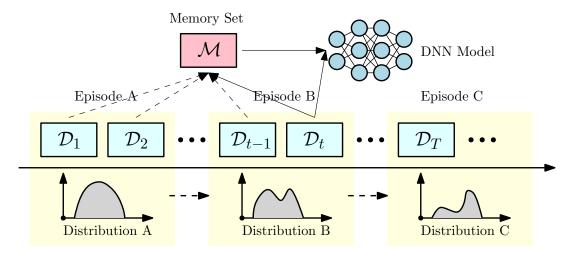


Figure 1: Proposed CL framework for the episodically dynamic environment. The data is feeding in a sequential manner (thus the system can only access D_t at time t) with changing episodes and distributions, and the model has a limited memory set \mathcal{M} (which cannot store all data D_1 to D_t). To maintain the good performance over all experienced data from D_1 to D_t , the proposed framework optimizes the data-driven model at each time t, based on the mixture of the current data D_t and the memory set \mathcal{M} . The memory set \mathcal{M} is then updated to incorporate the new data D_t .

where the environment changes in *episodes*, and within each episodes the distribution of the CSIs stays relatively stationary. Our goal is to design a learning model which can seamlessly and efficiently adapt to the changing environment, without knowing the episode boundary, and without incurring too much performance loss if a previously experienced episode reappears. We propose a CL framework for wireless systems, which incrementally adapts the DNN models by using data from the new episode as well as a limited but carefully selected subset of data from the previous episodes; see Fig. 1. Compared with the existing boundary-free CL algorithms (which are often general-purpose schemes without a formal optimization formulation) [31–33], our approach comes with a clearly defined and tailored optimization formulation that serves for the wireless resource allocation problem. In particular, our CL method is based on a min-max formulation which selects a small set of important data samples into the working memory according to a carefully designed *data-sample fairness* criterion. Moreover, we demonstrate the effectiveness of our proposed framework by customizing it to two popular DNN based models (one for power control and the other for multi-user beamforming). We test the proposed CL approach using both synthetic and real data. To advocate the reproducible research, the code of our implementation will be made available online at https://github.com/Haoran-S/ICASSP2021.

2 Literature Review

2.1 Deep learning (DL) for Wireless Communication (WC)

Recently, deep learning has been used to generate high-quality solutions for non-trivial wireless communication tasks [6–18]. These approaches can be roughly divided into following two categories:

2.1.1 End-to-end Learning

For the classic resource allocation problems such as power control, the work [11] shows that deep neural networks (DNNs) can be exploited to learn the optimization algorithms such as WMMSE [34], in an

end-to-end fashion. Subsequent works such as [12] and [13] show that unsupervised learning can be used to further improve the model performance. Different network structures, such as convolutional neural networks (CNN) [12] and graph neural networks (GNN) [14,35], and different modeling techniques, such as reinforcement learning [36], are also studied in the literature. Nevertheless, all the above mentioned methods belong to the category of end-to-end learning, where a black-box model (typically deep neural network) is applied to learn either the structure of some existing algorithms, or the optimal solution of a communication task.

2.1.2 Deep Unfolding

End-to-end learning based algorithms usually perform well, but their black-box nature makes such kind of models suffer from poor interpretability, i.e., the mechanism of black-box model is usually hard to understand. Further, their performance heavily depends on the quality and amount of the training samples. Alternatively, deep unfolding based methods [37] unfold existing optimization algorithms iteration by iteration and approximate the per-iteration behavior by one layer of the neural network. The rationale is that, by properly training the network parameters, the multi-layer neural networks can match the behavior of iterative optimization algorithms. The hope is that the trained network uses much fewer iterations (layers), while preserving good interpretability and generalization performance. In the machine learning community, well-known works in this direction include the unfolding of the iterative soft-thresholding algorithm (ISTA) [38], unfolding of the non-negative matrix factorization methods [39], and the unfolding of the alternating direction method of multipliers (ADMM) [40]. Recently, the idea of unfolding has been used in communication task such as MIMO detection [41–43], channel coding [44], resource allocation [45], channel estimation [46], and beamforming problems [47]; see a recent survey paper [37].

2.2 Continual Learning (CL)

Continual learning is originally proposed to improve reinforcement learning tasks [48] to help alleviate the catastrophic forgetting phenomenon, that is, the tendency of abruptly losing the knowledge about the previously learned task(s) when the current task information is incorporated [22]. It has later been broadly used to improve other machine learning models, and specifically the DNN models; see a recent review article on this topic [21]. Generally speaking, the CL paradigm can be roughly classified into the following three main categories.

2.2.1 Regularization Based Methods

Based on the Bayesian theory and inspired by synaptic consolidation in Neuroscience, the regularization based methods penalize the most important parameters to retain the performance on old tasks. Some most popular regularization approaches include Elastic Weight Consolidation (EWC) [23], Learning without Forgetting (LwF) [24], and Synaptic Intelligence (SI) [25]. However, regularization or penalty based methods naturally introduce trade-offs between the performance of old and new tasks. For example, if a large penalty is applied to prevent the model parameters from moving out of the optimal region of old tasks (i.e., a set of parameters that incur the lowest training/testing cost), the model may be hard to adapt to new tasks. On the other hand, if a small penalty is applied, it may not be sufficient to force the parameters to stay in the optimal region to retain the performance on old tasks. There are also criticisms on the performance loss of regularization based methods for a long chain of tasks [49,50].

2.2.2 Architectures Based Methods

By associating neurons with tasks (either explicitly or not), many different types of dynamic neural network architectures are proposed to address the catastrophic forgetting phenomenon [26,27]. However, due to the nature of the parameter isolation, architecture based methods usually require the knowledge of the task boundaries, and thus they are not suitable for wireless settings, where the environment change is often difficult to track.

2.2.3 Rehearsal Based Methods

Tracing back to the 1990s [51], the rehearsal based (aka. memory-based) methods play an important role in areas such as reinforcement learning [52]. As its name suggests, rehearsal based methods use a small set of samples for later training rehearsal, either through selecting and storing the most represented samples such as the incremental classifier, representation learning (iCaRL) [30] and gradient episodic memory (GEM) [28], or use generative models to generate representative samples [29]. However, both GEM and iCaRL require the knowledge of the task boundaries (i.e., where a new episode starts), which does not suit wireless communication problems well. The authors of [33] propose boundary-free methods by selecting the samples through random reservoir sampling, which fills the memory set with data that is uniformly randomly sampled from the streaming episodes. More complex mechanisms are also introduced recently to further increase the sampling diversity, where the diversity is measured by either samples' stochastic gradient directions [32] or samples' Euclidean distances [31,53].

2.3 Comparison of CL with other methods

In this section, we will introduce a few related methods which also deal with streaming data, and compare them with the CL approach.

2.3.1 Online Learning

As opposed to batch learning, online learning deals with the learning problems where the training data comes sequentially, and data distribution over time may or may not be consistent [54]. The ultimate goal of online learning is to minimize the cumulative loss over time, utilizing the previously learned knowledge. In particular, when data sampling is independently and identically distributed, online gradient descent is essentially the stochastic gradient descent method, and all classic complexity results can be applied. On the other hand, when data sampling is non-stationary and drifts over time, online learning methods are more likely to adapt to the most recent data but losing the performance on past data, i.e., suffering from the catastrophic forgetting [55].

2.3.2 Transfer Learning

Different from online learning, transfer learning is designed to apply the knowledge gained from one problem to another problem, based on the assumption that related tasks will benefit each other, see a complete survey article [56]. By transferring the learned knowledge from old tasks to new tasks, transfer learning can quickly adapt to new tasks with fewer samples and less labeling effort. A typical application is the model fine-tuning on a (potentially small) user-specific problem (e.g., MNIST classification) based on some offline pre-trained model using a comprehensive dataset (e.g. ImageNet dataset). By applying the gained knowledge from the original dataset (e.g. ImageNet dataset), the model can adapt to the new dataset

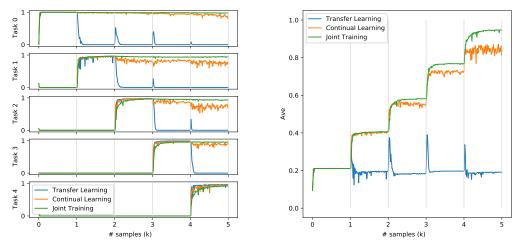


Figure 2: MNIST handwritten digits classification example: (left) testing accuracy on each task, which contains two out of ten digits (with no overlapping digits between tasks), (right) testing accuracy on all ten digits. We can observe that naive transfer learning on one task will significantly degrade the performance on the rest of the tasks, while joint training using accumulated dataset can preforms well on all experienced tasks.

(e.g., MNIST classification) more quickly with fewer samples. Similar ideas have been applied in wireless settings recently [15, 57, 58] to deal with scenarios that network parameters changes. However, since the model is purely fine-tuned or optimized on the new dataset, after the knowledge transfer, the knowledge from the original model may be altered or overwritten, resulting in significant performance deterioration on the original problem [23].

2.3.3 MNIST Example

In summary, though being successful in many application scenarios, both transfer learning and online learning are prone to deviate from the original models. We illustrate this phenomenon using the MNIST data set on a handwritten digits classification task. We first split the training data into five disjoint subset, each containing two out of ten digits and 1,000 samples. We then train a deep neural network model using sequential data (progressively available 10 samples per time), one task after another. Three training strategies are compared in our simulation, they are transfer learning based on fine-tuning (update the model using previous model as initialization), continual learning (introduce a small memory set for later training rehearsal) and joint training (training from scratch using all accumulated data). Since each time only 10 samples are available to the model while each task contains 1,000 samples, the transfer learning implementation can also be viewed as an online learning method. Testing accuracy on each task and on the entire testing set are reported in Fig. 2. One can observe that the naive online learning and/or transfer learning on one task can easily destroy the previously learned parameters and can result in significant performance degradation on the rest of the tasks, while joint training using accumulated dataset can preforms well on all experienced tasks. Meanwhile, CL can serve as an efficient alternative to joint training, with much reduced computation and storage overhead.

3 The Episodic Wireless Environment

The focus of this paper is to design learning algorithms in a *dynamic* wireless environment, so that the data-driven models we built can seamlessly, efficiently, quickly and continuously adapt to new tasks. This section provides details about our considered *dynamic* environment, and discuss potential challenges.

Specifically, we consider an "episodically dynamic" setting where the environment changes relatively slowly in "episodes", and during each episode the learners observe multiple batches of samples generated from the same stationary distribution; see Fig. 1. We will use \mathcal{D}_t to denote a small batch of data collected at time t, and assume that each episode k contains a set of T_k batches, and use $\mathcal{E}_k = \{\mathcal{D}_t\}_{t \in T_k}$ to denote the data collected in episode k. To have a better understanding about the setting of this work, let us consider the following example.

3.1 A Motivating Example

Suppose a collection of base stations (BSs) run certain DNN based resource allocation algorithm to provide coverage for a given area, e.g., a shopping mall. The users' activities in the area can contain two types of dynamic patterns: 1) regular but gradually changing patterns – such as regular daily commute for the employees and customers, and such pattern could slowly change from week to week (e.g., the store that people like to visit in the summer is different in winter); 2) irregular but important patterns – such as large events (e.g., promotion during the anniversary season), during which the distribution of user population (and thus the CSI distribution) will be significantly different compared with their usual distributions, and more careful resource allocation has to be performed. The episode, in this case, can be defined as "an unusual period of time that includes a particular event", or a "usual period of time".

Suppose that each BS solves a weighted sum-rate (WSR) maximization problem for single-input single-output (SISO) interference channel, with a maximum of K transmitter and receiver pairs. Let $h_{kk} \in \mathbb{C}$ denote the direct channel between transmitter k and receiver k, and $h_{kj} \in \mathbb{C}$ denote the interference channel from transmitter j to receiver k. The power control problem aims to maximize the weighted system throughput via allocating each transmitter's transmit power p_k . For a given snapshot of the network, the problem can be formulated as the following:

$$\max_{p_1, \dots, p_K} R(\mathbf{p}; \mathbf{h}) := \sum_{k=1}^K \alpha_k \log \left(1 + \frac{|h_{kk}|^2 p_k}{\sum_{j \neq k} |h_{kj}|^2 p_j + \sigma_k^2} \right)
\text{s.t.} \quad 0 \le p_k \le P_{\text{max}}, \ \forall \ k = 1, 2, \dots, K,$$
(1)

where P_{max} denotes the power budget of each transmitter; $\{\alpha_k > 0\}$ are the weights. Problem (1) is known to be NP-hard [59] but can be effectively approximated by many optimization algorithms [34]. The data-driven methods proposed in recent works such as [11–15] train DNNs using some pre-generated dataset. Here \mathcal{D}_t can include a mini-batch of channels $\{h_{kj}\}$, and each episode can include a period of time where the channel distribution is stationary.

At the beginning of each week, a DNN model for solving problem (1) (pretrained using historical data, say \mathcal{D}_0) is preloaded on the BSs to capture the regular patterns in the shopping mall area. Note that one can preload multiple models on each BS, say one model for every six hours. The question is, what should the BSs do when the unexpected but irregular patterns appear? Specifically, say at 7:00 am a model is loaded to allocate resources up until noon. During this time the BSs can collect batches of data $\mathcal{D}_t, t = 1, 2, \cdots$. Then shall the BS update its model during the six hour period to capture the dynamics

of the user/demand distribution? If so, shall we use the entire data set, including the historical data and the real-time data, to re-train the neural network (which can be time-consuming), or shall we use transfer learning to adapt to the new environment on the fly (which may result in overwritten the basic 7:00 am model)?

To address the above issues, we propose to adopt the notion of continual learning, so that our model can incorporate the new data \mathcal{D}_t on the fly, while keeping the knowledge acquired from $\mathcal{D}_{0:t-1}$. In the next section, we will detail our proposed CL formulation to achieve such a goal.

4 CL for Learning Wireless Resource Allocation

4.1 Memory-based CL

Our proposed method is based upon the notion of the memory-based CL proposed in [28, 31, 32], which allows the learner to collect a small subset of historical data for future re-training. The idea is, once \mathcal{D}_t is received, we fill in a memory \mathcal{M}_t (with fixed size) with the most representative samples from all experienced episodes $\mathcal{D}_{0:t-1}$, and then train the neural network at each time t with the data $\mathcal{M}_t \cup \mathcal{D}_t$. Several major features of this approach are listed below:

- The learner does not need to know where a new episode starts (that is, the boundary-free setting) it can keep updating \mathcal{M}_t and keep training as data comes in.
- If one can control the size of the memory well (say comparable to the size of each \mathcal{D}_t), then the training complexity will be made much smaller than performing a complete training over the entire data set $\mathcal{D}_{0:t}$, and will be comparable with transfer learning approach which uses \mathcal{D}_t .
- If the size of a given data batch \mathcal{D}_t is very small, the learner is unlikely to *overfit* because the memory size is kept as fixed during the entire training process. This makes the algorithm relatively more robust than the transfer learning technique.

One of the most intuitive boundary-free memory-based CL methods is the random reservoir sampling algorithm [33], which fills \mathcal{M}_t with data uniformly randomly sampled from each of the past episodes. Other mechanisms include the works that aim to increase the sampling diversity, and the diversity is either measured by stochastic gradient directions with respect to the sample loss functions [32] or samples' Euclidean distances [31,53]. However, these works have a number of drawbacks. First, for the reservoir sampling approach, if certain episode only contains a very small number of samples, then samples from this episode will be poorly represented in \mathcal{M}_t because the probability of having samples from an episode in the memory is only dependent on the size of the episode. Second, for the diversity based methods, the approach is again heuristic, since it is not clear how the "diversity" measured by large gradient or Euclidean distances can be directly linked to the quality of representation of the dataset. Third, and perhaps most importantly, the ways that the memory sets are selected are independent of the actual learning tasks at hand. This last property makes these algorithms broadly applicable to different kinds of learning tasks, but also prevents them from exploring application-specific structures. It is not clear whether, and how well these approaches will work for the wireless communication applications of interest in this paper.

4.2 The Proposed Approach

In this work, we propose a new memory-based CL formulation that is customized to the wireless resource allocation problem. Our approach is a departure from the existing memory-based CL approaches discussed in the previous subsection, because we intend to directly use features of the learning problem at hand to build our memory selection mechanism.

To explain the proposed approach, let us begin with presenting two common ways of formulating the training problem for learning optimal wireless resource allocation. First, one can adopt an unsupervised learning approach, which directly optimizes some native measures of wireless system performance, such as the throughput, Quality of Service, or the user fairness [60,61], and this approach does not need any labeled data. Specifically, a popular DNN training problem is given by

$$\min_{\mathbf{\Theta}} \sum_{i \in \mathcal{D}_{0:T}} \ell(\mathbf{\Theta}; \mathbf{h}^{(i)}) \tag{2}$$

where $\mathbf{h}^{(i)}$ is the ith sample; $\boldsymbol{\Theta}$ is the DNN weight to be optimized; $\ell(\cdot)$ is the negative of the per-sample sum-rate function, that is: $\ell(\mathbf{\Theta}; \mathbf{h}^{(i)}) = -R(\pi(\mathbf{\Theta}; \mathbf{h}^{(i)}); \mathbf{h}^{(i)})$, where R is defined in (1) and $\pi(\mathbf{\Theta}; \mathbf{h}^{(i)})$ is the output of DNN which predicts the power allocation. The advantage of this class of unsupervised learning approach is that the system performance measure is directly built into the learning model, while the downside is that this approach can get stuck at low-quality local solutions due to the non-convex nature of DNN [62].

Secondly, it is also possible to use a supervised learning approach. Towards this end, we can generate some labeled data by executing a state-of-the-art optimization algorithm over all the training data samples [11]. Specifically, for each CSI vector $\mathbf{h}^{(i)}$, we can use algorithms such as the WMMSE [34] to solve problem (1) and obtain a high-quality solution $\mathbf{p}^{(i)}$. Putting the $\mathbf{h}^{(i)}$ and $\mathbf{p}^{(i)}$ together yields the *i*th labeled data sample. Such a supervised learning approach typically finds high-quality models [11,62], but often incurs significant computation cost since generating labels can be very time-consuming. Additionally, the quality of the learning model is usually limited by that of label-generating optimization algorithms.

Our idea is to leverage the advantages of both training approaches to construct a memory-based CL formulation. Specifically, we propose to select the most representative data samples $\mathbf{h}^{(i)}$'s into the working memory by using a sample fairness criteria, that is, those data samples that have relatively low system performance are more likely to be selected into the memory. Meanwhile, the DNN is trained by performing supervised learning over the selected data samples. Intuitively, the subset of under-performing data samples is selected to represent a given episode, and we believe that as long as a learning model can perform well on those under-performing samples, then it should perform well for the rest of the samples in a given episode.

To formulate such a notion of sample fairness, let us first assume that the entire dataset $\mathcal{D}_{0:T}$ is available. Let us use $\ell(\cdot)$ to denote a function measuring the per-sample training loss, $g(\cdot)$ a loss function measuring system performance for one data sample, Θ the weights to be trained, $\mathbf{x}^{(i)}$ the ith data sample and $\mathbf{y}^{(i)}$ the ith label. Let $\pi(\Theta; \mathbf{x}^{(i)})$ denote the output of the neural network. Let us consider the following bi-level optimization problem

$$\min_{\boldsymbol{\Theta}} \quad \sum_{i \in \mathcal{D}_{0:T}} \lambda_*^{(i)}(\boldsymbol{\Theta}) \cdot \ell(\boldsymbol{\Theta}; \mathbf{x}^{(i)}, \mathbf{y}^{(i)}) \tag{3a}$$
s.t.
$$\lambda_*(\boldsymbol{\Theta}) = \arg\max_{\boldsymbol{\lambda} \in \mathcal{B}} \sum_{i \in \mathcal{D}_{0:T}} \lambda^{(i)} \cdot g(\boldsymbol{\Theta}; \mathbf{x}^{(i)}, \mathbf{y}^{(i)}), \tag{3b}$$

s.t.
$$\lambda_*(\mathbf{\Theta}) = \arg \max_{\lambda \in \mathcal{B}} \sum_{i \in \mathcal{D}_{0,T}} \lambda^{(i)} \cdot g(\mathbf{\Theta}; \mathbf{x}^{(i)}, \mathbf{y}^{(i)}),$$
 (3b)

where \mathcal{B} denotes the simplex constraint

$$\mathcal{B} := \left\{ oldsymbol{\lambda} \mid \sum_{i \in \mathcal{D}_{0:T}} \lambda^{(i)} = 1, \quad \lambda^{(i)} \geq 0, \quad \forall \ i \in \mathcal{D}_{0:T}
ight\}.$$

In the above formulation, the upper level problem (3a) optimizes the weighted training performance across all data samples, and the lower level problem (3b) assigns larger weights to those data samples that have higher loss $g(\cdot)$ (or equivalently, lower system level performance). Clearly, the lower level problem has a linear objective value, so the optimal λ_* is always on the boundary of the simplex, and the non-zero elements in λ_* will all have the same weight ¹. Such a solution naturally selects a subset of data for the upper-level training problem to optimize. A few remarks about the above formulation is in order.

Remark 1. (Choices of Loss Functions) One important feature of the above formulation is that we decompose the training problem and the data selection problem, so that we can have the flexibility of choosing different loss functions according to the applications at hand. Below we discuss a few suggested choices of the training loss $\ell(\cdot)$ and the system performance loss $g(\cdot)$.

First, the upper layer problem trains the DNN parameters Θ , so we can adopt any existing training formulation we discussed above. For example, if supervised learning is used, then one common training loss is the MSE loss:

$$\ell_{\text{MSE}}(\boldsymbol{\Theta}; \mathbf{x}^{(i)}, \mathbf{y}^{(i)}) = \|\mathbf{y}^{(i)} - \pi(\boldsymbol{\Theta}, \mathbf{x}^{(i)})\|^2. \tag{4}$$

Second, we suggest that the lower level loss function $g(\cdot)$ be directly related to the system level performance. For example, we can set $g(\cdot)$ to be some adaptive weighted negative sum-rate for ith data sample

$$g_{WSR}(\boldsymbol{\Theta}, \mathbf{x}^{(i)}, \mathbf{y}^{(i)}) = -\alpha_i(\boldsymbol{\Theta}; \mathbf{x}^{(i)}, \mathbf{y}^{(i)}) \cdot R(\pi(\boldsymbol{\Theta}, \mathbf{x}^{(i)}); \mathbf{x}^{(i)}).$$
 (5)

If we choose $\alpha_i(\Theta; \mathbf{x}^{(i)}, \mathbf{y}^{(i)}) \equiv 1, \ \forall i$, then the channel realization that achieves the worse throughput by the current DNN model will always be selected. However, this may not be a good choice since the achievable rates at samples across different episodes can vary significantly (e.g., some episodes can have strong interference). Then it is likely that the selected data become concentrated to a few episodes. As an alternative, we can choose $\alpha_i(\Theta; \mathbf{x}^{(i)}, \mathbf{y}^{(i)}) = \frac{1}{R^{(i)}}$, where $R^{(i)}$ is the rate achievable by running some existing optimization algorithm on the sample $\mathbf{x}^{(i)}$. This way, the data samples that achieve the worst sum-rate "relative" to the state-of-the-art optimization algorithm is more likely to be selected. Typically the ratio $1/R^{(i)} \cdot R(\pi(\Theta, \mathbf{x}^{(i)}); \mathbf{x}^{(i)})$ is quite uniform across data samples [11], making it a good indication that a particular data sample is underperforming or not.

Remark 2. (Special Case) As a special case of problem (3), one can choose $\ell(\cdot)$ to be the same as $g(\cdot)$. Then the bi-level problem reduces to the following min-max problem, which optimizes the worst case performance (measured by the loss $\ell(\cdot)$)) across all samples:

$$\min_{\boldsymbol{\Theta}} \max_{\boldsymbol{\lambda}} \quad \sum_{i \in \mathcal{D}_{0:T}} \lambda^{(i)} \cdot \ell(\boldsymbol{\Theta}; \mathbf{x}^{(i)}, \mathbf{y}^{(i)})
\text{s.t.} \quad \sum_{i \in \mathcal{D}_{0:T}} \lambda^{(i)} = 1, \quad \lambda^{(i)} \ge 0 \quad \forall \ i \in \mathcal{D}_{0:T}.$$
(6)

¹In fact, this problem becomes closely related to the classical min-max fairness criteria; see [63] for detailed discussion

Algorithm 1: Fairness Based Continual Learning Framework

```
1 Input: Memory \mathcal{M}_0 = \emptyset, memory size M, max iterations R, step-sizes \alpha, \beta
 2 while receive \mathcal{D}_t do
               \mathcal{G}_t = \mathcal{M}_t \cup \mathcal{D}_t
               for r = 1 : R do
 4

\frac{\mathbf{\Theta} \leftarrow \mathbf{\Theta} - \alpha \sum_{i \in \mathcal{G}_t} \lambda_t^{(i)} \cdot \nabla \ell(\mathbf{\Theta}, \mathbf{x}^{(i)}, \mathbf{y}^{(i)})}{\boldsymbol{\lambda}_t \leftarrow \operatorname{proj} (\boldsymbol{\lambda}_t + \beta \cdot G(\mathbf{\Theta}; \mathbf{x}, \mathbf{y}))}

                                                                                                                                                                        // G(\boldsymbol{\Theta}; \mathbf{x}, \mathbf{y}) := \{q(\boldsymbol{\Theta}, \mathbf{x}^{(i)}, \mathbf{y}^{(i)})\}_{i \in \mathcal{G}_t}
               if |\mathcal{G}_t| < M then
 7
                \mathcal{M}_{t+1} = \mathcal{G}_t.
 8
 9
                      index = sort(\lambda_t)[1:M]
10
                                                                                                                                                                               // Pick top M largest weights
                     \mathcal{M}_{t+1} = \mathcal{G}_t[\text{index}]
                                                                                                                                                                               // Assign M samples to memory
11
```

When $\ell(\cdot)$ is taken as the negative per-sample sum-rate defined in (2), problem (6) is related to the classical min-max resource allocation [60,64,65], with the key difference that it does not achieve fairness across users, but rather to achieve fairness across data samples.

At this point, neither the bi-level problem (3) nor the min-max formulation (6) can be used to design CL strategy yet, because solving these problems requires the full data $\mathcal{D}_{0:T}$. To make these formulations useful for the considered CL setting, we make the following approximation. Suppose that at t-th time instance, we have the memory \mathcal{M}_t and the new data set \mathcal{D}_t available. Then, we propose to solve the following problem to identify worthy data candidates at time t:

$$\min_{\boldsymbol{\Theta}} \quad \sum_{i \in \mathcal{M}_t \cup \mathcal{D}_t} \lambda_t^{(i)}(\boldsymbol{\Theta}) \cdot \ell(\boldsymbol{\Theta}; \mathbf{x}^{(i)}, \mathbf{y}^{(i)})
\text{s.t.} \quad \boldsymbol{\lambda}_t(\boldsymbol{\Theta}) = \arg\max_{\boldsymbol{\lambda} \in \mathcal{B}_t} \sum_{i \in \mathcal{M}_t \cup \mathcal{D}_t} \lambda^{(i)} \cdot g(\boldsymbol{\Theta}; \mathbf{x}^{(i)}, \mathbf{y}^{(i)}),$$
(7)

where \mathcal{B}_t denotes the simplex constraint

$$\mathcal{B}_t := \left\{ \boldsymbol{\lambda} \mid \sum_{i \in \mathcal{M}_t \cup \mathcal{D}_t} \lambda^{(i)} = 1, \quad \lambda^{(i)} \ge 0, \ \forall \ i \in \mathcal{M}_t \cup \mathcal{D}_t \right\}.$$

Similarly, if $\ell(\cdot)$ and $g(\cdot)$ are the same, then the following min-max problem will be solved:

$$\min_{\mathbf{\Theta}} \max_{\mathbf{\lambda}} \quad \sum_{i \in \mathcal{M}_t \cup \mathcal{D}_t} \lambda^{(i)} \cdot \ell(\mathbf{\Theta}, \mathbf{x}^{(i)}, \mathbf{y}^{(i)})$$
subject to
$$\sum_{i \in \mathcal{M}_t \cup \mathcal{D}_t} \lambda^{(i)} = 1, \lambda^{(i)} \ge 0, \ \forall \ i \in \mathcal{M}_t \cup \mathcal{D}_t.$$
(8)

More specifically, at a given time t, the above problem will be approximately solved, and we will collect M data points $j \in \mathcal{M}_t \cup \mathcal{D}_t$ whose corresponding $\lambda^{(j)}$ are the largest. These data points will form the next memory \mathcal{M}_{t+1} , and problem (7) (or (8)) will be solved again; See Algorithm 1 for details. Note that in

the table, we have defined λ and $G(\Theta, \mathbf{x}, \mathbf{y})$ as concatenations of all local weights and loss functions in $\mathcal{G}_t := \mathcal{M}_t \cup \mathcal{D}_t$, i.e., $\lambda := \{\lambda^{(i)}\}_{i \in \mathcal{M}_t \cup \mathcal{D}_t}$ and $G(\Theta; \mathbf{x}, \mathbf{y}) := \{g(\Theta, \mathbf{x}^{(i)}, \mathbf{y}^{(i)})\}_{i \in \mathcal{M}_t \cup \mathcal{D}_t}$, respectively.

Note that line 4-6 in Algorithm 1 is motivated by a class of alternating projected gradient descent ascent (PGDA) algorithm recently proposed for solving non-convex/concave min-max problem (8), see [66] for a recent survey about related algorithms. Such an iteration converges to some stationary solutions for the min-max formulation (8) (see the claim below), while it is not clear if it works for the bi-level formulation (7). Nevertheless, we found that in practice Algorithm 1 works well for both formulations.

Claim 1. [67, Theorem 4.2] Consider the min-max problem (8) with loss function $\ell(\cdot)$ that continuously differentiable in both Θ and λ , and there exists constants L_1 and L_2 such that for every Θ_1 , Θ_2 and data points \mathbf{x}, \mathbf{y} , we have

$$\|\ell(\mathbf{\Theta}_{1}, \mathbf{x}, \mathbf{y}) - \ell(\mathbf{\Theta}_{2}, \mathbf{x}, \mathbf{y})\| \le L_{1} \|\mathbf{\Theta}_{1} - \mathbf{\Theta}_{2}\|,$$
$$\|\nabla_{\mathbf{\Theta}} \ell(\mathbf{\Theta}_{1}, \mathbf{x}, \mathbf{y}) - \nabla_{\mathbf{\Theta}} \ell(\mathbf{\Theta}_{2}, \mathbf{x}, \mathbf{y})\| \le L_{2} \|\mathbf{\Theta}_{1} - \mathbf{\Theta}_{2}\|.$$

For each time instance t, if we apply the algorithm with the following iterative updates

$$\boldsymbol{\Theta} \leftarrow \boldsymbol{\Theta} - \alpha \sum_{i \in \mathcal{M}_t \cup \mathcal{D}_t} \lambda_t^{(i)} \cdot \nabla \ell(\boldsymbol{\Theta}; \mathbf{x}^{(i)}, \mathbf{y}^{(i)}),$$
$$\boldsymbol{\lambda}_t \leftarrow \operatorname{proj} (\boldsymbol{\lambda}_t + \beta \cdot G(\boldsymbol{\Theta}; \mathbf{x}, \mathbf{y})),$$

then there exists some appropriately chosen step-sizes α, β such that the algorithm converges to a stationary point of the problem (8).

5 Experimental Results

In this section, we illustrate the performance of the proposed CL framework using two different applications:

1) power control for weighted sum-rate (WSR) maximization problem with single-input single-output (SISO) interference channel defined in (1), where the end-to-end learning based DNN is used; see [11] for details about the architecture; 2) beamforming for WSR maximization problem with multi-input single-output (MISO) interference channel, where deep unfolding based DNN architecture is used; see [17] for details about the architecture. By demonstrating the effectiveness of the proposed methods on these two approaches, we believe that our approach can be extended to many other related problems as well.

5.1 Simulation Setup

The experiments are conducted on Ubuntu 18.04 with Python 3.6, PyTorch 1.6.0, and MATLAB R2019b on one computer node with two 8-core Intel Haswell processors and 128 GB of memory. The codes are made available online through https://github.com/Haoran-S/ICASSP2021.

5.2 Randomly Generated Channel

We first demonstrate the performance of our proposed framework using randomly generated channels, for a scenario with K=10 transmitter-receiver pairs. We choose three standard types of random channels used in previous resource allocation literature [11,13] stated as following:

Rayleigh fading: Each channel coefficient h_{ij} is generated according to a standard normal distribution, i.e.,

$$\operatorname{Re}(h_{ij}) \sim \frac{\mathcal{N}(0,1)}{\sqrt{2}}, \quad \operatorname{Im}(h_{ij}) \sim \frac{\mathcal{N}(0,1)}{\sqrt{2}}, \quad \forall i,j \in \mathcal{K}.$$
 (9)

Rician fading: Each channel coefficient h_{ij} is generated according to a Gaussian distribution with 0dB K-factor, i.e.,

$$\operatorname{Re}(h_{ij}) \sim \frac{1 + \mathcal{N}(0, 1)}{2}, \quad \operatorname{Im}(h_{ij}) \sim \frac{1 + \mathcal{N}(0, 1)}{2}, \quad \forall i, j \in \mathcal{K}.$$

Geometry channel: All transmitters and receivers are uniformly randomly distributed in a R×R area. The channel gains $|h_{ij}|^2$, which follow the pathloss function in [68], can be written as

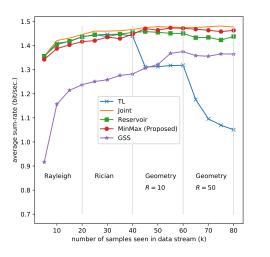
$$|h_{ij}|^2 = \frac{1}{1 + d_{ij}^2} |f_{ij}|^2, \ \forall \ i, j$$

where f_{ij} is the small-scale fading coefficient follows $\mathcal{CN}(0,1)$, d_{ij} is the distance between the *i*th transmitter and *j*th receiver.

Then, we use these coefficients to generate four different episodes as marked in Fig. 3, namely the Rayleigh fading channel, the Ricean fading channel, and the geometry channel (with nodes distributed in a $10m \times 10m$ and a $50m \times 50m$ area). Note that this way of generating data and organizing the episodes may not make sense in practice. It is only used at this point as a "toy" data set. Later we will utilize real data to generate more practical scenarios. For each episode, we generate 20,000 channel realizations for training and 1,000 for testing. We also stacked the test data from all episodes to form a mixture test set, i.e., containing 4,000 channel realizations. During the training stage, a total of 80,000 channel realizations are available. A batch of 5,000 realizations is revealed each time, and the memory space contains only 2,000 samples from the past. That is, $|\mathcal{D}_{1:16}| = 80,000$, $|\mathcal{D}_t| = 5,000$, $|\mathcal{M}_t| = 2,000$, $\forall t$.

For the data-driven model, we use the end-to-end learning based fully connected neural network model as implemented in [11]. For each data batch of 5,000 realizations at time t, we train the model Θ_t using following five different approaches for 100 epochs, using previous model Θ_{t-1} as initialization

- 1. Transfer learning (TL) [15] update the model using the current data batch (a total of 5,000 samples);
- 2. Joint training (Joint) update the model using all accumulated data up to current time (up to 80,000 samples);
- 3. Reservoir sampling based CL (Reservoir) [31] update the model using both the current data batch and the memory set (a total of 7,000 samples), where data samples in the memory set are uniformly randomly sampled from the streaming episodes;
- 4. Gradient sample selection based CL (GSS) [32] update the model using both the current data batch and the memory set (a total of 7,000 samples), where data samples in the memory set are selected by maximizing the samples' stochastic gradient directions diversity;
- 5. Proposed fairness based CL (MinMax) in Algorithm 1 update the model using both the current data batch and the memory set (a total of 7,000 samples), where data samples in the memory set are selected according to the proposed data-sample fairness criterion (7). Unless otherwise specified, as



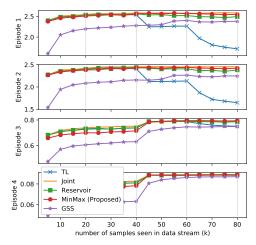


Figure 3: Testing performance on randomly generated channels. (left) Average sum-rate performance comparison on the mixture test set of all four episodes. (right) Average sum-rate performance comparison for each individual episode. The gray line indicates the change of episode, which is unknown during training time.

suggested in Section 4.2, the training loss $\ell(\cdot)$ is chosen as the MSE loss (4), the system performance loss $g(\cdot)$ is chosen as the adaptive weighted negative sum-rate loss (5), and the weights is chosen as the sum-rate achievable by the WMMSE method [34].

The simulation results of five different approaches are compared and shown in Fig. 3 - 5. In specific, Fig. 3 shows the trends of the average sum-rate performance, Fig. 4 shows how the above sum-rate compared with those rates obtained by the WMMSE algorithm [34], i.e., approximation ratios, and Fig. 5 shows the distribution of the approximation ratios for the mixture test set at the last timestamp. Note that the joint training method uses up to 80,000 in memory spaces and thus violates our memory limitation (i.e., 7,000 in total), and the transfer learning method adapts the model to new data each time and did not use any additional memory spaces. We can observe that the proposed fairness based CL method performs well over all tasks (cf. Fig. 3 - 5), nearly matching the performance of the joint training method, whereas the transfer learning suffers from some significant performance loss as the "outlier" episode comes in (i.e., geometry channel in our case). Furthermore, we can observe that the proposed fairness based CL method outperforms other CL-based methods (Reservoir and GSS), in terms of both the average sum-rate (cf. Fig. 3) and sample fairness (cf. Fig. 5) (i.e., fewer samples are distributed in the low sum-rate region) over all tasks. This validates that the proposed approach indeed incorporates the problem structure and advocates fairness across the data samples.

We also demonstrate the drawbacks of the reservoir sampling method using a two-episode unbalanced dataset (each with 2,000 and 18,000 samples, respectively) with a small memory size (100 samples). We consider two episodes where episode 1 has 2,000 samples, and it follows the original Rayleigh fading described in (9); episode 2 has 18,000 samples, and it follows the Rayleigh fading but with strong direct channels (the diagonal entries has magnitude 5 times larger). Simulation results in Fig. 6 on both the sum-rate and the sum-rate ratio (compared with WMMSE) show that the reservoir sampling performs worse in unbalancing scenarios, since the episodes with fewer samples will be poorly represented in the memory set (i.e., memory set contains around $\frac{2,000}{2,000+18,000} \times 100 = 10$ samples from episode 1, out of

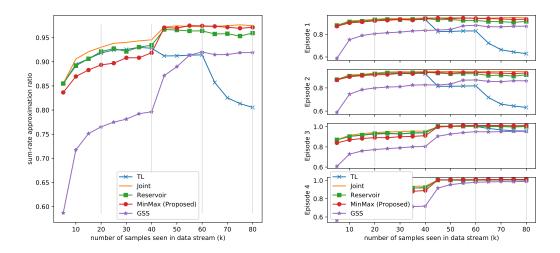


Figure 4: Testing performance on randomly generated channels. (left) Average sum-rate approximation ratios comparison on the mixture test set of all four episodes. (right) Average sum-rate approximation ratios comparison for each individual episode. The ratio is computed by the achievable sum-rate of the indicated method divided by the sum-rate achieved by the WMMSE algorithm. The gray line indicates the change of episode, which is unknown during training time.

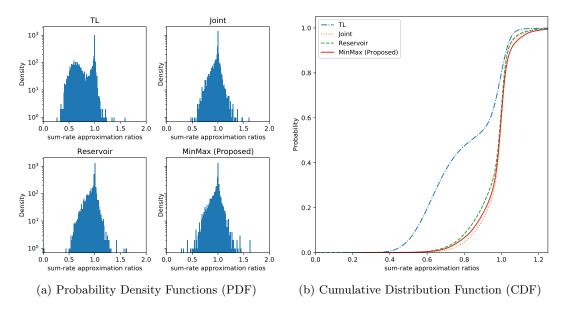


Figure 5: Distribution of the per-sample sum-rate approximation ratio evaluated at the last time stamp on the mixture test set of all four episodes. The ratio is computed by the achievable sum-rate of the indicated method divided by the sum-rate achieved by the WMMSE algorithm.

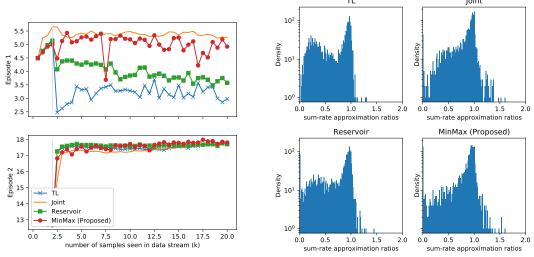


Figure 6: Example on unbalanced data (2,000 vs 18,000) and small memory size (0.1k) to showcase of drawbacks of the naive reservoir sampling. Simulation results on (left) sum-rate and (right) sum-rate ratio compared with WMMSE, show that the reservoir sampling performs similar to transfer learning while the proposed fairness based methods almost matches the best performance achieved by joint training with accumulated data.

a 100 memory size). On the other hand, the proposed fairness based method almost matches the best performance achieved by joint training with accumulated data.

5.3 Real Measured Channel

To validate our approach under more realistic scenarios, we further consider the outdoor 'O1' ray-tracing scenario generated from the DeepMIMO dataset [69]. The used dataset consists of two streets and one intersection, with the top-view showed in Fig. 7. The horizontal street is 600m long and 40m wide, and the vertical street is 440m long and 40m wide. We consider six different episodes that users gathering locations (i.e., red rectangles in Fig. 7) gradually shift following the order indicated by the numbers. For each episode, we generate 20,000 channel realizations for training and 1,000 for testing. For each channel realization, we generate the channel based on the same set of BSs 1 to 10 (i.e., blue circles in Fig. 7), and randomly pick K=10 user locations from the indicated user area i (i.e., red rectangles in Fig. 7). The BS is equipped with single antenna and has the maximum transmit power $p_k=30$ dBm. The noise power is set to -80 dBm. The simulation results are shown in Fig. 8 - 10. From Fig. 8-9 we can observed that, after experiencing all six episodes, our proposed fairness based method is able to perform much better than transfer learning and reservoir sampling-based CL approach on the mixture test set. This is because both transfer learning and reservoir sampling-based CL approach suffer from forgetting issues for old episodes (i.e. episode 1 and 2), but the proposed fairness based method can perform well on episode 1 and 2 due to its fairness design – focusing on under-performing episodes (i.e. episode 1 and 2) while relaxing on outperforming episodes (i.e. episode 5 and 6). One interesting observation is that the proposed method is able to outperform the joint training with accumulated data in terms of the average sum-rate (cf. LHS of Fig. 8). This is because the joint training will treat all samples equally, thus if we have two episodes in opposite directions, the model will generate something in the middle. Instead, our proposed fairness based method will be more focus on one episode that generates the highest cost thus achieve higher average performance due to its fairness design.

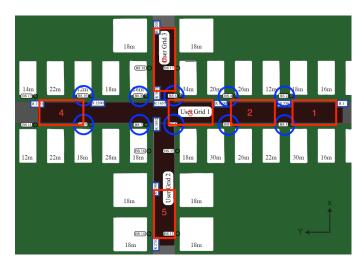


Figure 7: Top view of the DeepMIMO ray-tracing scenario [69], showing the two streets, the buildings, the 10 base stations (blue circle), and the user x-y grids (red rectangular).

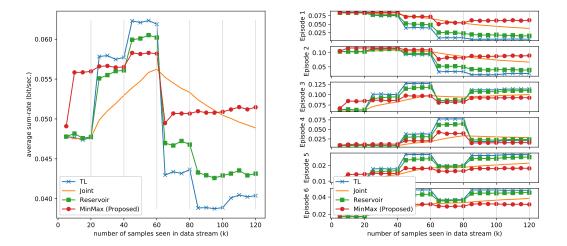


Figure 8: Average sum-rate performance comparison on real measured channels. (left) performance on mixture of all six episodes. (right) per-task sum-rate performance.

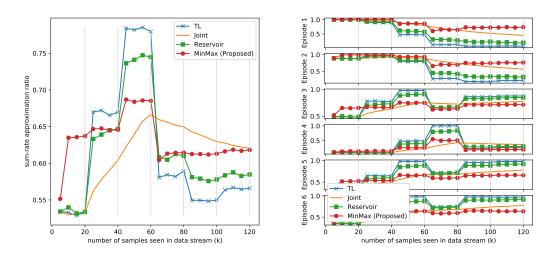


Figure 9: Average sum-rate approximation ratios comparison on real measured channels. (left) performance on mixture of all six episodes. (right) per-task performance.

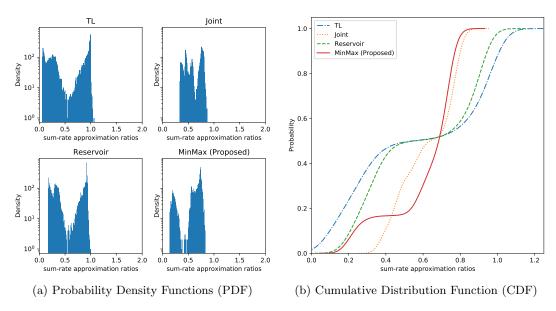


Figure 10: Distribution of the per-sample sum-rate ratio evaluated at the last time stamp on the mixture test set of all six episodes.

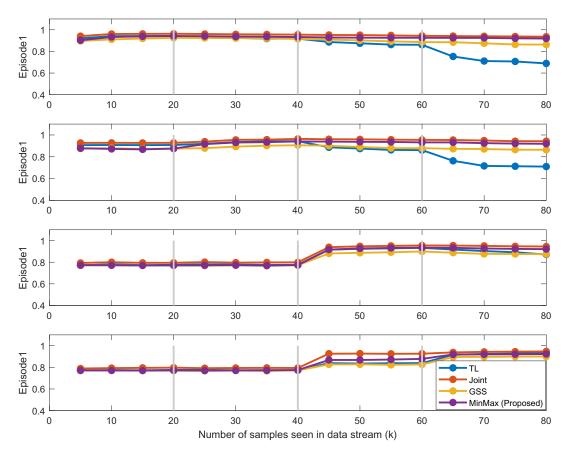


Figure 11: Achievable accuracy comparison for beamforming problem over the test set under randomly generated channels, as the model learned progressively over time.

5.4 Beamforming Experiments

Next, we further validate our CL based approach in the beamforming problems – weighted sum-rate (WSR) maximization in multi-input and single-output (MISO) interference channels. Different from the previous sections where the deep learning model is based on the end-to-end learning [11], we consider a deep unfolding based model proposed in [17] and [70]. Specifically, by unfolding the parallel gradient projection (PGP) method, the parallel GP based recurrent neural network (RNN-PGP) method [70] builds a recurrent neural network (RNN) architecture to learn the beamforming solutions efficiently. By leveraging the low-dimensional structures of the optimal beamforming solution, the constructed neural network can achieve high WSRs with significantly reduced runtime, while exhibiting favorable generalization capability with respect to the antenna number, BS number and the inter-BS distance.

In our simulation, we consider the same four types of channels described in Section 5.2, but in the form of cellular network described in [70, Section 2] with the number of base stations (BS) K = 19, number of neighbor considered c = 18 and transmit antennas $N_t = 36$. For the two episodes with geometry channels, we set the cell radius d as 0.5 km and 1 km, respectively. The simulation results over four different approaches are compared and reported in Fig. 11, where the x-axis represents the number of data samples has been observed, and the y-axis denotes the ratio of the (weighted) sum-rate achieved by the RNN-PGP [70] and that achieved by the WMMSE algorithm [34]. For our proposed approach (MinMax),

we use the minmax formulation (8) where both the training loss $\ell(\cdot)$ and the system performance loss $g(\cdot)$ are chosen as the MSE loss (4). It can be observed that the proposed algorithm (MinMax) almost matches the joint training performances, with only limited memory usage, and it outperforms the GSS approach. On the other hand, the naive transfer learning can suffer from significant performance loss on old episodes (Episode 1 and 2 in our case).

6 Conclusion and Future Works

In this work, we design a new "learning to continuously optimize" framework for optimizing wireless resources in dynamic environments, where parameters such as CSIs keep changing. By building the notion of continual learning (CL) into the modeling process, our framework is able to seamlessly and efficiently adapt to the episodically dynamic environment, without knowing the episode boundary, and most importantly, maintain high performance over all the previously encountered scenarios. The proposed approach is validated through two popular wireless resource allocation problems (one for power control and one for beamforming), two popular DNN based models (one for end-to-end learning and one for deep unfolding), and uses both synthetic and real data sets. Our empirical results make us believe that our approaches can be extended to many other related problems.

Our work represents a preliminary step towards understanding the capability of deep learning for wireless problems with dynamic environments. There are many interesting questions to be addressed in the future, and some of them are listed below:

- Is it possible to design continual learning strategy with theoretical guarantees?
- Is it possible to design theoretical results for the proposed stage-wise bi-level optimization problem (7) or even the global bi-level optimization problem (3)?
- Is it possible to quantify the generalization performance of the proposed fairness framework?
- Is it possible to extend our frameworks to other wireless tasks such as signal detection, channel estimation and CSI compression?

References

- [1] H. Sun, W. Pu, M. Zhu, X. Fu, T.-H. Chang, and M. Hong, "Learning to continuously optimize wireless resource in episodically dynamic environment," in submitted to ICASSP, 2021.
- [2] A. Voulodimos, N. Doulamis, A. Doulamis, and E. Protopapadakis, "Deep learning for computer vision: A brief review," Computational Intelligence and Neuroscience, 2018.
- [3] T. Young, D. Hazarika, S. Poria, and E. Cambria, "Recent trends in deep learning based natural language processing," IEEE Computational Intelligence Magazine, vol. 13, no. 3, pp. 55–75, 2018.
- [4] H. Wang, N. Wang, and D.-Y. Yeung, "Collaborative deep learning for recommender systems," in Proceedings of the 21th ACM International Conference on Knowledge Discovery and Data Mining (SIGKDD), 2015, pp. 1235–1244.
- [5] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," Nature, vol. 521, no. 7553, pp. 436–444, 2015.
- [6] H. Ye, G. Y. Li, and B.-H. Juang, "Power of deep learning for channel estimation and signal detection in OFDM systems," IEEE Wireless Communications Letters, vol. 7, no. 1, pp. 114–117, 2017.
- [7] H. Sun, A. O. Kaya, M. Macdonald, H. Viswanathan, and M. Hong, "Deep learning based preamble detection and TOA estimation," in proceedings of the 2019 IEEE Global Communications Conference (GLOBECOM).
- [8] E. Nachmani, E. Marciano, L. Lugosch, W. J. Gross, D. Burshtein, and Y. Be'ery, "Deep learning methods for improved decoding of linear codes," <u>IEEE Journal of Selected Topics in Signal Processing</u>, vol. 12, no. 1, pp. 119–131, 2018.
- [9] C.-K. Wen, W.-T. Shih, and S. Jin, "Deep learning for massive MIMO CSI feedback," <u>IEEE Wireless</u> Communications Letters, vol. 7, no. 5, pp. 748–751, 2018.
- [10] H. Sun, Z. Zhao, X. Fu, and M. Hong, "Limited feedback double directional massive MIMO channel estimation: From low-rank modeling to deep learning," in proceedings of the IEEE 19th International Workshop on Signal Processing Advances in Wireless Communications (SPAWC), 2018.
- [11] H. Sun, X. Chen, Q. Shi, M. Hong, X. Fu, and N. D. Sidiropoulos, "Learning to optimize: Training deep neural networks for interference management," <u>IEEE Transactions on Signal Processing</u>, vol. 66, no. 20, pp. 5438–5453, 2018.
- [12] W. Lee, M. Kim, and D.-H. Cho, "Deep power control: Transmit power control scheme based on convolutional neural network," IEEE Communications Letters, vol. 22, no. 6, pp. 1276–1279, 2018.
- [13] F. Liang, C. Shen, W. Yu, and F. Wu, "Towards optimal power control via ensembling deep neural networks," IEEE Transactions on Communications, vol. 68, no. 3, pp. 1760–1776, 2019.
- [14] M. Eisen and A. R. Ribeiro, "Optimal wireless resource allocation with random edge graph neural networks," IEEE Transactions on Signal Processing, 2020.
- [15] Y. Shen, Y. Shi, J. Zhang, and K. B. Letaief, "LORM: Learning to optimize for resource management in wireless networks with few training samples," <u>IEEE Transactions on Wireless Communications</u>, vol. 19, no. 1, pp. 665–679, 2019.

- [16] H. Huang, Y. Peng, J. Yang, W. Xia, and G. Gui, "Fast beamforming design via deep learning," <u>IEEE</u> Transactions on Vehicular Technology, vol. 69, no. 1, pp. 1065–1069, 2019.
- [17] M. Zhu and T.-H. Chang, "Optimization inspired learning network for multiuser robust beamforming," in proceedings of the IEEE 11th Sensor Array and Multichannel Signal Processing Workshop (SAM), 2020.
- [18] W. Cui, K. Shen, and W. Yu, "Spatial deep learning for wireless scheduling," <u>IEEE Journal on Selected</u> Areas in Communications, vol. 37, no. 6, pp. 1248–1261, 2019.
- [19] S. Dörner, S. Cammerer, J. Hoydis, and S. Ten Brink, "Deep learning based communication over the air," IEEE Journal of Selected Topics in Signal Processing, vol. 12, no. 1, pp. 132–143, 2017.
- [20] T. O'Shea and J. Hoydis, "An introduction to deep learning for the physical layer," <u>IEEE Transactions</u> on Cognitive Communications and Networking, vol. 3, no. 4, pp. 563–575, 2017.
- [21] G. I. Parisi, R. Kemker, J. L. Part, C. Kanan, and S. Wermter, "Continual lifelong learning with neural networks: A review," Neural Networks, vol. 113, pp. 54–71, 2019.
- [22] M. McCloskey and N. J. Cohen, "Catastrophic interference in connectionist networks: The sequential learning problem," in Psychology of Learning and Motivation. Elsevier, 1989, vol. 24, pp. 109–165.
- [23] J. Kirkpatrick, R. Pascanu, N. Rabinowitz, J. Veness, G. Desjardins, A. A. Rusu, K. Milan, J. Quan, T. Ramalho, A. Grabska-Barwinska et al., "Overcoming catastrophic forgetting in neural networks," Proceedings of the National Academy of Sciences, vol. 114, no. 13, pp. 3521–3526, 2017.
- [24] Z. Li and D. Hoiem, "Learning without forgetting," <u>IEEE Transactions on Pattern Analysis and Machine Intelligence</u>, vol. 40, no. 12, pp. 2935–2947, 2017.
- [25] F. Zenke, B. Poole, and S. Ganguli, "Continual learning through synaptic intelligence," <u>Proceedings</u> of Machine Learning Research, vol. 70, p. 3987, 2017.
- [26] A. A. Rusu, N. C. Rabinowitz, G. Desjardins, H. Soyer, J. Kirkpatrick, K. Kavukcuoglu, R. Pascanu, and R. Hadsell, "Progressive neural networks," arXiv:1606.04671, 2016.
- [27] J. Yoon, E. Yang, J. Lee, and S. J. Hwang, "Lifelong learning with dynamically expandable networks," in <u>International Conference on Learning Representations</u>, 2018.
- [28] D. Lopez-Paz and M. Ranzato, "Gradient episodic memory for continual learning," in <u>Advances in Neural Information Processing Systems</u>, 2017, pp. 6467–6476.
- [29] H. Shin, J. K. Lee, J. Kim, and J. Kim, "Continual learning with deep generative replay," in <u>Advances</u> in Neural Information Processing Systems, 2017, pp. 2990–2999.
- [30] S.-A. Rebuffi, A. Kolesnikov, G. Sperl, and C. H. Lampert, "icarl: Incremental classifier and representation learning," in <u>Proceedings of the IEEE conference on Computer Vision and Pattern Recognition</u>, 2017, pp. 2001–2010.
- [31] D. Isele and A. Cosgun, "Selective experience replay for lifelong learning," in <u>Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)</u>, 2018, pp. 3302–3309.

- [32] R. Aljundi, M. Lin, B. Goujaud, and Y. Bengio, "Gradient based sample selection for online continual learning," in Advances in Neural Information Processing Systems, 2019, pp. 11816–11825.
- [33] D. Rolnick, A. Ahuja, J. Schwarz, T. Lillicrap, and G. Wayne, "Experience replay for continual learning," in Advances in Neural Information Processing Systems, 2019, pp. 350–360.
- [34] Q. Shi, M. Razaviyayn, Z.-Q. Luo, and C. He, "An iteratively weighted MMSE approach to distributed sum-utility maximization for a MIMO interfering broadcast channel," <u>IEEE Transactions on Signal Processing</u>, vol. 59, no. 9, pp. 4331–4340, 2011.
- [35] Y. Shen, Y. Shi, J. Zhang, and K. B. Letaief, "A graph neural network approach for scalable wireless power control," in proceedings of the 2019 IEEE Globecom Workshops (GC Wkshps).
- [36] Y. S. Nasir and D. Guo, "Multi-agent deep reinforcement learning for dynamic power allocation in wireless networks," <u>IEEE Journal on Selected Areas in Communications</u>, vol. 37, no. 10, pp. 2239–2250, 2019.
- [37] A. Balatsoukas-Stimming and C. Studer, "Deep unfolding for communications systems: A survey and some new directions," in 2019 IEEE International Workshop on Signal Processing Systems (SiPS). IEEE, 2019, pp. 266–271.
- [38] K. Gregor and Y. LeCun, "Learning fast approximations of sparse coding," in proceedings of the 27th International Conference on Machine Learning (ICML), 2010, pp. 399–406.
- [39] J. R. Hershey, J. L. Roux, and F. Weninger, "Deep unfolding: Model-based inspiration of novel deep architectures," arXiv preprint arXiv:1409.2574, 2014.
- [40] P. Sprechmann, R. Litman, T. B. Yakar, A. M. Bronstein, and G. Sapiro, "Supervised sparse analysis and synthesis operators," in Advances in Neural Information Processing Systems, 2013, pp. 908–916.
- [41] N. Samuel, T. Diskin, and A. Wiesel, "Deep mimo detection," arXiv preprint arXiv:1706.01151, 2017.
- [42] —, "Learning to detect," <u>IEEE Transactions on Signal Processing</u>, vol. 67, no. 10, pp. 2554–2564, 2019.
- [43] H. He, C.-K. Wen, S. Jin, and G. Y. Li, "A model-driven deep learning network for mimo detection," in proceedings of the IEEE Global Conference on Signal and Information Processing (GlobalSIP), 2018, pp. 584–588.
- [44] S. Cammerer, T. Gruber, J. Hoydis, and S. Ten Brink, "Scaling deep learning-based decoding of polar codes via partitioning," in proceedings of the IEEE Global Communications Conference (GLOBECOM), 2017, pp. 1–6.
- [45] M. Eisen, C. Zhang, L. F. Chamon, D. D. Lee, and A. Ribeiro, "Learning optimal resource allocations in wireless systems," IEEE Transactions on Signal Processing, vol. 67, no. 10, pp. 2775–2790, 2019.
- [46] M. Borgerding, P. Schniter, and S. Rangan, "Amp-inspired deep networks for sparse linear inverse problems," IEEE Transactions on Signal Processing, vol. 65, no. 16, pp. 4293–4308, 2017.

- [47] Q. Hu, Y. Cai, Q. Shi, K. Xu, G. Yu, and Z. Ding, "Iterative algorithm induced deep-unfolding neural networks: Precoding design for multiuser mimo systems," <u>IEEE Transactions on Wireless Communications</u>, 2020.
- [48] M. B. Ring, "Continual learning in reinforcement environments," Ph.D. dissertation, University of Texas at Austin Austin, Texas 78712, 1994.
- [49] S. Farquhar and Y. Gal, "Towards robust evaluations of continual learning," <u>arXiv preprint</u> arXiv:1805.09733, 2018.
- [50] M. De Lange, R. Aljundi, M. Masana, S. Parisot, X. Jia, A. Leonardis, G. Slabaugh, and T. Tuytelaars, "A continual learning survey: Defying forgetting in classification tasks," <u>arXiv preprint</u> arXiv:1909.08383, 2019.
- [51] A. Robins, "Catastrophic forgetting, rehearsal and pseudorehearsal," <u>Connection Science</u>, vol. 7, no. 2, pp. 123–146, 1995.
- [52] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. Riedmiller, "Playing atari with deep reinforcement learning," arXiv preprint arXiv:1312.5602, 2013.
- [53] T. L. Hayes, N. D. Cahill, and C. Kanan, "Memory efficient experience replay for streaming learning," in 2019 International Conference on Robotics and Automation (ICRA). IEEE, 2019, pp. 9769–9776.
- [54] S. Shalev-Shwartz et al., "Online learning and online convex optimization," 2011.
- [55] R. Aljundi, E. Belilovsky, T. Tuytelaars, L. Charlin, M. Caccia, M. Lin, and L. Page-Caccia, "Online continual learning with maximal interfered retrieval," in <u>Advances in Neural Information Processing Systems</u>, 2019, pp. 11849–11860.
- [56] S. J. Pan and Q. Yang, "A survey on transfer learning," <u>IEEE Transactions on Knowledge and Data Engineering</u>, vol. 22, no. 10, pp. 1345–1359, 2009.
- [57] Y. Yuan, G. Zheng, K.-K. Wong, B. Ottersten, and Z.-Q. Luo, "Transfer learning and meta learning based fast downlink beamforming adaptation," arXiv preprint arXiv:2011.00903, 2020.
- [58] A. Zappone, M. Di Renzo, M. Debbah, T. T. Lam, and X. Qian, "Model-aided wireless artificial intelligence: Embedding expert knowledge in deep neural networks for wireless system optimization," IEEE Vehicular Technology Magazine, vol. 14, no. 3, pp. 60–69, 2019.
- [59] Z.-Q. Luo and S. Zhang, "Dynamic spectrum management: Complexity and duality," <u>IEEE Journal of Selected Topics in Signal Processing</u>, vol. 2, no. 1, pp. 57–73, 2008.
- [60] J. Mo and J. Walrand, "Fair end-to-end window-based congestion control," <u>IEEE/ACM Transactions</u> on Networking,, vol. 8, no. 5, pp. 556 –567, 2000.
- [61] M. Hong and Z.-Q. Luo, "Signal processing and optimal resource allocation for the interference channel," in Academic Press Library in Signal Processing. Academic Press, 2013.
- [62] B. Song, H. Sun, W. Pu, S. Liu, and M. Hong, "To supervise or not to supervise: How to effectively learn wireless interference management models?" in submitted to ICASSP, 2021.

- [63] S. Lu, I. Tsaknakis, M. Hong, and Y. Chen, "Block alternating optimization for non-convex min-max problems: Algorithms and applications in signal processing and communications," <u>IEEE Transactions</u> on Signal Processing, 2019, accepted.
- [64] M. Bengtsson and B. Ottersten, "Optimal downlink beamforming using semidefinite optimization," in Proceedings of the 37th Annual Allerton Conference, 1999.
- [65] M. Razaviyayn, M. Hong, and Z.-Q. Luo, "Linear transceiver design for a MIMO interfering broadcast channel achieving max-min fairness," Signal Processing, vol. 93, no. 12, pp. 3327–3340, 2013.
- [66] M. Razaviyayn, T. Huang, S. Lu, M. Nouiehed, M. Sanjabi, and M. Hong, "Nonconvex min-max optimization: Applications, challenges, and recent theoretical advances," <u>IEEE Signal Processing Magazine</u>, vol. 37, no. 5, pp. 55–66, 2020.
- [67] M. Nouiehed, M. Sanjabi, T. Huang, J. D. Lee, and M. Razaviyayn, "Solving a class of non-convex min-max games using iterative first order methods," in <u>Advances in Neural Information Processing</u> Systems, 2019, pp. 14934–14942.
- [68] H. Inaltekin, M. Chiang, H. V. Poor, and S. B. Wicker, "On unbounded path-loss models: effects of singularity on wireless network performance," <u>IEEE Journal on Selected Areas in Communications</u>, vol. 27, no. 7, pp. 1078–1092, 2009.
- [69] A. Alkhateeb, "DeepMIMO: A generic deep learning dataset for millimeter wave and massive MIMO applications," in proceedings of the Information Theory and Applications Workshop (ITA), San Diego, CA, 2019.
- [70] M. Zhu, T.-H. Chang, and M. Hong, "Learning to beamform in heterogeneous massive mimo networks," arXiv:2011.03971, 2020.