



<http://www.diva-portal.org>

This is the published version of a paper published in .

Citation for the original published paper (version of record):

Hellström, H., B. da Silva Jr., J M., Amiri, M M., Chen, M., Fodor, V. et al. (2022)

Wireless for Machine Learning: A Survey

Foundations and Trends in Signal Processing, 15(4): 290-399

<https://doi.org/10.1561/20000000114>

Access to the published version may require subscription.

N.B. When citing this work, cite the original published paper.

Permanent link to this version:

<http://urn.kb.se/resolve?urn=urn:nbn:se:kth:diva-313006>

Foundations and Trends® in Signal Processing

Wireless for Machine Learning: A Survey

Suggested Citation: Henrik Hellström, José Mairton Barros da Silva Jr., Mohammad Mohammadi Amiri, Mingzhe Chen, Viktoria Fodor, H. Vincent Poor and Carlo Fischione (2022), “Wireless for Machine Learning: A Survey”, Foundations and Trends® in Signal Processing: Vol. 15, No. 4, pp 290–399. DOI: 10.1561/2000000114.

Henrik Hellström

KTH Royal Institute of Technology
hhells@kth.se

José Mairton B. da Silva Jr.

KTH Royal Institute of Technology

Mohammad Mohammadi Amiri

Massachusetts Institute of Technology

Mingzhe Chen

Princeton University

Viktoria Fodor

KTH Royal Institute of Technology

H. Vincent Poor

Princeton University

Carlo Fischione

KTH Royal Institute of Technology

This article may be used only for the purpose of research, teaching, and/or private study. Commercial use or systematic downloading (by robots or other automatic processes) is prohibited without explicit Publisher approval.

now

the essence of knowledge

Boston — Delft

Contents

1	Introduction	291
1.1	Related Work	293
1.2	Notation and Organization	296
2	Primer on Distributed Machine Learning	300
2.1	Problem Formulation for Centralized Machine Learning . . .	301
2.2	Problem Formulation for Distributed Machine Learning . . .	304
2.3	Federated Learning	305
2.4	Summary	308
3	Analog Over-the-air Computation	309
3.1	Primer	309
3.2	Over-the-air Computation For Distributed Machine Learning	313
3.3	Review of SISO Over-the-air Computation	316
3.4	Review of MIMO Over-the-air Computation	329
4	Digital Communications	335
4.1	Primer	335
4.2	Digital Communications for Distributed Machine Learning .	337
4.3	Review of Importance-aware Communications	340
4.4	Review of Radio Resource Management for Federated Learning	345

5	Open Problems	362
5.1	Over-the-air Computation	362
5.2	Digital Communications	365
5.3	Problems Relevant to Analog and Digital Communications .	367
6	Applications	368
6.1	Smart City	368
6.2	Vehicular Communication	370
6.3	Augmented and Virtual Reality	372
6.4	Edge Caching	373
6.5	Unmanned Aerial Vehicles	375
7	Conclusions	377
	References	379

Wireless for Machine Learning: A Survey

Henrik Hellström¹, José Mairton Barros da Silva Jr.¹,
Mohammad Mohammadi Amiri², Mingzhe Chen³, Viktoria Fodor¹,
H. Vincent Poor³ and Carlo Fischione¹

¹*School of Electrical Engineering and Computer Science, KTH Royal Institute of Technology, Sweden; hhells@kth.se*

²*MIT Media Laboratory, Massachusetts Institute of Technology, USA*

³*Department of Electrical and Computer Engineering, Princeton University, USA*

ABSTRACT

As data generation increasingly takes place on devices without a wired connection, Machine Learning (ML) related traffic will be ubiquitous in wireless networks. Many studies have shown that traditional wireless protocols are highly inefficient or unsustainable to support ML, which creates the need for new wireless communication methods. In this monograph, we give a comprehensive review of the state-of-the-art wireless methods that are specifically designed to support ML services over distributed datasets. Currently, there are two clear themes within the literature, analog over-the-air computation and digital radio resource management optimized for ML. This survey gives an introduction to these methods, reviews the most important works, highlights open problems, and discusses application scenarios.

Henrik Hellström, José Mairton Barros da Silva Jr., Mohammad Mohammadi Amiri, Mingzhe Chen, Viktoria Fodor, H. Vincent Poor and Carlo Fischione (2022), “Wireless for Machine Learning: A Survey”, *Foundations and Trends® in Signal Processing*: Vol. 15, No. 4, pp 290–399. DOI: 10.1561/20000000114.

©2022 H. Hellström *et al.*

1

Introduction

With the increasing popularity of mobile devices and the continuous growth of Internet of Things (IoT), we are having increasing access to vast amounts of distributed data. According to a recent report from Ericsson, the global number of connected IoT devices will rise to 4.1 billion by 2024 [49], which is four times the 1 billion observed in 2019. Simultaneously, breakthroughs in Machine Learning (ML) are allowing us to analyze the data of edge devices so as to solve a wide range of complex problems, such as image recognition [66], language processing [39], and predictive modeling [23]. However, since ML was originally conceived in centralized settings where all data must be aggregated at a common location, the application of ML on distributed datasets over wireless networks is generating new challenges for the wireless networks, namely:

- **Privacy:** Many ML applications require the use of privacy-sensitive data. In these cases, it is either desirable or necessary that the training dataset cannot be inferred by eavesdropping upon the ML updates being transferred wirelessly [150];
- **Security:** When an ML model is trained distributively, a bad actor can corrupt the final model by transmitting malicious model updates [159]. Wireless protocol design should inhibit an attacker's ability to do so;

- **Communication and Energy Efficiency:** Distributed ML (DML) requires the communication of high-dimensional model updates for hundreds or thousands of iterations before the model has converged. This communication of updates generally forms the performance bottleneck of the training process, imposing the risk of excessively draining the batteries of training devices and overwhelming the capacity of the wireless network [144].

To address these challenges, a new approach toward communication protocol design has emerged [198]. This new approach considers the design of new wireless methods for carrying data needed for the ML tasks. Unlike traditional wireless protocol design, the objective of Wireless for ML is not to deliver bits as efficiently as possible, but to distill the intelligence carried within the data. The traditional communication protocols that are designed to maximize data rate and minimize bit errors have been shown to be greatly inefficient for carrying ML related data [9], [35], [100], [118], [200]. Instead, Wireless for ML offers new methods that are better aligned with the ML objective and invites us to rethink how wireless communication protocols are designed. Among the novel methods that have been proposed, two major themes arise, namely analog over-the-air computation (AirComp) and radio resource management (RRM) optimized for ML. In AirComp, the long-standing doctrine of interference avoidance is questioned and novel interference-promoting protocols are proposed while in RRM for ML, the new objectives lead to solutions that are fundamentally different from what is used today.

The idea of wireless protocols customized for ML, although not yet available in the current cellular wireless standards, is compatible with the current standard specifications. The new cellular standard 5G has introduced the concept of network slicing to improve flexibility and scalability [130]. Network slicing allows independent sets of network protocols to run on common physical infrastructure, to support services with conflicting requirements. As an example, video streaming requires high data rates and accepts high latency, while critical IoT usually requires low latency and high reliability while accepting low data rates. Prior to the emergence of 5G, these services could not be supported using the same protocols, but with network slicing, they can be implemented on the same physical infrastructure [15]. Going beyond 5G, the demand for ML services is projected to grow significantly and discussions

have begun on a dedicated network slice for ML in future-generation cellular networks such as beyond-5G and 6G [60], [131], [151], [191]. Given this possibility, the investigation of Wireless for ML becomes relevant not only for local-area networks but also for large-scale cellular networks.

1.1 Related Work

Although the general intersection of ML and wireless communications is currently a prolific field of research that has already generated multiple surveys, there are fewer works reviewing Wireless for ML. The current surveys can roughly be classified into three categories: *ML for Wireless Communications*, *Wireless for ML*, and *Communication-Efficient DML*. We list a set of representative surveys in Table 1.1. A brief description of the three areas follows.

1. **Wireless for ML** uses wireless communication protocols as a method to enable or significantly improve ML training over wireless networks. Unlike in traditional wireless communication, the communication system is not oblivious to the meaning that the bits convey. Instead, Wireless for ML is a task-oriented communication philosophy, where the goal of the communication system is to distill the intelligence carried within the data.
2. **Communication-efficient DML** has the same goal as Wireless for ML but uses different methods. Instead of customizing the wireless protocols, advancements are made by modifying or redesigning the ML algorithm. The results of these works are agnostic to the communication protocol so that they can be applied regardless of the specific technologies used to transmit data.
3. **ML for wireless** uses ML as a method to design wireless communication protocols or other elements for general communication services. Therefore, its goal is the same as in traditional wireless communications, i.e., efficient and reliable transfer of arbitrary data. The communication system should support a wide variety of services and is therefore deliberately oblivious to the semantics of transmitted bits.

Table 1.1: Surveys written within the intersection of ML and communications. The topics of ML for Communications and Communication-efficient DML have been covered in many surveys, unlike Wireless for ML. At most, Wireless for ML has been covered briefly in conjunction with Communication-efficient DML.

Year	Journal	Ref.	Research Area from Figure 1.1
2017	IEEE Communication Surveys and Tutorials	[109]	3
2018	Proceedings of the IEEE	[120]	2
2019	Proceedings of the IEEE	[194]	2
2020	IEEE Communication Surveys and Tutorials	[73]	3
2020	IEEE Communication Surveys and Tutorials	[162]	3
2020	IEEE Internet of Things Journal	[40]	Mostly 2 with some 1
2020	IEEE Communication Surveys and Tutorials	[164]	2
2020	IEEE Internet of Things Journal	[3]	2
2020	IEEE Communication Surveys and Tutorials	[178]	Mostly 2 with some 1
2021	IEEE Internet of Things Journal	[74]	2
2021	Elsevier High-Confidence Computing	[170]	2
2021	arXiv	[54]	Mostly 1 with some 2
This survey			1

In addition to the three categories above, their intersections can be considered as areas of their own, illustrated in Figure 1.1. The intersection of Wireless for ML and Communication-efficient DML considers the co-design of the ML algorithm and the wireless protocol. With such an approach, researchers attempt to reach some global optimality, which is lost when the two problems are treated in isolation. Additionally, one can consider the intersection between Wireless for ML and ML for Wireless, where ML would be used as a tool to design a wireless protocol with the goal of supporting distributed ML services. However, as far as we are aware, no works have been published in this direction. In this survey, we consider all works within Wireless for ML, including its intersections, symbolized by the green crescent in Figure 1.1.

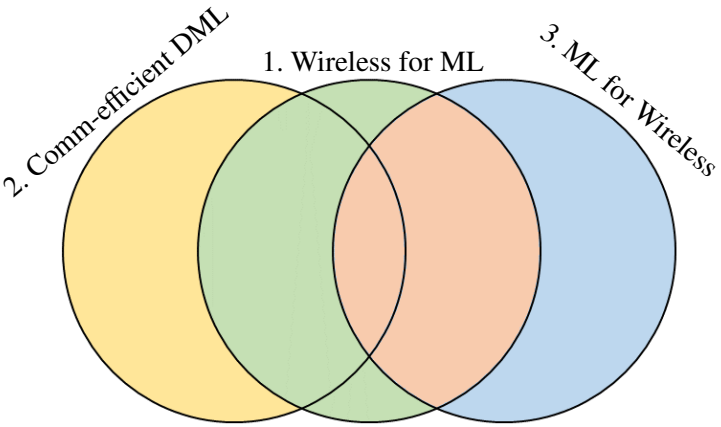


Figure 1.1: Illustration of the relationship between Wireless for ML and related fields. The first circle corresponds to Communication-efficient DML, the second to Wireless for ML, and the third to ML for Wireless. The blue area corresponds to pure ML for Wireless, which is a very prolific field of research that has already generated a large number of review articles. Likewise, the yellow area corresponds to pure Communication-efficient DML which is also a well-covered area. In this survey, we focus on the green moon, i.e., pure Wireless for ML and its intersection with Communication-efficient DML. As far as we are aware, there are no published works in the red area.

Some of the papers in Table 1.1 discuss Wireless for ML, but the treatments there are not extensive since that is not the main purpose of these papers. The closest match to our survey is [54]. However, despite describing some works within Wireless for ML, the paper is not a comprehensive survey of the field, instead its purpose is to introduce a new framework to describe Federated

Learning. We believe that due to this gap, there is currently no one-stop survey that offers an overview of the Wireless for ML literature, which motivates us to write this survey with the following contributions:

- We provide an introduction to important concepts necessary to understand the field as a whole, such as DML, over-the-air computation, and the distinction between generic wireless communication protocols and Wireless for ML;
- We describe the most important works of the field in a concise way to offer a thorough overview of the state-of-the-art, both for analog over-the-air computation and digital communications;
- We discuss several important open problems and future research directions within Wireless for ML;
- We describe a number of application areas where Wireless for ML can provide a benefit to society, such as vehicular communications and virtual reality, and describe the challenges associated with those applications.

1.2 Notation and Organization

All the contributions that we survey are essentially concerned with the solution to a basic problem, namely the training of a classifier over a wireless communication network constrained by the natural characteristics of the wireless channel. Throughout this survey, we assume a centralized architecture where there is a central controller or parameter server (PS) able to make decisions such as user selection, bandwidth allocation, and aggregation frequency control. Such an architecture is representative of most of the wireless networks used today, from large scale mobile to personal area networks. The communication channel is wireless and is thus subject to fading, additive noise, and bandwidth restrictions. The training dataset is always carried by user devices and the training algorithms will always be chosen to minimize a loss based on the global dataset. Unless specified otherwise, the network consists of one PS, i.e., the base station (BS) or the access point (AP), and K user devices, e.g., IoT devices, user equipments (UEs), or other wireless devices. Each device (say the k^{th}) carries a subset \mathcal{D}_k of the global dataset \mathcal{D} and the PS carries no data. The global dataset consists of N training samples and corresponds to the

union of data available at all the user devices. For communication, the uplink h_k and downlink g_k channel coefficients corresponding to the k^{th} UE are of particular importance. Figure 1.2 illustrates the setup, a full list of notation is given in Table 1.2, and relevant abbreviations are given in Table 1.3.

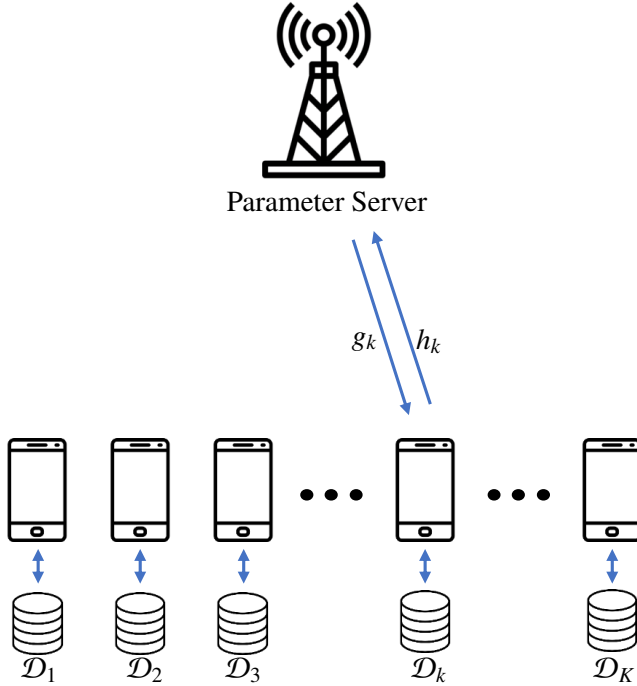


Figure 1.2: Illustration of the PS and wireless network setup used throughout this survey. Current wireless communication protocols substantially hinder or completely block distributed training over this setup. The Wireless for ML paradigm is an approach to tackle such hindrances and blockages.

The rest of this survey is organized as follows: Section 2 provides a primer on DML and in particular Federated Learning (FL). In Sections 3 and 4, we survey the Wireless for ML literature for over-the-air computation and digital communication, respectively. In Section 5, we discuss the open problems in Wireless for ML within both analog over-the-air computation and digital communications. Then, in Section 6, we discuss applications supported by Wireless for ML. Finally, we give some concluding remarks in Section 7.

Table 1.2: Reference list of commonly used variables in this survey. Ordered alphabetically and by case.

Variable	Interpretation
B	Bandwidth available to the learning system
\mathcal{D}_k	Dataset carried by device k
E	Number of epochs
K	Number of user devices
M	Number of antennas at the parameter server
N	Number of data samples in the global dataset
N_k	Number of data samples stored at device k
\mathcal{S}^t	Set of selected devices at iteration t
T_{round}	Time for federated learning communication round
β	Learning rate
η	Post-transmission scalar
$\nabla f(\mathbf{w})$	Gradient of function f evaluated at \mathbf{w}
b_k	Ratio of total bandwidth allocated to device k
d	Number of model parameters in \mathbf{w}
$f(\mathbf{w})$	Empirical risk function of the global model \mathbf{w}
g_k	CSI in downlink direction from server to device k
h_k	CSI in uplink direction from device k to server
$l(\mathbf{w})$	Loss function for parameter \mathbf{w}
p_k	Uplink power allocated to device k
v	Additive white Gaussian noise
\mathbf{w}^t	Global model parameters at iteration t
\mathbf{w}_k^t	Local model parameters for device k at iteration t
\mathbf{x}	Input or feature of data sample
\mathbf{y}	Output or label of data sample

Table 1.3: Reference list of most abbreviations used in this survey.

Acronym	Phrase
ADMM	Alternating Direction Method of Multipliers
AirComp	Over-the-air Computation
BAA	Broadband Analog Aggregation
BPSK	Binary Phase-Shift Keying
BS	Base Station
CML	Centralized Machine Learning
CoCoA	Comm-efficient distributed dual Coordinate Ascent
CoMAC	Computation over Multiple-Access Channels
CSI	Channel State Information
DML	Distributed Machine Learning
DP	Differential Privacy
DSGD	Distributed Stochastic Gradient Descent
ESN	Echo State Network
FD	Federated Distillation
FedAvg	Federated Averaging
FL	Federated Learning
IID	Independent and Identically Distributed
IRS	Intelligent Reflective Surface
IoT	Internet of Things
LTE	Long Term Evolution
MIMO	Multiple Input Multiple Output
ML	Machine Learning
MSE	Mean Square Error
OFDMA	Orthogonal Frequency Division Multiple Access
PS	Parameter Server
RRM	Radio Resource Management
SGD	Stochastic Gradient Descent
SISO	Single Input Single Output
SNR	Signal to Noise Ratio
QoE	Quality of Experience
UAV	Unmanned Aerial Vehicle
VR	Virtual Reality
ZF	Zero-Forcing

2

Primer on Distributed Machine Learning

In conventional ML, model training is considered to take place in centralized settings, where the processing capability and training datasets are locally available within one computational device. Therefore, in the context of a wireless network, centralized machine learning (CML) models and algorithms require that all training data must be transmitted from the user devices to a central server. While possible, such an approach has two major practical problems. Firstly, this approach relies on a complete sacrifice of privacy since all the user devices must be willing to reveal their entire datasets to the server. In many cases, this lack of privacy renders training impossible, since the users may not be willing to share their data, it would be considered immoral to collect the data, or the privacy of the users is legally protected. Secondly, the size of training datasets is an important factor in determining the performance of ML models, where larger datasets generally generate better results [152]. This naturally leads to a desire of training with massive datasets, which is very challenging to communicate over a wireless network [71]. Recently, DML has been proposed as a means to overcome these challenges. Differently from CML, DML works over a dataset distributed among many devices, and optionally performs even distributed training.

In DML methods, the training can be distributed entirely across the devices, which represents the decentralized architecture; or it can be done jointly by a central PS and the devices, which represents the centralized architecture. In this survey, we focus on the centralized architecture within DML because it provides strong guarantees in terms of communication bandwidth usage, latency, parameter update frequency, and desired fault tolerance [14]. Figure 1.2 shows the centralized architecture, in which the K devices communicate only with the PS, which usually has higher computational power than the other devices and is not necessarily represented by a single server (see [14, Section 7] for other PS infrastructures). Notice that the centralized architecture with PS is similar to the operation of current cellular networks, Wi-Fi, and IoT networks with a central controller that could be an app, router, or an IoT device. In DML, the training goal is global, i.e., all the participating devices have a common goal.

The purpose of this section is to introduce the basic concepts in DML, which we will use and will refer to often in the rest of the survey, especially for what concerns the mathematical concepts of ML and their relation to wireless communication protocols. In the following, we discuss the learning goal of CML methods before specifically explaining the learning goal of DML methods, and then we introduce FL methods.

2.1 Problem Formulation for Centralized Machine Learning

We discuss herein the general CML problem of supervised learning, i.e., the problem of labeling unseen data based on information from a set of labeled training data [21]. The common learning goal is to represent a prediction function $h : \mathcal{X} \rightarrow \mathcal{Y}$ from an input space \mathcal{X} to an output space \mathcal{Y} such that, given $\mathbf{x} \in \mathcal{X}$, the value $h(\mathbf{x})$ offers an accurate prediction about the true output $y \in \mathcal{Y}$. Hence, the prediction function h should minimize a risk measure over an adequately selected family of prediction functions, termed \mathcal{H} . Instead of optimizing over a generic family of prediction functions, it is commonly assumed that the prediction function h has a fixed form and is parametrized by a real vector $\mathbf{w} \in \mathbb{R}^d$ with dimension d .

Then, for some $h(\cdot; \cdot) : \mathbb{R}^{d_x} \times \mathbb{R}^{d_y} \rightarrow \mathbb{R}^{d_y}$, the family of prediction functions is $\mathcal{H} \triangleq \{h(\cdot; \mathbf{w}) : \mathbf{w} \in \mathbb{R}^d\}$, where d_x and d_y are the dimensions of \mathbf{x} and \mathbf{y} , respectively.

To meet the learning goal, it is necessary to obtain the prediction function in the family \mathcal{H} that minimizes the losses due to inaccurate predictions. To this end, we assume a loss function $l : \mathbb{R}^{d_x} \times \mathbb{R}^{d_y} \rightarrow \mathbb{R}$ that given an input-output pair (\mathbf{x}, \mathbf{y}) , yields the loss $l(h(\mathbf{x}; \mathbf{w}), \mathbf{y})$ [21]. Notice that $h(\mathbf{x}; \mathbf{w})$ and \mathbf{y} represent the predicted and true outputs, respectively. The model parameter \mathbf{w} is chosen such that the expected loss incurred from any input-output pair is minimized. The loss functions $l(\cdot; \mathbf{w})$ can be either convex on \mathbf{w} , such as when used for linear regression or binary classification (linear support vector machine (SVM)), or nonconvex, such as when used for image classification using neural networks with several layers. Let us assume that the losses are measured with respect to a probability distribution $\Pr(\mathbf{x}, \mathbf{y})$ in the input-output space $\mathbb{R}^{d_x} \times \mathbb{R}^{d_y}$, i.e., $\Pr : \mathbb{R}^{d_x} \times \mathbb{R}^{d_y} \rightarrow [0, 1]$. Then, the objective function we want to minimize is

$$R(\mathbf{w}) = \int_{\mathbb{R}^{d_x} \times \mathbb{R}^{d_y}} l(h(\mathbf{x}; \mathbf{w}), \mathbf{y}) d\Pr(\mathbf{x}, \mathbf{y}) = \mathbb{E}[l(h(\mathbf{x}; \mathbf{w}), \mathbf{y})], \quad (2.1)$$

in which $R : \mathbb{R}^d \rightarrow \mathbb{R}$ is the expected risk given a parameter vector \mathbf{w} with respect to the probability distribution $\Pr(\mathbf{x}, \mathbf{y})$. The minimum expected risk, denoted by $R(\mathbf{w}^\star)$ with $\mathbf{w}^\star := \arg \min_{\mathbf{w}} \{R(\mathbf{w})\}$, is also known as the *test* or *generalization error*. Therefore, the common learning goal in ML can be understood as the minimization of the *test error* [64].

To minimize the expected risk in Eq. (2.1), it is necessary to have complete information about the probability distribution $\Pr(\mathbf{x}, \mathbf{y})$ of the input-output pair. However, such minimization is not possible in most situations because complete information of $\Pr(\mathbf{x}, \mathbf{y})$ is not available. Therefore, the practical learning goal becomes the minimization of an estimate of the expected risk R . To this end, we assume that there are $N \in \mathbb{N}$ independently drawn input-output data samples $\mathcal{D} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N \subseteq \mathbb{R}^{d_x} \times \mathbb{R}^{d_y}$, and we define the empirical risk function $R_N : \mathbb{R}^d \rightarrow \mathbb{R}$ as

$$R_N(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N l(h(\mathbf{x}_i; \mathbf{w}), \mathbf{y}_i). \quad (2.2)$$

With the empirical risk, the optimization problem is as follows:

$$\underset{\mathbf{w}}{\text{minimize}} \quad \frac{1}{N} \sum_{i=1}^N l(h(\mathbf{x}_i; \mathbf{w}), \mathbf{y}_i), \quad (2.3)$$

in which the minimization of R_N is the practical optimization problem that needs to be solved when performing supervised learning. The minimum empirical risk is also known as the *training error* and can be understood as an estimation of the *test error* [64].

To solve optimization problem (2.3), several optimization algorithms have been proposed using stochastic optimization methods, such as stochastic gradient descent (SGD), with or without the use of data partitioning into batches [21]. A general SGD method solves iteratively optimization problem (2.3), with iterations given by

$$\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t - \beta \sum_{i \in \mathcal{S}^t} \nabla f_i(\mathbf{w}^t), \forall t \in \mathbb{N} \quad (2.4)$$

where $\mathbf{w}^t \in \mathbb{R}^d$, β is the learning rate, $f_i(\mathbf{w})$ is the composition of the loss function l and h evaluated at sample i , and \mathcal{S}^t is a set with cardinality N^t . The sum in (2.4) depends on the set \mathcal{S}^t and may represent pure SGD, batch gradient descent, or a joint approach with minibatch SGD [21]. For $N^t = 1$, Eq. (2.4) represents the pure SGD method, and the unique element of the set \mathcal{S}^t corresponds to the seed ξ^t of the sample pair $(\mathbf{x}^t, \mathbf{y}^t)$, which is chosen randomly from $\{1, \dots, N\}$. For $N^t = N$, Eq. (2.4) represents the batch gradient descent method, in which the gradient is evaluated for all samples N and taken into account at each iteration t . For $1 < N^t < N$, Eq. (2.4) represents the minibatch SGD method, in which N^t is termed batch size and all N^t elements of \mathcal{S}^t are chosen randomly at each iteration t . The iterations are evaluated until they reach a minimizer of the empirical risk R_N .

In practice, the training error is evaluated by solving optimization problem (2.3) with N samples; whereas the test error is evaluated by comparing the prediction function $h(\mathbf{x}_j; \mathbf{w})$ using previously unseen input $\mathbf{x}_j \in \mathcal{X}$ to predict the corresponding output $\mathbf{y}_j \in \mathcal{Y}$. Specifically to classification problems, the classification accuracy is the ratio between the number of correct predictions and the number of incorrect predictions given by the learning model. Throughout the survey, the learning performance of ML algorithms is related to the training and test errors. Specifically to classification problems, we refer to the performance as classification accuracy.

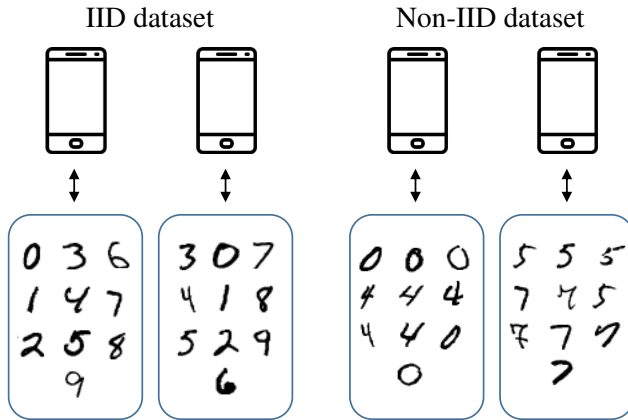


Figure 2.1: Illustration of IID vs. non-IID data for the MNIST dataset. Non-IID data is common when data is generated by user devices, making model convergence of DML harder.

2.2 Problem Formulation for Distributed Machine Learning

Differently from traditional ML methods, in DML the N samples are either split or generated at different K devices. For simplicity, we assume throughout the survey that the samples are generated at K devices. Let us denote by \mathcal{D}_k the dataset owned by device k and $N_k = |\mathcal{D}_k|$ as the cardinality of \mathcal{D}_k . As a consequence of generating data in a distributed fashion, local data distributions at each device can be skewed in comparison to the global dataset. Consider the classic scenario of digit recognition. In the global MNIST dataset, we have 10% representation of each digit 0-9 [89], and the digits are independent and identically distributed (IID). If the digit distribution of the local datasets does not match the global one, the data can be non-IID, see Figure 2.1.

With the splitting of the data across the devices, the empirical risk function can be rewritten as

$$f(\mathbf{w}) = \sum_{k=1}^K \frac{N_k}{N} F_k(\mathbf{w}) = \sum_{k=1}^K \frac{N_k}{N} \sum_{i \in \mathcal{D}_k} f_i(\mathbf{w}). \quad (2.5)$$

When the dataset owned by the K devices are IID, then $\mathbb{E}_{\mathcal{D}_k}[F_k(\mathbf{w})] = f(\mathbf{w})$, where the expectation $\mathbb{E}_{\mathcal{D}_k}[\cdot]$ is taken over the dataset of device k . If the datasets owned by the K devices are non-IID, the loss function $F_k(\cdot)$ at device

k could be an arbitrarily bad approximation of the function $f(\cdot)$ [59], thus harming the convergence.

Similar to the traditional ML methods, DML methods use many optimization techniques to minimize the empirical risk in Eq. (2.5), such as distributed stochastic gradient descent (DSGD) [203], consensus optimization [113], and the alternating direction method of multipliers (ADMM) [22]. For both data distributions, the centralized DML architecture needs to exchange information about the parameters between the K devices and the PS. Depending on the optimization technique used, this information, commonly referred to as just *model*, can be the parameter \mathbf{w} , the gradient $\nabla F_k(\mathbf{w})$, the gradient update $\nabla F_k(\mathbf{w}^t) - \nabla F_k(\mathbf{w}^{t-1})$, or the parameter update $\mathbf{w}^t - \mathbf{w}^{t-1}$. In this survey, we will use *model* to refer specifically to the parameter variable \mathbf{w} , which can be local for each device, \mathbf{w}_k , or global, \mathbf{w} .

To improve the applicability of DML methods, there are still many challenges for both DML architectures and different optimization solvers. Some of these challenges are communication efficiency, system and statistical heterogeneity, and privacy loss [86], [92]. The communication efficiency is related to the massive number of messages that need to be exchanged between the PS and a large number of devices, which may cause high latency and increase the convergence time. The systems heterogeneity is related to the different storage, computing, and communication capability of each device; whereas the statistical heterogeneity is related to the different distribution of the data each device may have, which makes the sample distribution among the devices non-IID. Privacy loss can happen when the devices have sensitive data that they do not wish to expose to other devices and/or the PS.

Some algorithms to tackle the challenges above have been proposed [77], [86], [107], including communication efficient distributed dual coordinate ascent (CoCoA) and CoCoA+ algorithms [77], [107] that address challenges related to communication efficiency. One of these algorithms is FL, which has been proposed as a solution aimed at solving all the challenges mentioned above and thus differs from the CoCoA and CoCoA+ algorithms.

2.3 Federated Learning

In FL methods [86], a common global model is trained in a distributed manner using the PS within the centralized architecture of DML. The common sce-

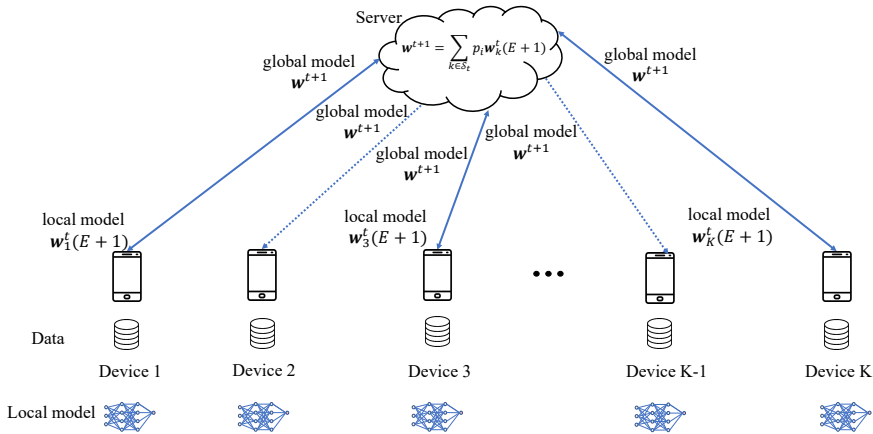


Figure 2.2: Federated learning scenario with K devices and a PS. Only the federated devices, the ones that belong to set S_t , participate in the learning at communication round t and not all K devices, represented by the solid and dotted lines, respectively.

nario in FL is the one in which the number of participating devices is typically large and they have slow or unstable connections; the devices do not want to share their raw data with the PS or other devices; and there is heterogeneity in the data across the devices and in the computation capabilities. Note that this scenario implies that FL methods must address challenges in terms of communication efficiency, privacy, and systems/statistical heterogeneity, which are the challenges common to DML methods mentioned in Section 2.2.

Figure 2.2 shows a FL scenario in which K devices and the PS use FL towards a common global goal, which is to minimize the empirical risk. Notice that only the federated devices that belong to the set S^t participate in the learning at communication round t and not all the K devices, represented by the solid and dotted lines in Figure 2.2, respectively. The raw data is kept locally at each device, and the devices participating in the training minimize their local functions F_k , which means having SGD updates similar to Eq. (2.4) for E local iterations. Then, the devices send to the PS their local models $\mathbf{w}_k^t(E+1)$ that minimize the local functions at communication round t . The PS aggregates the local models with proper scaling p_k and broadcasts the global model \mathbf{w}^{t+1} to all participating devices. Therefore, FL improves the communication efficiency by avoiding many communication rounds with the PS due to the transmission of the updated model only after the local iterations;

takes into account the devices heterogeneity by the possibility of different number of local iterations at the devices as well as the selection of devices to participate in the training; and finally, it improves privacy by not sharing the raw data.

The first FL method proposed was federated averaging (FedAvg) [110, Algorithm 1]. With FedAvg, the PS randomly selects a fraction C , $0 < C \leq 1$, of the K devices to participate in the training at global iteration t , i.e., the set \mathcal{S}^t has cardinality $\max(\lceil CK \rceil, 1)$. Each device $k \in \mathcal{S}^t$ minimizes the local function F_k by computing the gradient $\nabla F_k(\mathbf{w}_k^t)$ and iterating E local iterations (named epochs) applying the updates as

$$\mathbf{w}_k^t(i+1) \leftarrow \mathbf{w}_k^t(i) - \beta \nabla F_k(\mathbf{w}_k^t(i)), \forall i = 1, \dots, E. \quad (2.6)$$

Notice that the gradient $\nabla F_k(\mathbf{w}_k^t(i))$ can be obtained using SGD with different batch sizes. After E epochs, device k sends $\mathbf{w}_k^t(E+1)$ to the PS, which aggregates the local models of the participating devices at iteration t to generate the updated global model as

$$\mathbf{w}^{t+1} \leftarrow \sum_{k \in \mathcal{S}^t} \frac{N_k}{N} \mathbf{w}_k^t(E+1). \quad (2.7)$$

Then, the PS sends the updated global model \mathbf{w}^{t+1} to all participating devices, and the iterative process between the devices and the PS continues until global convergence is achieved, at which $\mathbf{w}^{t+1} = \mathbf{w}^*$. To measure the rate of convergence, we use

$$f(\mathbf{w}^{t+1}) - f(\mathbf{w}^*), \quad (2.8)$$

which may be a decreasing function in t . For FL, there is no closed-form expression for this convergence rate, so model performance cannot be predicted before training. However, for certain scenarios there are theoretical guarantees based on the upper bounding of the convergence rate [35].

Since FedAvg was proposed in [110], many other FL methods have been proposed and investigated for many scenarios, including sparse and/or quantized FL [34], [127], [146], private FL using differential privacy [166], fair FL [94], and FL over wireless communications [158]. For an in-depth overview of recent FL methods and applications, we refer the reader to [84], [92], [116], [178], [180].

2.4 Summary

DML methods overcome some challenges of traditional ML methods, and similarly, FL methods overcome some challenges of more general DML methods. Recently, FL has been investigated due to its robustness participation by massive numbers of users, privacy-enhancing properties, and both statistical and device heterogeneity. However, there are still several challenges that DML and FL need to overcome when applied to Wireless for ML.

In the following sections, we will discuss novel wireless protocols that are specifically designed to address the challenges of Wireless for ML. Specifically, we will discuss in Section 3 the use of analog over-the-air computation and in Section 4 the use of digital RRM for ML.

3

Analog Over-the-air Computation

3.1 Primer

A prominent theme in the Wireless for ML literature is a method called either AirComp or computation over multiple-access channels (CoMAC) [112]. We dedicate this subsection to explain the basics of CoMAC.

In wireless communications, significant attention is placed on the avoidance of interference. As an example, orthogonal frequency division multiple access (OFDMA) splits the wireless spectrum into small blocks of time and frequency and allocates these blocks to different users in the network. Such a system achieves nearly interference-free communication at the cost of significantly reducing the available transmission time and bandwidth for each user. In contrast, CoMAC actively promotes interference. Multiple users are allocated the same time and frequency resources, causing their signals to combine in the air. By carefully designing precoding functions at the transmitting devices, the signal superposition property can be leveraged to calculate functions of the transmitted messages over-the-air [190]. CoMAC addresses applications when the receiver does not need the individual messages from the transmitting devices, but only some function them, for example their sum or average.

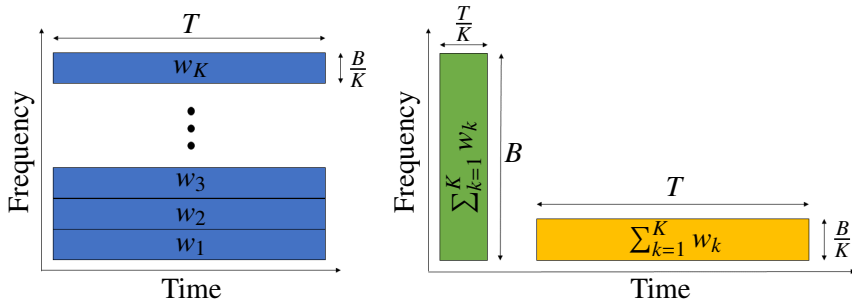


Figure 3.1: Communication-efficiency comparison of FDMA and over-the-air computation. In the left sub-figure, FDMA (blue bars) communicates the K messages over separate frequency bands, and the sum is calculated locally at the parameter server. In contrast, over-the-air computation communicates all K messages jointly by leveraging the interference of simultaneous transmissions. Since over-the-air computation allows complete bandwidth sharing, this results in either K times less latency (green vertical bar) or bandwidth (yellow horizontal bar).

As CoMAC does not allocate orthogonal radio resources, it could be mistaken for the recently proposed non-orthogonal multiple access (NOMA) schemes. However, unlike CoMAC, NOMA needs to enable the reconstruction of the individual messages, and thus employs successive interference cancellation to eliminate interfering signals. This cancellation is possible only by introducing diversity in either the power or code domain [76], [132].

Since CoMAC does not allocate orthogonal resources or introduce additional diversity, the spectral efficiency grows linearly with the number of devices in the network [1]. Consider the network setup from earlier in Figure 1.2 and that the server wants to calculate a sum of K messages w_k carried by the user devices; over-the-air computation would then require approximately K times less resources to communicate this sum. As illustrated in Figure 3.1, the protocol designer can choose to crystallize this resource efficiency to reduce latency and/or bandwidth consumption.

3.1.1 Sum function example

Herein, we will demonstrate how CoMAC calculates a sum function over the air. We follow the system model illustrated in Figure 1.2 where h_k denotes the channel from device k to the PS. If all devices transmit simultaneously over the same frequency band, the server will receive a linear combination of these

signals due to the additive nature of simultaneously arriving electromagnetic waves. Denote the signal transmitted by device k to be w_k . Then, the received signal r at the server is

$$r = \sum_{k=1}^K h_k w_k + v, \quad (3.1)$$

where v is an additive white Gaussian noise (AWGN) term. Here we assume that the antenna of the PS does not saturate. Because of the fading, the received sum is weighted by different weights for each device, and the server is unable to reconstruct the desired function $\sum_{k=1}^K w_k$. A possible solution is to let the user devices pre-equalize their channel. So instead of transmitting w_k directly, they transmit $z_k = w_k/h_k$. This way, the server would receive

$$r = \sum_{k=1}^K h_k z_k + v = \sum_{k=1}^K w_k + v. \quad (3.2)$$

Except for the noise term, this corresponds to the desired function. Considering that the signal strength is the sum of K signals, while the noise v is the same as if a single device transmitted, r is generally a good estimator for the desired sum. While this simple description illustrates the basic idea of CoMAC, there are several simplifying assumptions that must be dealt with in practice. We discuss these next.

Channel State Information

To pre-equalize the channel in Eq. (3.2), the user device must know the channel state information (CSI) of h_k , which cannot be estimated at the device directly. If classical channel estimation were to be employed at device k , the estimated value would be g_k in the downlink direction. A naive solution to this problem is to let the server estimate h_k by having the mobile devices transmit individual preamble signals in the uplink direction, and then feed the CSI back to the mobile devices. However, the transmission of these preambles would require orthogonal transmission of the uplink signals, negating the benefits of over-the-air computation.

Instead, [1] presents a solution based on channel reciprocity. The underlying claim is that forward and reverse channels are the same up to a constant multiplier due to differences in hardware between the transmit and receive

chains. By introducing a calibration stage in which the up- and downlink channels are measured for each sensor device k , this constant multiplier can be found as $c_k = h_k(0)/g_k(0)$, where $h_k(0)$ and $g_k(0)$ are the uplink and downlink channels at time 0, respectively. Since the multiplier remains constant, subsequent communication rounds can calculate the uplink channel using the downlink CSI measured with the broadcast from the server as $h_k(t) = c_k g_k(t)$. However, this solution has only been tested for stationary nodes. For dynamic scenarios, such as cellular and vehicular communications, other solutions have been proposed [43], [82], [195]. In the interest of brevity, we refrain from discussing these other methods here, but blind over-the-air computation is discussed in Section 3.4. Note that this calibration stage is not required in traditional digital communications since the channel can be equalized at the receiver, so the transmitters do not require knowledge of h_k .

Synchronization

A second problem arises from an inherent assumption in Eq. (3.2), which is that the transmitted signals arrive simultaneously at the server. Even small synchronization errors can lead to major estimation errors because the sum is calculated with an analog signal. Synchronization would be required at a symbol-level, which may be difficult to achieve with traditional synchronization. To overcome such a problem, multiple novel approaches have been proposed, including: dedicated hardware that transmits sinusoidal tones [2], longer transmission blocks to reduce the synchronization requirement [56], or the "timing advance" functionality of Long Term Evolution (LTE) networks [156]. Additionally, if some devices are far away from the PS, there might be a need to estimate the propagation latency and compensate at the transmission.

Power control

The pre-equalization scheme from (3.2) assumes that the devices have the capability to transmit $z_k = w_k/h_k$. However, if the device is experiencing a deep fade, h_k will be a very small number, thereby requiring a tremendous amount of power for pre-equalization. With practical devices, the peak power is constrained, and such a scheme is unfeasible. To get around this constraint,

several researchers have formulated power control problems [29], [103] which introduce a post-transmission scalar $\sqrt{\eta}$. This scalar is applied by the PS after receiving the sum, yielding

$$r = \sum_{k=1}^K \frac{h_k z_k}{\sqrt{\eta}} + \frac{v}{\sqrt{\eta}}. \quad (3.3)$$

With the newly introduced η , the amplitude required for pre-equalization of the channel changes to $z_k = \sqrt{\eta} w_k / h_k$. If the post-transmission scalar is selected to be $\eta < 1$, the required transmission power is reduced, enabling more devices to invert their channels. However, a reduction of η also leads to an increase in the relative noise power. This tradeoff leads to the power control problem, which aims to optimally select z_k and η without exceeding transmission power constraints. We discuss this problem further in Section 3.3.8.

3.1.2 Summary

By promoting interference, CoMAC allows all devices to share the electromagnetic spectrum without allocating orthogonal radio resources to each user. This technique achieves throughput gains approximately proportional to the number of participating devices, which is a tremendous improvement even with a relatively small number of users. The main drawback of the method is that the individual messages cannot be reconstructed at the receiver, which limits the application to scenarios where a function of the messages is sufficient. In the preceding paragraphs, we gave a simple example that demonstrates how channel pre-equalization CoMAC can be used to calculate the sum function. However, this method has several practical issues such as strong demands on CSI, stringent synchronization requirements, and limited peak transmission powers at the user devices.

3.2 Over-the-air Computation For Distributed Machine Learning

As explained in Section 2.3, the model aggregation step of FL consists of transmitting multiple local models from the user devices to the PS and then computing a weighted mean of these updates to generate the next iteration of the model, see (2.7). The individual local models are not needed at any point of the FL algorithm, only this weighted sum. As such, the sum can be directly

computed over-the-air instead of separately transmitting each model vector and then averaging at the PS. This basic idea has served as the foundation of a large body of work that explores the impact of CoMAC on DML and that extend the idea further.

As illustrated in the CoMAC example from Section 3.1, non-uniform fading across the network is a major challenge for estimating the desired function. We presented a power modulation solution based on channel reciprocity and inversion, which is the standard method to overcome this challenge for single input single output (SISO) networks. However, for multiple input multiple output (MIMO) networks, alternative solutions have been proposed, such as the blind CoMAC which utilizes channel-hardening to avoid the CSI acquisition problem. Additionally, the consideration of MIMO comes with other interesting CoMAC-solutions such as beamforming, cell-free massive MIMO, and intelligent reflective surface (IRS)-assisted CoMAC. With this in mind, the remainder of this section is split into two parts, SISO and MIMO. A comprehensive list of papers on CoMAC for ML is given in Table 3.1 and 3.2.

Table 3.1: Summary of the SISO AirComp for ML literature. The papers are ordered according to when they are covered in the survey.

Topic	Ref.	Summary
Broadband Analog Aggregation	[200]	FL using AirComp over a broadband channel with truncated channel inversion to handle fading.
Gradient Sparsification	[9]	Sparsification of gradients combined with error accumulation for compression before transmitting.
	[10]	Extension of [9] to consider fading channels, uses truncated channel inversion.
	[11]	Performance comparison of [10] scheme, sequential digital transmission, and BAA.
	[50]	Utilization of temporal structures in the gradient updates to form a Bayesian prior in the gradient estimation step.
Federated Distillation	[4]	Trains by communicating model outputs instead of model parameters. Over-the-air computation is used to combine model output vectors for each class.

Topic	Ref.	Summary
Training with Noisy Gradients	[137]	Proposal of gradient-based multiple-access scheme that does not cancel the fading effect but operates directly with noisy gradients.
	[138]	Convergence rate analysis for gradient-based multiple-access.
Data Sharing	[153]	DSGD training using combined gradients. Introduces data redundancy to combat non-IID. data.
Analog Federated ADMM	[47]	Second-order training algorithm with CoMAC communication.
Digital Aggregation	[196]	First digital over-the-air computation method using one-bit quantization of gradients.
	[81]	Clustered digital over-the-air computation that minimizes the probability of incorrect gradient sign estimation.
Power Control	[28]	Optimal selection of pre- and post-processing scalars using FL bounds.
	[189]	Estimation of gradient statistics to improve power control for Federated Learning.
Retransmissions	[67]	Proposal of retransmission-based model update scheme that enables an estimation-communication tradeoff.
	[68]	Development of heuristic to predict the optimal number of retransmissions.
Differential Privacy	[96]	Uses the noise added naturally by the wireless channel to enhance data privacy for free.
Byzantine Attacks	[147]	Considers the grouping of participating devices to mitigate Byzantine attacks.
Device-to-Device Communication	[173]	First decentralized machine learning scheme using over-the-air computation.
	[143]	Decentralized SGD with gradient tracking and variance reduction.
Bayesian Learning	[97]	Proposes the channel-driven Monte-Carlo sampling method that leverages channel noise to estimate the posterior distribution of ML parameters.

Table 3.2: Summary of the MIMO AirComp for ML literature. The papers are ordered according to when they are covered in the survey.

Topic	Ref.	Summary
Blind Learning	[106]	The assumption of channel knowledge at the user devices is lifted. Instead, multiple antennas at the PS is employed to alleviate the fading effect.
	[6]	Extension of [106] to consider imperfect channel estimation at the PS.
Nonlinear Estimator	[79]	Recovering the average of local models sent from the devices using their sparsity with a nonlinear estimator.
Cell-Free Massive MIMO	[160]	FL in a cell-free massive MIMO framework with CSI estimation using CoMAC pilot transmission.
Beamforming and User Selection Co-Design	[179]	Optimal user scheduling based on tradeoff between maximizing participation and limiting distortion from aggregation error.
Intelligent Reflective Surfaces	[165]	Optimized beamforming, user selection, and phase-shift control via intelligent reflective surfaces (IRSs) to maximize device participation.
	[102]	Optimization over upper bound on FL loss to find proper phase-shift control, device selection, and beamforming for IRS FL.
	[101]	Channel state information free transmission via IRS.
	[70]	Energy minimization with IRS-assisted over-the-air computation.

3.3 Review of SISO Over-the-air Computation

3.3.1 Broadband analog aggregation

The first paper to suggest CoMAC as multiple access for FL appears to be [197]. This paper presents a short case study that compares the latency of orthogonal transmission with CoMAC under identical conditions. The case study displays a significant reduction in latency, ranging from one to three orders of magnitude, with minor sacrifices in terms of classification accuracy. Later on, the same group presented a fully-fledged scheme called broadband analog aggregation (BAA) in [200]. Similar to LTE, the BAA scheme divides

the spectrum into resource blocks (RBs). However, instead of dedicating each RB to a single user, the blocks are dedicated to one element of the model update vector. This way, all K users can transmit their model updates simultaneously over the same RBs to calculate the weighted sum of model updates from (2.7) over-the-air.

As we explained in Section 3.1.1, channel pre-equalization is used to generate the sum function (3.2). As a consequence of this scheme, the receive SNR is identical for every user, because devices with weaker channels compensate by transmitting at higher powers. In BAA, devices with sufficiently weak channels are excluded from training, since they are unable to pre-equalize their channels. With this in mind, we consider the inclusion of a post-transmission scalar $\sqrt{\eta}$ as in (3.3). If $\sqrt{\eta}$ is reduced, more devices are able to invert their channels, which increases device participation. In the context of FL, higher participation means a larger training dataset. As such, the reduction of $\sqrt{\eta}$ increases data quantity. However, the receive signal to noise ratio (SNR) is:

$$\text{SNR} = \eta \frac{\left(\sum_{k=1}^K w_k\right)^2}{\sigma_z^2}, \quad (3.4)$$

which is proportional to η . Therefore we have a tradeoff between data quantity and receive SNR. In [200], the authors isolate this tradeoff and coin the term communication (SNR)-learning (data quantity) tradeoff. This tradeoff appears in many CoMAC-FL systems and is important to consider when optimizing such systems.

3.3.2 Gradient sparsification

Although the BAA scheme significantly reduces the communication load for FL, it does not consider improvements in terms of the ML algorithm. In contrast, the next paper we discuss utilizes gradient sparsification together with CoMAC to further reduce the communication cost. Gradient sparsification is based on the observation that up to 99.9% of the gradient exchange in DSGD is nearly redundant [95]. Therefore, a majority of the gradients can be discarded with minimal reductions to learning accuracy.

In [9], the combination of gradient sparsification and CoMAC appeared for the first time. In this paper, a simple channel model without fading was considered. In [10], the same scheme was extended to consider fading chan-

nels, where truncated channel pre-equalization was used to generate the sum. Finally in [11], an experimental comparison of three different FL approaches (orthogonal transmission, BAA, and gradient sparsification with CoMAC) is conducted. The study focuses on training an MNIST classifier, it assumes a limited transmission budget in terms of time slots, and compares the final test accuracy after the transmission budget is out. The results reveal that both CoMAC approaches outperform orthogonal communication with up to 40% better classification accuracy. The study also indicates that the inclusion of gradient sparsification has substantial benefits, with up to 10% classification accuracy over BAA.

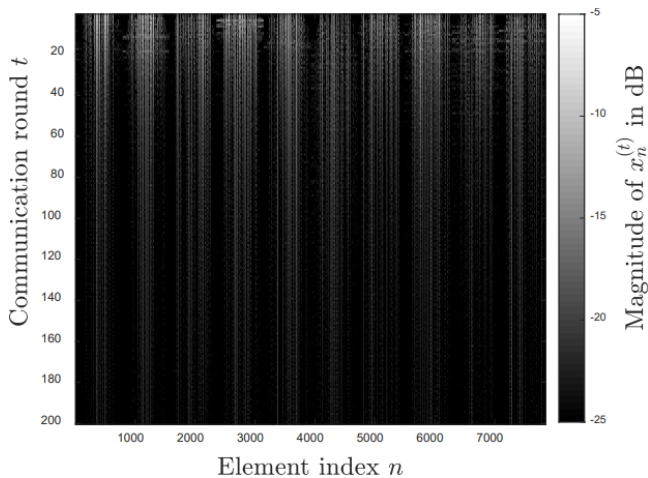


Figure 3.2: Illustration of the temporal structure of gradient updates in over-the-air FL. The amplitudes of the gradient elements are encoded in grayscale over time. We can see that the sparsity of the gradient is roughly retained through time, which can be exploited to improve the estimation of the local models. Source: [50]

In a more recent work [50], the authors noticed a predictable structure in the aggregated gradients. From Figure 3.2, we can see that the amplitude of the different gradient elements changes slowly over time, more or less retaining the sparsity structure through the entire training process. To model this structure, [50] uses two independent Markov chains for the support and amplitude. By combining this simple model and the stored gradients from previous communication rounds, a prior belief on the gradient can be formed. As explained in Section 3.1, over-the-air computation always results in noise, therefore the PS must estimate the gradient after receiving the uplink signal.

If there is no prior information, the best estimate is to just directly use the received signal. Instead, this paper uses Bayesian estimation with the prior belief from the Markov chain model to make a better estimation. In the numerical comparisons of [50], this approach strictly outperforms the results from [11].

3.3.3 Federated distillation

As explained in Section 2, FL achieves consensus by sharing locally trained models with the PS. These local models can become enormous when considering deep neural networks with millions of neurons, such as the VGG models that consist of $d=130$ -140 model parameters [149]. With this in mind, there have been attempts to develop an alternative to FL called federated distillation (FD). In FD, model outputs are communicated instead of the model parameters [80]. In other words \mathbf{w}_k^t from Eq. (2.7) is replaced with the average of the local model outputs, thus communicating an \mathbb{R}^{d_y} vector instead of an \mathbb{R}^d vector. Often classification problems have $d_y \leq 100$ labels but millions of parameters d , causing a reduction in the number of transmitted bits by many orders of magnitude. Upon receiving these model outputs, the server calculates their average and communicates it back in the downlink. These average model outputs can then be used by the devices to train their ML models. As we wish to focus on the communication protocol, we refrain from explaining how these model outputs are used for training and refer the interested reader to [80].

In [4] FD is combined with CoMAC. First, each device combines the model outputs over multiple training samples, generating one value for each label. Then, for each label, a global average is calculated over-the-air. This can lead to massive reductions in communication cost but unlike the gradient sparsification schemes, FD does have a noticeable drop in classification accuracy. The numerical study conducted in [4] suggests it can be between 1-20% lower than FL.

3.3.4 Training with noisy gradients

Unlike all papers we have surveyed so far, which used channel inversion to combat fading, see Eq. (3.2), the authors of [137] suggest just transmitting without doing any precoding. Such a scheme has the advantage of not requiring a channel estimate, and a generally simpler implementation. However, since

fading is not inverted, the received local models at the PS represent noisy and distorted versions of the transmitted local models. The corresponding distorted average is then used to perform the FL update directly. An important contribution of this paper is an upper bound on the FL loss, arguably the first bound that considers AirComp. An extended convergence and numerical analysis is given in [138] containing simulation results based on the Million Song Dataset [17]. The results reveal comparable or slightly worse FL loss compared to a digital scheme but with significantly reduced energy consumption.

3.3.5 Data sharing

In Section 2.1, we explained that there is an important distinction to make between IID and non-IID training data on the devices. With non-IID data, there is no guarantee that the locally trained models resemble the global model, which can significantly harm FL performance. In extreme examples, non-IID data can harm the classification accuracy by up to 55% [193].

Realistically, we should always expect FL data distributions to be non-IID. For instance, environmental monitoring devices will have different data distributions depending on sensor location, text prediction algorithms depend on user behavior, and body sensor systems depend on the physiology of the host. To combat this, [153] introduces a data sharing phase into CoMAC for ML, where each user device shares its dataset with a small number of neighbors before training begins. The study considers the same communication scheme as in BAA, but performs data sharing before training begins. The included numerical study on the MNIST dataset considers highly non-IID data distributions where each device only carries samples of one digit. They show that classification accuracy goes from 72% to 82% by having each user device share its dataset with just one neighbor.

3.3.6 Analog federated ADMM

When over-the-air computation is used to calculate a sum, channel pre-equalization is employed to counteract heterogeneous fading over the network, as explained in Section 3.1.1. When all devices perform pre-equalization, the over-the-air computation will produce the desired result. However, some devices may be unable to pre-equalize their channel due to limited transmission power. To solve this problem, all papers surveyed up to this point simply

exclude those devices from participating. Instead, [47] proposes the first over-the-air computation algorithm that overcomes channel perturbations without pre-equalization; the method is based on a novel FL framework rooted in ADMM, which the authors call analog federated alternating direction method of multipliers (A-FADMM).

For the sake of inclusivity, we will not assume that the reader is familiar with ADMM and avoid mentioning specifics of the ADMM algorithm in this subsection. Instead, we focus on the model update that is communicated by A-FADMM, because it differs significantly from what we see in Section 3.1.1 and has the interesting property of avoiding channel pre-equalization. For the reader who wants a deeper look into ADMM, we refer to [22].

The global model update of A-FADMM is given as follows:

$$\mathbf{w}^{t+1} \leftarrow \frac{1}{\sum_{k \in S^t} |h_k|^2} \sum_{k \in S^t} \frac{N_k}{N} \left(|h_k|^2 \mathbf{w}_k^t(E+1) + h_k \lambda_k^t(E+1)/\rho \right). \quad (3.5)$$

Compared to the standard FL update in (2.7) there are two major differences. Firstly, the channels h_k have been directly incorporated into the FL problem formulation, and secondly there are now two new variables $\lambda_k^t(E+1)$ and ρ which represents the dual variable of the ADMM algorithm and a penalty variable respectively. The semantics of these variables can be ignored for the sake of this discussion. Notice that the channel h_k is a factor both for the local model $\mathbf{w}_k^t(E+1)$ and the dual variable $\lambda_k^t(E+1)$. This means that the user devices can transmit

$$N_k(h_k^* \mathbf{w}_k^t(E+1) + \lambda_k^t(E+1)/\rho), \quad (3.6)$$

where h_k^* is the conjugate of h_k . Then, using over-the-air computation, the PS receives

$$\sum_{k \in S^t} N_k \left(|h_k|^2 \mathbf{w}_k^t(E+1) + h_k \lambda_k^t(E+1)/\rho \right) + v. \quad (3.7)$$

If this expression is multiplied by $1/(N \sum_{k \in S^t} |h_k|^2)$ it generates the desired function from (3.5) in expectation. Therefore, channel pre-equalization is not required and A-FADMM has the advantage of avoiding device exclusion completely. This directly increases the training data quantity, which should improve learning performance. On the other hand, one could argue that the multiplication by h_k leads to weak transmission signals, thereby potentially reducing the SNR compared to channel pre-equalization.

In addition to proposing A-FADMM, the authors of [47] prove that the algorithm converges for convex functions under time-varying channels. The convergence rate is also evaluated numerically by training with the MNIST dataset. The results suggest that A-FADMM converges faster than both traditional FL with over-the-air computation as well as digital ADMM without over-the-air computation.

3.3.7 Digital aggregation

Current telecommunications infrastructure is almost exclusively designed for digital communications. Because of this, the implementation of analog CoMAC in large scale networks becomes problematic. To avoid constructing new analog chipsets at large scale, [196] proposes an adaptation to over-the-air aggregation which would be compatible with current transceivers. The proposed protocol is based on 1-bit SGD [135] which uses single-bit compression of gradient descent updates. Specifically, each element of the user devices' gradient vectors $\nabla F_k(\mathbf{w}^l)$ takes one of two values (1 or -1). These binary gradient vectors are combined to form an element-wise majority vote at the PS.

The CoMAC protocol represents these SignSGD gradients using one of the two Binary Phase-Shift Keying (BPSK) symbols. Because the two BPSK waveforms are inverted versions of the other, wireless superposition will correctly represent the addition of +1 and -1. In other words, the sum of a +1 BPSK waveform and a -1 BPSK waveform will be zero, given that their amplitudes are identical. Therefore, the CoMAC sum function would directly calculate the desired element-wise majority vote over the air.

The performance of one-bit digital CoMAC is compared to BAA in [200] by training a classifier for MNIST. The results suggest that the classification accuracy of digital CoMAC is nearly identical to BAA, with less than 1% loss of accuracy. This result indicates that CoMAC could potentially be implemented in cellular networks without requiring significant change in the hardware. As most CoMAC schemes, perfect synchronization is assumed in both theoretical analysis and numerical simulation. One could argue that one-bit digital CoMAC is more sensitive to synchronization errors since it relies upon cancellation of two opposite BPSK waveforms, unlike analog CoMAC which only requires additive powers.

A second digital CoMAC scheme was proposed in [81] which introduces a clustered structure for the majority vote operation of the network. Rather than having all devices communicate directly with the server and thereby casting their votes in a “direct democracy” system, the authors propose intermediate relays that serve as representatives. This breaks the vote into two stages, where all the devices first cast their votes to their closest relay, which uses majority vote to generate a new gradient vector. Then in the second stage, the relays vote to the PS in a “representative democracy” system. The selection of relays can be done in a smart way so that the relays have similar channel strengths to the server. If the strengths are similar, the channels do not necessarily need to be inverted, which alleviates the need for CSI estimation and improves the probability of success in the final majority vote. Simulation results show improvement over a cluster-free system both in terms of the gradient estimation at the PS and the classification accuracy.

3.3.8 Power control

As explained in Section 3.1.1, there is a power control problem associated with CoMAC. The problem arises because the pre-equalization of the channel is restricted by limited transmission power at the user devices. In this section, we discuss the problem of optimal power control.

In all the papers mentioned up to this point, sub-optimal power control was used. Specifically, devices with fading below a certain threshold were excluded from participation and the remaining devices perfectly inverted their channels. Instead, [29], [103], [184] study the problem more rigorously to minimize the estimation error under transmission power constraints. They consider the problem (3.8) to minimize the mean squared error between the received signal and the desired sum:

$$\begin{aligned} \min_{\mathbf{p}, \eta} \quad & \mathbb{E} \left[\left(\sum_{k=1}^K \frac{h_k p_k \Delta \mathbf{w}_{n,k}}{\sqrt{\eta}} + \frac{v}{\sqrt{\eta}} - \sum_{k=1}^K \mathbf{w}_k \right)^2 \right] \\ \text{s.t.} \quad & p_k \leq P_{\max}, \forall k \end{aligned} \quad (3.8)$$

where p_k is the transmission power of device k , P_{\max} is the peak power constraint, and the remaining variables are defined in (3.3). There are two sources of error, one is the misalignment error caused by devices being unable to pre-equalize their channels and the second is the error induced by the

AWGN ν . The post-processing scalar η controls the tradeoff between the two, where a higher η reduces the noise-induced error directly, but indirectly worsens the misalignment error by making it harder to invert the channel. The specific problem posed in (3.8) is solved to a global minimum in both [29] and [103], given certain simplifying assumptions.

In the context of FL, the power-control problem affects both the convergence rate and final accuracy of the ML model. In [28] (later extended in [27]), a similar set-up to (3.8) is used, with pre- and post-processing scalars for power control, but with the objective function replaced by an upper bound on FL convergence. The proposed scheme vastly outperforms the device-exclusion scheme in terms of prediction accuracy.

Another work [189] considers the use of gradient statistics to evaluate the expectation in (3.8). For known gradient statistics, the authors find the optimal solution in closed form using the mean squared norm and the squared multivariate coefficient of variation. In a practical scenario, these statistics would be unknown, but the solution can be used in conjunction with live estimates of the statistics to determine good pre- and post-processing scalars.

3.3.9 Retransmissions

In digital communications, there is a well-known tradeoff between communication rate and error probability. For example, the modulation order n determines the number of bits $\log_2(n)$ that can be transmitted in a single symbol. Simultaneously, a higher modulation order makes the demodulation problem harder, thereby increasing the probability of error. As such, the modulation order acts as a tradeoff between communication rate and error probability. Similarly, forward error-correcting codes can be used to correct erroneously demodulated bits at the receiver, but simultaneously introduce redundant bits which reduces the rate of communication. In contemporary digital communication protocols, it is common practice to adaptively select the modulation order and coding rate with respect to the estimated channel [58] but in analog CoMAC such a practice does not exist. With this in mind, the authors of [67] consider a retransmission-based scheme to analyze the tradeoff communication rate and estimation error for over-the-air FL.

The scheme presented in [67] is similar to the power control proposals [29], [103] except that the model update $\Delta \mathbf{w}_{n,k}$ is transmitted M times in

the uplink instead of just once. At the receiver, these M transmissions are collected and their arithmetic mean is used to generate the next iteration of the global model update. This way, the signal part of the transmission combines constructively, while the noise part is random and can therefore combine destructively. This scheme is analyzed by proving an upper bound on the FL loss, which reveals that the convergence rate is strictly increasing in M . To make a fair comparison between transmission with $M = 1$ and $M > 1$, the authors perform a simulation study in which the uplink transmission budget is fixed to \bar{C} , such that only \bar{C}/M communication rounds can be performed. Despite using M times fewer communication rounds, the simulation study indicates that there are scenarios in which $M > 1$ achieves higher classification accuracy after consuming the communication budget. Therefore, the performance of Over-the-Air FL can be improved by including retransmissions, without incurring additional costs in terms of latency or energy consumption.

In a follow-up study [68], the optimal choice of M is studied further and a heuristic is developed to predict M^* before training begins. Numerical results indicate that the heuristic is generally successful at identifying M^* , including the case when $M^* = 1$. As such, the system can predict when the conditions are not right for retransmissions and select one-shot uplink transmission.

3.3.10 Differential privacy

Compared to CML, FL makes a step towards data privacy by keeping the data local at the users. However, sharing the local models or the gradients may reveal sensitive information about the users data [30], [111]. Adding a level of uncertainty to the local models or the gradients computed at the users can enhance the privacy of user data at the cost of lower utility. Differential privacy (DP) is a privacy measure that quantifies the amount of information leakage about individual data points by measuring the sensitivity of the revealed statistics to a change at a single data point, and it is widely adopted as a promising privacy measure.

It is shown in [136] that the additive nature of the wireless multiple access channel from the user devices to the PS provides local DP guarantees for the devices where the privacy leakage per device is scaled with $1/\sqrt{K}$. If the channel noise is not sufficient to satisfy the DP target, a subset of the devices

add power constrained artificial noise that benefit all the devices. Alternatively, [85] introduces an energy efficient differentially private approach for FL over wireless networks by scaling down the transmit power rather than adding noise to the transmit signals at the devices. In general, a certain level of DP can be achieved for free with the analog transmission from the devices due to the noise added by the wireless multiple access channel which can act as a privacy-inducing mechanism [96].

3.3.11 Byzantine attacks

An unfortunate consequence of the distributed and privacy-preserving nature of FL is that malicious users can transmit modified model updates with the intention of disrupting the training process [159]. Even a single client can seriously harm the performance of the end model [19]. These malicious clients are called Byzantine, and their attacks are called Byzantine attacks. As a countermeasure, a recent idea has emerged for distributed computation among agents called “coded computing”. This idea consists in transforming the client’s information by functions which on the one side hide the client’s information, and on the other side can add robustness to the computation because the PS applies another function that attempts to minimize the effect of the Byzantine attacks [13], [18], [53], [124]. However, these countermeasures generally rely upon detecting anomalies in individual model updates, which is difficult for AirComp where the average model updates are calculated directly over the air.

This gap in security for over-the-air FL is a serious concern. A first step to address this concern can be found in [147]. Specifically, [147] proposes that the participating devices are split into G groups, with K/G devices per group. Each group is allocated its own time slot for over-the-air computation, thereby generating G received model updates at the PS. With these G vectors, the PS can apply coded computing methods to mitigate potential Byzantine Attacks. In [147], the authors prove that the proposed algorithm converges to a neighborhood of the optimal \mathbf{w} when the number of attackers are less than $G/2$. A such, the choice of G acts as a tradeoff between communication efficiency and security.

3.3.12 Device-to-device communications

Up until this point of the survey, we have only considered distributed ML over star networks, which can be modeled by the multiple-access channel and therefore leverage AirComp. In this subsection, we briefly discuss work on device-to-device communication over more general network topologies. For such topologies, there is no dedicated PS and the devices are only able to communicate with their immediate neighbors in a single hop. Therefore, the FedAvg algorithm cannot be directly applied for ML training. However, there are other methods, such as decentralized SGD, which are guaranteed to converge under assumptions of noiseless communication, convexity and connectivity [154].

In [173], the problem of decentralized SGD with over-the-air computation was studied for the first time. They consider a connectivity graph model with probabilistic blockages due to shadowing, where unblocked channels are described by Rayleigh fading and AWGN. To enable AirComp in such a network, they propose a scheduling policy that aims to select as many non-interfering subnetworks with star topologies as possible for each time slot. Once the subnetworks are identified, a two-step iterative procedure is initiated. In the first step, over-the-air computation is leveraged to communicate the average gradient to the center of each star network. In the second step, all centers broadcast the received gradient average to the edge devices. This way, every device in the subnetwork knows the arithmetic mean of the gradients after two time slots. This scheme is evaluated numerically by training an MNIST classifier for $K = 8$ devices with randomly generated connectivity graphs. The results suggest that over-the-air computation converges with significantly fewer communication blocks than orthogonal digital communication, but reaches a lower accuracy as the number of communication blocks approach infinity.

In [143], a similar setup to [173] is considered, but with the added consideration of gradient tracking [125] and variance reduction [172]. Gradient tracking refers to the introduction of an auxiliary variable into the optimization problem of decentralized SGD that tracks the average gradient of all devices in the network. With such an auxiliary variable, linear convergence can be guaranteed with a constant step size [125]. With variance reduction, an iterative estimator of the batch gradient is designed, whose variance progressively approaches zero as the parameter vector approaches a local minimizer. With

variance reduction, the error floor of SGD is eliminated even with a constant step size, which is not possible for vanilla SGD. In [143], the proposed decentralized scheme is proven to converge linearly under standard convexity assumptions, fully-connected graphs, and bounded gradients.

3.3.13 Bayesian learning

While ML has displayed impressive accuracy for many classification tasks, ML models are not perfect and will occasionally make mistakes. For certain applications, such mistakes could have unwanted consequences that limit the applicability of ML. To mitigate the harm caused by ML mistakes, it is desirable to consider models with the ability of assessing the certainty of its predictions. Consider the application of fall detection. Multiple accelerometers can be attached to a patient's body with the intent of detecting falls and alerting the patient's medical assistant [169]. A common problem with these systems is that alerts are communicated to assistants for normal, healthy activity which causes unnecessary and unwanted visits [25]. If the alerts were sent to the assistant together with a measure of the model's uncertainty, the assistant could make a better decision on whether they should intervene. In statistics, the term *predictive uncertainty* is used to describe this virtue, where many state-of-the-art ML methods, such as neural networks, are poor at quantifying predictive uncertainty, and tend to produce overconfident predictions [88].

Bayesian learning is a popular method to quantify the predictive uncertainty of neural networks, in which a prior distribution is specified upon the parameters of a neural network and then, given the training data, the posterior distribution over the parameters is computed. If we compare this to traditional ML, we can say that traditional ML generates a point estimate of the parameters, i.e., one instantiation of the weights and biases of the neural network, while Bayesian learning attempts to generate a full distribution over the parameters, i.e., the posterior distribution. Exact calculation of the posterior distribution is in general intractable, so approximate methods are used to generate an estimate of the distribution, such as Monte-Carlo sampling [88]. Once estimated, the posterior distribution is leveraged to quantify the uncertainty of any given prediction. The interested reader can refer to [168] for a detailed description of the uncertainty quantification.

In [97], distributed Bayesian learning is brought into the wireless setting using over-the-air computation. The main contribution of the paper is the introduction of an idea called *channel-driven Monte-Carlo sampling* where the channel noise is utilized as an integral part of the sampling for estimating the posterior distribution. If accounted for, the channel noise combined with the analog transmissions in over-the-air computation may not cause harm to the performance of the learning. This is in contrast to FL, where the noise generally slows down convergence and should be compensated for, as discussed in Sections 3.3.8 and 3.3.9. In [97], the channel-driven Monte-Carlo method is analyzed analytically by means of a convergence proof and numerically by extensive simulations.

3.4 Review of MIMO Over-the-air Computation

3.4.1 Blind learning

Similar to traditional MIMO communications, the channel estimation effort of CoMAC systems is in the opposite direction of traditional SISO communication, since equalization is performed at the transmitter instead of at the receiver. This is problematic, because while the downlink channel can be estimated using the model broadcast of FL, the uplink channel can not. To solve this problem, one can use channel reciprocity together with a calibration factor to estimate the uplink channel [1] but this is both more expensive (it requires a calibration stage) and less precise than downlink channel estimation. In CoMAC, this problem is exacerbated since the CSI knowledge is used to achieve signal alignment, and poor channel estimation will result in distorted function computation [201]. With this in mind, the channel hardening phenomenon of MIMO communications carries particular importance for CoMAC. In [57], channel hardening is leveraged to perform over-the-air computation without deterministic channel knowledge at any node in the network. Specifically, the authors quantify the gap in performance between a system with full CSI knowledge and one with only statistical knowledge at the PS and no CSI knowledge at the user devices. For a network with $M > 1$ antennas at the PS and $K > 1$ single-antenna sensor devices, they prove that this performance gap approaches zero as $KM \rightarrow \infty$.

In the previous section, we highlighted [9] that introduces gradient sparsification to over-the-air FL. In [106], this scheme is extended to consider blind learning. The main contributions of this work are to propose a CoMAC-based FL technique that requires no transmit CSI from the devices and to provide insights into how the number of antennas affect learning accuracy. The numerical results show that for $M = 2K^2$, the accuracy nearly matches a non-fading channel. For a lower number of antennas $M = 2K$ the accuracy drop compared to the non-fading channel is about 5%.

The work in [106] is then further extended in [6] to consider imperfect CSI at the PS. The authors show that the lack of perfect CSI results in an additional zero-mean interference term with a variance proportional to $1/M$. Similarly, worst-case analysis shows that the imperfect CSI results in slower convergence but that the effect is inversely proportional to the number of antennas. Finally, numerical analysis on the MNIST and CIFAR-10 datasets reveal significant performance improvement as M increases with a more pronounced effect when channel estimation is not perfect.

3.4.2 Nonlinear estimator

One challenge in FL over wireless networks is the presence of a noisy shared wireless medium from the typically abundant users to the PS, over which the users transmit their local models or gradients. The goal is to deliver users' signals to the PS as accurately as possible. Equipping the PS with multiple antennas can improve communication reliability between the users and the PS, where multi-antenna transmission and/or reception beamforming techniques can be employed [6], [160], [179]. However, the above works consider only linear beamforming techniques at the multi-antenna PS to estimate the signals transmitted from the users.

In general, a linear beamforming technique at a multi-antenna receiver does not lead to any optimal estimation performance [79]. Instead, the authors in [79] design an estimator based on the sparsity of the gradient vectors computed at the users. Motivated by this sparsity, a compressive sensing approach in the user domain is employed, where the gradient vectors at different users are permuted using different patterns such that only a small subset of the users transmit non-zero entries at each dimension. This results in a sparse transmitted signal from the users, and using this sparsity, the PS employs a

nonlinear estimator to recover the average of the gradients almost accurately. This approach is extended in [78] by employing the gradient compression technique introduced in [9], [11] to reduce the transmission bandwidth over the wireless multiple access channel from the users to the PS.

3.4.3 Cell-free massive MIMO

Recently, a new architecture for multi-user MIMO, called cell-free massive MIMO, has emerged. In cell-free massive MIMO, a large number of APs collaboratively serve users over the same time/frequency resources [114]. All APs collaborate through a backhaul network, enabling fine synchronization that can be used for conjugate beamforming in the downlink and matched filtering in the uplink. The main advantage of the cell-free architecture is the broad coverage due to the high number of APs. This is especially important for over-the-air FL since the communication quality of CoMAC for ML is determined by the device with the worst channel [200].

In [160], a comprehensive scheme combining cell-free massive MIMO and FL was proposed. The FL process is divided into four steps, starting with CSI acquisition and ending in global model aggregation at the centralized PS. Unlike the previous subsection, the proposed scheme does not utilize blind transmission but it is able to estimate the channel using non-orthogonal transmission. By making all sensor devices transmit their pilot sequences simultaneously over the same bandwidth, the channels can be estimated using multiple measurements received by the large number of APs. Numerical results show that cell-free massive MIMO can reduce training time by up to 33% when compared to massive MIMO with collocated antennas.

3.4.4 Beamforming and user selection co-design

Due to the communication-learning tradeoff, see Section 3.3.1, user selection should be made to strike a balance between receive SNR and data quantity. The solution to this problem in the SISO case was to set a fading threshold based on a power constraint and only include users below that threshold. By introducing multiple antennas at the AP, [179] instead proposes receive beamforming to maximize the participating users while ensuring that the aggregation error is constrained. The proposed user selection and beamforming scheme is compared to a semidefinite relaxation baseline and a global optimization approach

with exponential time complexity. In terms of probability of feasibility, the proposed approach was significantly better than semidefinite relaxation and nearly identical to the global optimum. Additionally, the approach was used to train on the CIFAR-10 dataset [87] and the proposed approach achieved nearly double the relative classification accuracy of semidefinite relaxation.

3.4.5 Intelligent reflective surfaces

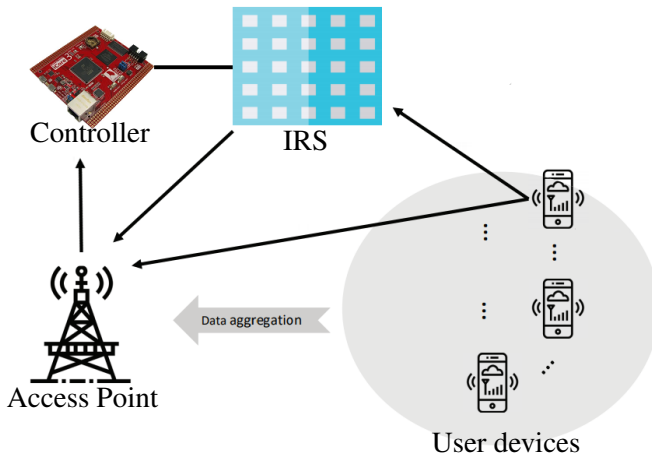


Figure 3.3: Illustration of an IRS for over-the-air computation. In this figure, the uplink data transmission is assisted by the reflective surface to improve the wireless channel. To control the phase shifts of the reflected signals, the AP communicates with a controller attached to the reflecting surface.

The IRS is a recent technological development that has received strong interest from both academia and industry [192]. The purpose of an IRS is to introduce a “mirror” for electromagnetic waves that can be tuned to reflect incident signals toward the intended receiver. The surfaces consist of passive reflective elements which adjust the phase shift of the incoming signal, effectively creating a MIMO effect. In addition to the reflecting elements, a controller is installed that allows for APs to configure the phase shifts, illustrated in Figure 3.3. For the case of CoMAC, [82] was the first paper to propose a joint beamforming and IRS phase shift design to minimize the aggregation error. The paper showed very significant potential with up to

4 orders of magnitude lower estimation error than an IRS-free propagation environment.

For the case of FL, the IRS can be used to enable higher user participation. As we know from the previous Subsection 3.4.4, it is natural to constrain the number of participating devices to ensure that the aggregation error falls below an acceptable level. In [165], the authors proposed a joint beamforming, user selection, and IRS phase shift design to maximize the number of participating devices. The resulting scheme was able to approximately double the number of participating users compared to an equivalent system without an IRS, which can improve the test accuracy by up to 20% under the right conditions.

The maximization of user participation clearly has a positive impact on learning, but it is a rough proxy for the classification accuracy, which is the metric of interest. To address this issue, [102] found an upper bound on the FL loss under the IRS over-the-air setup and proposed an optimization problem that incorporates the loss function, thereby targeting the accuracy more directly. The simulation results of [102] reached nearly the same test accuracy as training over an error-free channel, outperforming [82] by 3%-30% depending on the experimental setup.

Besides improving the classification accuracy, IRSs can also be used to enable blind transmissions without having a large antenna array [101]. When using an IRS, blind transmission can be achieved even with single-antenna devices and single-antenna APs. However, the system is not completely blind, as it still requires receive CSI at the PS. Since there is no CSI at the transmitter, the devices cannot invert their channel before transmitting. Instead, [101] proposes that the devices transmit with maximum power, and the PS configures the IRS phase shift vector to achieve the desired function over-the-air. Such an approach achieves a significantly worse aggregation error than a system with CSI at the transmitter, but the error is still sufficiently low to achieve a comparable classification accuracy. Since FL works well with some level of noisy updates, the 4 orders of magnitude reduction from the IRS design can be excessive, opening up for designs that are less efficient in terms of mean squared error (MSE).

In [70], the authors investigated the use of multiple IRSs and over-the-air computation to support the deployment of FL. In their considered model, the devices can directly transmit FL models to the BS or using IRS. The authors jointly optimized the device selection, phase shift matrix, decoding vector, and

power control so as to minimize the energy that the devices use to transmit and train FL models. Simulation results comparing communication with and without an IRS reveal that the energy consumption of the FL training can be reduced by approximately an order of magnitude by transmitting via an IRS.

4

Digital Communications

4.1 Primer

The CoMAC systems discussed in the previous section provide an attractive solution to the DML problem. However, the technology is dependent on prerequisites that can be difficult to realize in practical scenarios, such as very stringent synchronization and customized hardware. Due to the challenges with CoMAC, digital communications still has to be considered as a basis for DML. Within digital communications, we consider orthogonal communication methods that leave the physical layer as it is. Then, the attention is placed on the data link and network layer, with a particular emphasis on RRM protocols for DML.

As explained in Section 1, the problem of DML differs in several ways from that of general data communication. These differences result in new constraints in terms of computational complexity, training time, training data, and more. In this setting, general data communication protocols perform poorly, motivating the design of digital communication protocols tailored to support DML. In this primer, we will discuss some of these differences in more detail to better understand why new digital protocols are needed.

4.1.1 Fairness

In traditional RRM, the well-known water-filling method [182] allocates more transmission power to users experiencing a good channel. This method leads to very efficient spectrum utilization, but generally leads to some users with no allocated power. Therefore, despite utilizing spectrum less efficiently, max-min-fairness protocols are often used to ensure a minimum level of service for all users in the network [176]. This sacrifice is not reasonable for FL since the participation of every user is not necessary to train a good model. In fact, if our goal is to maximize the classification accuracy of the ML model, the data-importance discussion in the previous section indicates that we should be deliberately treating users in a discriminatory manner, contradicting the demands on user fairness. Even if data-importance is not considered, there is no reason to sacrifice spectrum utilization to ensure user fairness for FL.

4.1.2 Training data

It is well-known that supervised ML performance is intricately connected to the quality and size of the training dataset. Therefore, we would ideally utilize every collected data point to train machine learning models. However, over resource-constrained wireless networks, this is not always possible. Therefore, we are posed with the problem of optimally selecting which data points to utilize. In centralized machine learning, this problem is related to which data points are communicated to the server, and in DML, the problem is related to which devices should participate (and thereby their datasets). One useful metric to guide such a selection is data importance (discussed further in Section 4.3), which can be utilized to value one data sample over another.

4.1.3 Computational capability

Since FL is traditionally a synchronous algorithm, it suffers from a problem known as the *straggler effect*, i.e., the effect where the slowest device acts as a bottleneck while remaining users idly wait for the next communication round [61]. Therefore, the heterogeneity of communication and computational capabilities becomes an important factor to consider for device scheduling and RRM. As an example, more bandwidth could be allocated to slow devices, thereby helping them to compensate for their slow training by communicating their local models quicker.

4.1.4 Energy

Most DML algorithms rely on multiple rounds of communication to reach convergence in the model training process, each of which consumes a significant amount of energy. Additionally, each communication round is associated with a computational task of training the model, which leads to further energy costs. To maintain an acceptable battery level at the training devices, the energy-efficiency of this process is of critical importance. There are specific properties of the FL algorithm that can be utilized to either consume less energy or transfer power from the base station to the user devices. As an example, there is a period of naturally occurring radio silence in FL, when the user devices are doing their local training. During this time it is possible to perform power transfer from the BS to the devices.

4.2 Digital Communications for Distributed Machine Learning

In this section, we have divided the digital DML literature into two categories: importance-aware communication and RRM for FL. The first category considers prioritization schemes that select users based on how valuable their training data is to the ML model. The second category tries to optimize RRM algorithms for FL. A comprehensive list of papers for digital DML methods can be found in Table 4.1 and 4.2.

Table 4.1: Summary of the Importance-Aware Communications literature. Papers that consider both importance-aware communications and radio resource management is covered in Table 4.2. The papers are ordered according to when they are covered in the survey.

Topic	Ref.	Summary
Centralized Learning	[99]	Retransmission protocol with data-sample prioritization.
	[98]	Extension of [99] to consider more advanced ML models such as convolutional neural networks.
	[100]	User selection protocol.
Federated Learning	[55]	Importance-aware user selection step.
	[91]	Comparison of different data importance metrics for user selection step.

Table 4.2: Summary of the Radio Resource Management for Machine Learning literature. The papers are ordered according to when they are covered in the survey.

Topic	Ref.	Summary
Participation Maximization	[118]	Client selection scheme that aims to maximize the number of participants in the Federated Learning training step.
	[174]	Joint client selection and bandwidth allocation considering the <i>later-is-better</i> phenomenon of FL.
	[175]	Joint time slot and bandwidth allocation with multiple co-existing FL services that share wireless resources.
Energy Efficiency	[186]	Joint client selection and bandwidth allocation scheme that aims to minimize the energy consumed for FL training.
	[181]	Joint time slot allocation, bandwidth allocation, and transmit power allocation.
	[42]	Joint time slot allocation, clock frequency optimization, and local accuracy optimization.
Packet Error Impact	[35]	Performs convergence analysis on the impact of packet errors in FL training. Utilizes the resulting upper bound to perform client selection, resource block allocation, and power allocation.
	[133]	Client selection scheme that weighs the ML model update contribution of individual devices based on their probability of successful transmission.
	[83]	Analyzes the convergence of SignSGD-based distributed learning.
Total Time Minimization	[142]	Joint client selection and bandwidth allocation to minimize the total time spent training the ML model.
	[32]	Joint client selection and resource block allocation.
Empirical Classification Error	[163]	Attempts to estimate the classification error empirically and uses this estimate to guide power allocation.
Federated Distillation	[119]	Combines Federated Distillation in the up-link with Federated Learning in the down-link. Also employs data sample mixing to enhance user privacy.

Topic	Ref.	Summary
Batch Size Selection	[129]	Treats hyperparameters of the machine learning algorithm as decision variables for the RRM problem. Specifically, a joint batch size selection and time-slot allocation scheme is developed.
Importance RRM	[128]	Combines importance-aware communication and RRM for FL by considering a client selection scheme. Specifically, the gradient divergence is used to guide the selection of participating devices.
	[32]	Considers update staleness and update drift to develop a joint client selection and resource block allocation scheme.
Energy Harvesting/Power Transfer	[185]	Joint batch size selection, clock frequency optimization, and learning-wireless power transfer tradeoff.
	[148]	Joint local number of iterations optimization and time slot allocation to transmit, compute and harvest energy.
Noisy Downlink	[12]	Digital downlink transmission of the global model is compared to analog transmission.
Federated Meta-Learning	[183]	The combination of meta-learning and FL is considered in a wireless network, where users are scheduled based on a convergence bound.
Empirical Classification Error	[163]	Attempts to estimate the classification error empirically and uses this estimate to guide power allocation.
Federated Distillation	[119]	Combines Federated Distillation in the uplink with Federated Learning in the downlink. Also employs data sample mixing to enhance user privacy.
Batch Size Selection	[129]	Treats hyperparameters of the machine learning algorithm as decision variables for the RRM problem. Specifically, a joint batch size selection and time-slot allocation scheme is developed.
Importance RRM	[128]	Combines importance-aware communication and RRM for FL by considering a client selection scheme. Specifically, the gradient divergence is used to guide the selection of participating devices.

Topic	Ref.	Summary
	[32]	Considers update staleness and update drift to develop a joint client selection and resource block allocation scheme.
Energy Harvesting/Power Transfer	[185]	Joint batch size selection, clock frequency optimization, and learning-wireless power transfer tradeoff.
	[148]	Joint local number of iterations optimization and time slot allocation to transmit, compute and harvest energy.
Noisy Downlink	[12]	Digital downlink transmission of the global model is compared to analog transmission.
Federated Meta-Learning	[183]	The combination of meta-learning and FL is considered in a wireless network, where users are scheduled based on a convergence bound.

4.3 Review of Importance-aware Communications

When traditional communication algorithms are designed to maximize data rate, they are implicitly assigning equal worth to each bit regardless of their information content. This makes sense in classical packet-switched networks since the abstraction of information in the OSI model prohibits the controller from interpreting the payload. However, in DML, data-importance is not uniform [140], thus if we consider that each bit has the same worth, resources are wasted to transmit low-importance data. The non-uniform data-importance for ML stems from two qualities: uncertainty and diversity [72]. Uncertainty refers to the confidence level with which the current model can classify a data sample, and diversity refers to the rarity of the label compared to the remaining training data set. Consider a system for classifying images of animals as in Figure 4.1. Low data-importance images would correspond to something that is easy to classify, such as a simple white background, a common animal, and a natural pose. As either the diversity (rarity of the animal) or uncertainty (difficult pose/background) increases, so does the data importance. By prioritizing samples with high data-importance, ML training is accelerated [167]. Since non-uniform data-importance is common, communication algorithms concerned with learning performance should incorporate uncertainty and diversity in their design by prioritizing high-importance data.

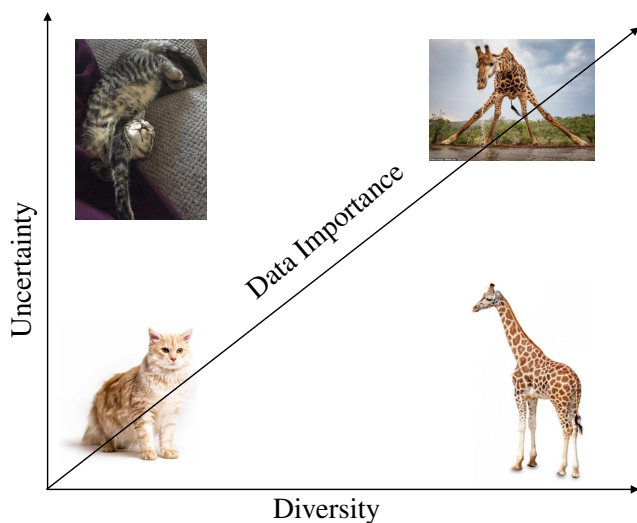


Figure 4.1: Data-importance illustration for image classification. Uncertainty measures how difficult a data sample is to classify, while diversity measures the rarity of the data sample label. The importance of training data is a good metric for prioritizing data samples when communication resources are limited.

This idea of evaluating data samples based on their importance during the training of the classifier model comes from a branch of ML called Active Learning [139]. The problem considered in Active Learning is with regard to the cost of labeling. Using a speech recognition example, the cost would come from having a human interpreter listening to recorded samples and transcribing labels to be used for the ML algorithm. In Wireless for ML, each sample is instead associated with a cost related to transmission, and since we consider supervised learning, the label is already available at the device. Although the fundamental goal is different, the metrics developed in Active Learning for evaluating data samples have been tested for the communication problem and have been shown to reduce the communication cost [55], [98].

4.3.1 Centralized learning

In this section, we consider centralized learning using data distributed over multiple devices in a network. In other words, training only occurs at the PS but wireless communication is still used to collect the data. This scenario is

relevant when the user devices do not have sufficient computational resources to perform local training but are still carrying data relevant for learning. In [98] and [99], the problem of developing an automatic repeat request (ARQ) protocol for ML is considered. Despite being orthogonal, the communication is with analog transmission, so there is always some distortion of the received sample. This distortion can be reduced by taking the mean of multiple transmissions of the same signal, thus improving effective receive SNR. Given a time slot budget, the goal is to maximize the final learning accuracy. Given that the time slot budget is not sufficient to upload every sample to the server, additional retransmissions reduce the total number of samples uploaded for training. This problem gives rise to a communication-learning tradeoff like earlier but based on retransmissions instead of participation.

On top of finding a balance between data quantity and quality, the protocols are designed to prioritize samples with higher importance. Three solutions are suggested in [99], as an example, we discuss "Importance ARQ for binary SVM classification" in detail. The protocol considers the acquisition of a data sample x from a user device. Using a first transmission, the PS estimates the data-importance and then the PS repeatedly requests the device to retransmit x until the effective receive SNR satisfies

$$\text{SNR}(T) > \min(\theta_0 \mathcal{U}_d(\hat{x}(T)), \theta_{\text{SNR}}), \quad (4.1)$$

where T is the number of retransmissions, θ_0 is a scaling factor, $\mathcal{U}_d(\hat{x}(T))$ is the uncertainty measure, and θ_{SNR} is the maximum SNR. The maximum SNR is there to prevent one sample from consuming too many transmissions, and $\mathcal{U}_d(\hat{x}(T))$ is defined as the distance to the SVM boundary, which is an uncertainty measure. The protocol based on Eq. (4.1) will allocate sufficiently many retransmissions for each device to reach their guaranteed minimum SNR. Devices carrying low-importance data are guaranteed lower minimum SNRs and are therefore given fewer retransmissions even if the channel is poor.

Apart from binary SVM classification, [99] contains extensions to multi-class SVM, generic classifiers, and convolutional neural networks (CNNs). According to experimental studies, the protocol outperforms purely channel-aware retransmission protocols in terms of classification accuracy by around 2-3% when training on the MNIST dataset.

In [100], importance-aware user selection is addressed. The devices are scheduled in a time-division manner and take turns to upload a data sample in each time slot. Once again, the radio resources are limited and the problem is to schedule devices in a manner that maximizes the final test accuracy. User selection is based on two factors, the channel quality of each user and the importance of their data. Devices experiencing lower fading are prioritized so that higher data rates are achieved, but only if their data is sufficiently important.

Unlike the retransmission case, it is not obvious how data-importance should be communicated to the PS. The problem lies in that both the model and the data samples are required to measure importance, and they are not present in the same entity. To solve this problem, [100] suggests using popular model compression methods [37], [202] to transmit a lighter version of the ML model to the user devices. This would reduce the size d of the local model $\mathbf{w} \in \mathbb{R}^d$. This way, the user device can evaluate their importance locally, and then inform the PS.

4.3.2 Federated learning

Since FL communicates local models or gradients instead of data samples, there is a need for data importance metrics that can be applied to gradients. In [55], the loss function is proposed as an importance metric. This metric is an uncertainty metric, as it directly describes how difficult a sample is to classify. Additionally, it is cheap to compute by performing inference on the already locally available ML model.

Using the loss function as a metric, the model importance is defined as

$$I_k = \frac{1}{\sqrt{N_k}} \sum_{i=1}^{N_k} l(h(\mathbf{x}_k^i; \mathbf{w}_k), \mathbf{y}_k^i), \quad (4.2)$$

where I_k is the importance of device k 's gradient to the global model, \mathbf{x}_k is the vector of data samples, \mathbf{y}_k are the labels for those samples, and N_k is the number of samples. The importance is evaluated locally at each user and is transmitted in the uplink together with the local model, illustrated in Figure 4.2. The PS takes advantage of I_k to determine which users to schedule for the upcoming communication round.

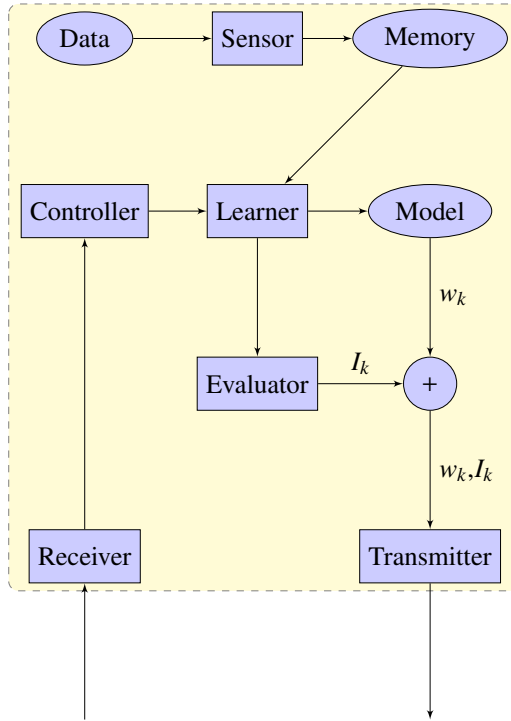


Figure 4.2: User device block diagram for an importance-aware FL system. The importance measurement is calculated by an evaluator block whose output is transmitted together with the local model to the PS.

Unlike the user selection case for CML, the ML model is now naturally present at the user device and the data-importance can be evaluated without the need for transmitting a compressed model to the user devices. As illustrated in Figure 4.2, the user device evaluates the importance locally and appends the data importance to the uplink packet containing the local model. At the start of each communication round, the PS selects a fixed number of users for participation. In vanilla FL the choice would be randomized, but this scheme proposes to select the users with the highest I_k . Using active FL, the proposal in [55] achieves the same performance as vanilla FL using 20-70% fewer epochs. Since the data size of the importance evaluation is small in comparison to the local model, this method also has a negligible overhead.

Rather than using just the loss as the data-importance metric, [91] opts to use a combination of the information entropy and the loss value. Specifically, the elements of the gradient vectors are assumed to follow a random distribution, and the entropy of the gradient elements is used as the data-importance metric. This way of quantifying gradient information originally comes from researchers aiming to perform completely different tasks such as fast tree approximation, community discovery [41], and autoencoding [48], similarly to how data-importance measures from Active Learning had completely different original purposes. In a simulation study, the authors of [91] compares user scheduling based on the gradient norm and gradient divergence to that of gradient entropy. The simulations indicate that the gradient entropy is superior to the norm and divergence when the dataset is non-IID.

4.4 Review of Radio Resource Management for Federated Learning

Because of the differences in objective between Wireless for ML and traditional data communications, direct application of the FL protocol without consideration of practical constraints in wireless communication systems, makes the overall training process inefficient [35], [118], [126]. Instead, RRM protocols should be customized for FL to enable efficient training of ML models using distributed data. Within RRM, we include the allocation of transmission power, bandwidth, time slots, and user scheduling. The objective of RRM for ML is a learning goal, such as the classification accuracy of a model, rather than a general data communication goal, such as data-rate maximization. This difference shapes Wireless for ML RRM in ways that might seem contradictory compared to traditional RRM.

In FL, multiple communication rounds have to be performed until the desired accuracy is reached. As is generally true for iterative algorithms, there is a tradeoff in FL between the computational complexity of each communication round versus the total number of rounds. Specifically, the time per communication round (T_{round}) and the loss decay per round (Δl) must be carefully balanced, as illustrated in Figure 4.3. Both T_{round} and Δl are impacted by the RRM decisions. Additionally, this need for balance leads to new decisions to be made by the PS such as:

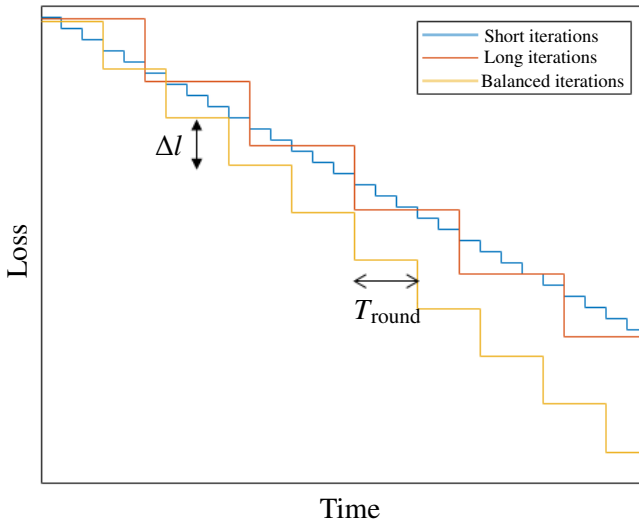


Figure 4.3: Qualitative plot illustrating FL convergence. There is a tradeoff between communication and computation, where slower iterations mean more computation and less communication. For optimal convergence, the RRM protocol should be able to adapt to available computational resources and channel qualities.

1. Deciding how many users will participate in each round;
2. Performing aggregation frequency control, which means to decide how many local training iterations each device performs before communicating its update;
3. Selecting the batch size of each user device's training algorithm.

In the following, we discuss important topics within RRM for FL that the literature has addressed.

4.4.1 Participation maximization

The first paper written about RRM for FL is on the topic of client selection [118]. In the original FL protocol, FedAvg, each communication round begins by the PS selecting a random fraction of clients and sending them the global classifier model [110]. The authors of [118] demonstrated the inefficiency of this selection over wireless networks, due to the heterogeneity of

channel conditions in the network. If clients with poor channels are selected, the uplink transmission is slow and the straggler problem will significantly slow down the training process. Alongside channel heterogeneity, computational resource heterogeneity will lead to the same problem.

To find a better client scheduling policy, an optimization problem is formed. Ideally, the objective function would be the classification accuracy on the test data, but there is currently no closed-form expression for this; see Eq. (2.8). Instead, the number of participants is maximized, which can serve as a rough proxy for the convergence rate [110]. The problem is constrained to exclude slow users, by introducing a deadline for the entire FL algorithm $T_{\text{round}} \in \mathbb{R}^+$. The PS is then subjected to a tradeoff in client selection between the number of participants in each round, and the time required to complete each round. A good T_{round} value is found experimentally, and then the following problem is solved:

$$\begin{aligned} \max_{\mathcal{S}} \quad & |\mathcal{S}| \\ \text{s.t.} \quad & T_{\text{round}} \geq T_{\text{cs}} + T_{\mathcal{S}}^d + T_c + T_{\mathcal{S}}^u + T_{\text{agg}} , \end{aligned} \quad (4.3)$$

where \mathcal{S} is the set of clients selected for the round, T_{cs} is the time to select clients, $T_{\mathcal{S}}^d$ is the downlink transmission time, T_c is the computational time to train the local model, $T_{\mathcal{S}}^u$ is the uplink transmission time, and T_{agg} is the time required to aggregate the local models at the PS. Using the proposed scheme, an extensive experimental study is conducted based on Long Term Evolution (LTE) networks in a mobile edge computing context. The studies indicate that the proposed solution consistently converges faster than out-of-the-box FL regardless of the choice of dataset (Fashion-MNIST or CIFAR-10), and the distribution of the data (IID or Non-IID).

In [174], the authors highlight a phenomenon in FL, termed *later-is-better*, in which the learning rounds are temporally interdependent and have varying significance towards the desired learning outcome. The authors show that, when using FL over a wireless network, it is important to take into account this phenomenon when designing resource management methods to support the FL task. To make use of these findings, the authors formulate a stochastic client selection and bandwidth allocation problem for a finite number of communication rounds while considering finite energy constraints on the clients. The problem aims to maximize the weighted sum of selected clients for a fixed number of communication rounds, whose weights depend on a temporal

parameter to capture the significance of selecting more clients in different communication rounds. The authors show that an increasing sequence of these temporal parameters often results in better FL performance due to a higher number of clients being selected in later rounds of the learning convergence. The constraints include a long-term energy budget on individual clients and feasibility constraints on the bandwidth allocation. Due to the time-varying and unpredictable wireless channel conditions, the authors use Lyapunov optimization to solve the optimization problem and propose an algorithm, named OCEAN, for online client selection and bandwidth allocation. In the results, the authors show that the OCEAN algorithm is adaptive to changing network environments and outperforms greatly other benchmarks that ignore the *later-is-better* effect of FL.

Differently from [118], [174], the authors in [175] consider a scenario with multiple FL services co-existing and sharing resources in a wireless network and propose bandwidth allocation to ensure sufficient client participation for each FL service. Specifically, they propose a two-level resource management framework comprising of intra- and inter-service resource allocation. The intra-service resource management problem aims to minimize the FL communication round time by optimizing the bandwidth allocation among the clients within each FL service. Subsequently, the inter-service resource management problem aims to distribute bandwidth resources among multiple simultaneous FL services. For both problems, the authors analyze both cooperative and non-cooperative FL service providers. For cooperative providers, they propose a distributed bandwidth allocation solution to optimize the overall performance of multiple FL services while considering the fairness among FL services and the privacy of clients and providers. For non-cooperative providers, they propose a new auction scheme with the FL providers as the bidders and the wireless server as the auctioneer, which is able to balance learning accuracy and fairness among the FL services. The bid is based on the bandwidth requested by the FL provider and the price it is willing to pay to get the requested bandwidth. The results show that the proposed solutions outperform other benchmarks, such as equal bandwidth allocation among clients or services, and bandwidth allocation proportional to the number of clients for each service, for various wireless network conditions.

4.4.2 Energy efficiency

Since FL over wireless networks are mostly concerned with either mobile or sensor devices, low energy consumption is critical. In [186], this topic is investigated in a joint bandwidth allocation and client selection scheme. Specifically, the energy consumption of transmitting the local model in the uplink is considered as

$$E_k^{up} = b_k B p_k t_k, \quad (4.4)$$

where b_k is the bandwidth allocation ratio, B is the total bandwidth, p_k is the power allocation in Watt/Hz, and t_k is the model uploading time. The joint bandwidth allocation and user selection scheme is then found by solving

$$\begin{aligned} \min_{b_k, t_k, I_k} \quad & \sum_{k=0}^{K-1} E_k^{up} - \lambda \sum_{k=0}^{K-1} I_k \\ \text{s.t.} \quad & \beta_k \in \{0, 1\}, \\ & \sum_{k=0}^{K-1} b_k = 1, \\ & 0 \leq t_k \leq T_k \end{aligned} \quad (4.5)$$

where I_k is an indicator function that is 1 if device k is selected, and T_k is a maximum time budget for each device. Similarly to the client selection scheme of the previous section, the number of participating devices has been used as a proxy for the convergence rate of the FL model. A numerical study on the MNIST dataset suggests that the proposed scheme outperforms a baseline of selecting every possible client in energy consumption by up to 25% with a 1-2% loss in classification accuracy.

In [181], the energy consumption for computation is considered in addition to transmission. The energy for computing the local model updates at device k is

$$E_k^c = \kappa A_k \log_2 \left(\frac{1}{a} \right) f_c^2, \quad (4.6)$$

where κ is the effective switch capacitance which depends on the chip architecture, A_k is an approximation of the energy consumption per training iteration, a is the local classification accuracy, and f_c is the computation capacity of device k measured in CPU cycles per second. To minimize this energy, the proposed

scheme allows the PS to control the local classification accuracy by selecting the number of local iterations per communication round and the computation capacity of user devices (presumably by giving the training task higher priority on their CPUs). It is worth noting that this paper does not consider FedAvg but uses the distributed approximate Newton-type method (DANE) [141] for training, in which the user devices implicitly uses the local Hessian to compute their updates. In [181], upper bounds on DANE convergence is used to determine the number of local iterations per communication round, thereby leading to different constraints of the radio resource management problem than for FedAvg. Simulation results suggest that the proposed scheme significantly outperforms baseline schemes of equal bandwidth allocation, fixed CPU frequency allocation, and fixed target accuracy allocation.

In [42], the authors propose a novel FL method, named FEDL, to handle heterogeneous user data and physical resources, and employ the proposed FL model to a resource management problem focused on the energy consumption and the communication round time. For the proposed FEDL model, the local model updates at the users minimize a surrogate function of the local objective function using the previous averaged global model and global gradient estimate. The authors provide the convergence analysis and establish the convergence rate of FEDL, which depends on the number of epochs and global iterations. For the resource management problem, the objective is to minimize the energy consumption and the communication round time while considering as variables the computation capacity of the users, the uplink communication time, the desired accuracy for the FEDL method, the controllable parameter for the local surrogate function, the communication time in one global round, and the time to compute one epoch. The proposed problem is non-convex and the authors provide a solution by decomposing the original problem into three subproblems. The numerical results indicate that FEDL outperforms FedAvg in various learning and wireless communication settings.

4.4.3 Packet errors

The studies discussed so far considered perfect CSI and error-free transmission. In [35] instead, the authors consider an outage model where packet-errors can occur, with error probability dependent on the allocated bandwidth and transmission power. The FL averaging step is updated using the outage model

to consider potential packet losses. With this new averaging step, an upper bound on the learning convergence is derived, that reveals the impact of packet errors on the training loss. Using this upper bound, the authors design a joint user selection and bandwidth/power allocation scheme, which converges despite the errors, but after convergence, the following optimality gap remains:

$$\mathbb{E} [l(\mathbf{w}^t)] - \mathbb{E} [l(\mathbf{w}^*)] = C \sum_{i=1}^K N_k (1 - I_k + I_k q_k(b_k, p_k)) , \quad (4.7)$$

where $l(\mathbf{w}^*)$ is the loss of the optimal model, C is a constant depending on the number of training samples and the Lipschitz parameter of the loss function, N_k is the number of training samples at device k , I_k is an indicator that is 1 if device k is scheduled, and $q_k(b_k, p_k)$ is the probability of packet error given the bandwidth and power allocation. This result shows that proper bandwidth and power allocation reduces the optimality gap, leading to better results after convergence.

Similar to [35], the authors in [133] consider a transmission success probability, complementary to the probability of error, which impacts the client scheduling policy and convergence analysis. The FL averaging step uses the success probability together with the scheduling policy and sends in the uplink the difference between the local model after E epochs, $\mathbf{w}_k^t(E)$, and the global model of the current communication round, \mathbf{w}^t . The transmission success probability for each device is derived using stochastic geometry tools in a cellular wireless network considering a fixed number of transmission attempts in the uplink. The authors study two scheduling policies to allocate M resource blocks, the first using uniform sampling of devices without replacement, and the second using a sampling of devices with predefined probability $\{\hat{q}_k\}$ with replacement. Subsequently, they also propose a suboptimal scheduling policy to improve the convergence rate. The authors derive the convergence analysis via an upper bound on the learning convergence and show that unsuccessful transmissions do not affect the convergence rate significantly after proper adjustment of the averaging step. They also show the impact of the number of local epochs, communication rounds, and transmission attempts on the convergence rate. Among the interesting results of [133], the authors prove and show numerical results that other schemes, which do not include the transmission success probability in the global model update step, may converge to the solution of a different FL problem, specifically biased towards the model

of devices with high success probabilities. To avoid such a bias, [133] proposes to weight the model update contribution of devices based on their probability of packet loss.

To improve the communication efficiency of [35], [83] adopts the idea of SignSGD over a lossy wireless network. This is similar to the DML algorithm that was considered in 3.3.7, but rather than to enable digital AirComp it is used to increase communication efficiency. Since only one bit per element of the gradient vector needs to be transmitted, SignSGD is over an order of magnitude more communication efficient than standard 32-bit elements. This efficiency comes at a cost of representing the gradient more coarsely, which intuitively should slow down convergence. However, such intuition is not always right. In fact, SignSGD has been proven to converge with a theoretical rate similar to or in some circumstances even better than standard SGD [16]. With this SignSGD scheme, the authors of [83] attempt to minimize the outage probabilities and maximize the number of communication rounds, while maintaining an energy consumption constraint. Simulation results show that the proposed scheme can achieve both higher classification accuracy (1-3%) and lower energy consumption (10-50%) than vanilla FedAVG.

4.4.4 Total time minimization

In the previous papers, the proposed RRM schemes were greedy algorithms in the sense that they only optimized for the current communication round. Instead, [142] proposes to minimize the total time of the entire FL process, from the first communication round until convergence.

The proposed solution is a joint bandwidth allocation and client scheduling protocol which is formed by minimizing the product of the total number of communication rounds and T_{round} . The problem is solved by decomposing the problem into one client scheduling sub-problem and one bandwidth allocation sub-problem. The reason for the decomposition is that the client scheduling problem is a combinatorial optimization problem, which is infeasible to solve exactly. Experimental results on the MNIST dataset compare the scheme to [118] and show that the classification accuracy can be significantly improved.

Differently, the authors in [32] aim to minimize the total convergence time which depends on the FL parameter transmission delay per iteration and the

number of iterations that FL requires to converge. In this problem, the authors consider a user selection matrix and a resource block allocation matrix as variables, which directly impacts whether or not the users participate in the training. The authors propose a probabilistic user selection, to schedule users that have a high impact in the global FL model, and an uplink resource block allocation, given the user selection. To further reduce the total convergence time, the authors use a neural network to estimate the local FL models of users that did not receive a resource block and use these estimated models to improve the convergence speed. The numerical results indicate a reduction in the FL convergence time of 56% and improvement in the accuracy of 20% when compared to an FL algorithm that randomly determines the subset of selected users and resource blocks allocated to each user for FL parameter transmission.

4.4.5 Empirical classification error

Although the ultimate goal of these RRM algorithms is to reach the highest possible classification accuracy under communication constraints, none of the protocols maximize the accuracy directly. There still exists a gap in ML theory, which is a closed-form expression for the relationship between the number of training samples and the classification accuracy. In this survey, we have seen multiple examples of getting around this gap by using other metrics as proxies for classification accuracy. In [163] instead, the use of an empirical function is proposed to model how the accuracy depends on the sample size. The empirical function of the classification error $\Theta(N)$ with respect to the number of training samples N is designed to satisfy three properties:

- The classification error is a percentage that must lie within $0 \leq \Theta(N) \leq 1$;
- More data provides more information, and thus $\Theta(N)$ should be a monotonically decreasing function of N ;
- As N increases, the magnitude of the derivative $\partial\Theta(N)/(\partial N)$ should gradually decrease and eventually go to zero, since infinitely increasing the sample size should not improve the classification error.

Based on these properties, the function $\Theta(N) = a \cdot N^{-b}$ is chosen, where a and b are tuning parameters. The function is then trained with a limited number

of training samples, the classification accuracy is tested, and this data point is used to fit the tuning parameters. By repeating this process, samples on classification accuracy are gathered, and the parameters are found via nonlinear least-squares fitting. After fitting, the function is used as the objective for an optimization problem to find an RRM scheme. In a numerical study, the resulting RRM scheme is compared to the standard max-min fairness and sum-rate maximization protocols. When training a classifier for the MNIST dataset, the proposed scheme outperformed both baselines by a classification accuracy of about 1-2%. If the same classification accuracy is targeted, the proposed scheme saves at least 30% transmission time compared to both baselines.

4.4.6 Federated Distillation

In Section 3.3.3, we discussed a DML scheme known as FD. There, the model outputs are combined in the uplink direction via AirComp to reach exceptional communication efficiency. In this section, we are instead considering a novel Federated Distillation approach that uses digital communication. In [119], the authors consider uplink-downlink asymmetric channels, where the uplink channel capacity is more limited than the downlink. Since the uplink channels are limited, the cheap communication of model outputs in FD makes sense. However, in the more powerful downlink channel, it would be better to communicate more information than what is contained in model updates, considering that pure FD sacrifices accuracy to pay for the communication efficiency.

Therefore, [119] proposes a scheme that communicates model outputs in the uplink (as in FD) and model parameters in the downlink (as in FL). To achieve this, a method known as FL after distillation [121] is utilized. Specifically, this means that the server converts uploaded model outputs to ML model parameters, using a process known as knowledge distillation [69]. In addition to model outputs, this process requires additional training samples from the user devices, which violates user privacy. Therefore, [119] utilizes a mixup scheme to obscure the original samples, in which the idea is to create locally superpositioned samples using the mixup algorithm [188] which provides realistic synthetic samples for the knowledge distillation process without sacrificing too much privacy. In a numerical study, the proposed mixup scheme achieves 42.4x smaller payload size than FL, which leads to

significantly more communication rounds for a fixed period of time. As a result, their proposal achieves up to 16.7% higher classification accuracy than FL.

4.4.7 Batch size selection

In [129], the authors include the selection of batch size among the decision variables for the RRM. The motivation is based on the aforementioned straggler effect (see Section 4.1), which causes the slowest device to act as a bottleneck. By giving the RRM control over the batch size, this situation can be improved in two ways. The fastest devices of the network can be asked to train with a larger batch size (N^t in Eq. (2.4)), thus increasing the accuracy of their gradients without decelerating the FL process. Similarly, the batch size of the slowest devices can be decreased, sacrificing some of their performance to accelerate the FL process. The scheme improved the classification accuracy by approximately 2% compared to both random selection of batch sizes and uniform selection of batch sizes.

4.4.8 Importance-aware radio resource management

The proposal in [128] is a user selection scheme taking both channel fading and data importance into account. Similar to how [118] uses the number of scheduled users as a proxy for convergence rate, [128] uses data importance. The optimal user selection is found by the following optimization problem:

$$\begin{aligned} \min_{p_1, p_2, \dots, p_K} \quad & \sum_{k=1}^K p_k (\rho(-I_k) + (1 - \rho)T_k) \\ \text{s.t.} \quad & \sum_{k=1}^K p_k = 1, \end{aligned} \tag{4.8}$$

where p_k is the probability that k is scheduled, I_k is the importance of the gradient at device k , and T_k is the time for device k to upload its gradient. In this case, the authors use the gradient divergence as the importance measurement, i.e., $I_k = \|\nabla F_k(\mathbf{w}^t) - \nabla f(\mathbf{w}^t)\|^2$. Note that the importance is negative since we want to maximize importance but minimize latency. The solution to this problem strikes a balance between data importance and channel quality,

where the weight between the two is controlled by $\rho \in [0, 1]$. When training an MNIST classifier, the channel and importance aware user scheduler outperformed a channel-based scheduler both in convergence rate and final classification accuracy. The simulation results suggest a decrease of less than half the convergence time and an improvement of up to 2% higher in the final accuracy.

Differently from [118], [128], the authors in [171] introduce scheduling policies that use novel update importance and latency policies for client scheduling to reduce the required number of communication rounds and the total time in a communication round. The update importance policy is based on two sub-metrics: update staleness and update drift. The update staleness measures the staleness associated with the local updates of each client and aims to keep the local updates as fresh as possible. The age of update rule on client k for communication round $t + 1$ is defined as $a_k(t + 1) = (a_k(t) + 1)(1 - s_k(t))$, where $a_k(t)$ is the age of the local update in round t , and $s_k(t)$ is a binary indicator that equals 1 if client k receives the global model in round t , i.e., if the wireless channel is above a predefined threshold for the signal detection, and 0 otherwise. The update drift is based on the distance, either the Manhattan or the Euclidean distance, between the local model and the global model.

For the latency-based policy, it considers a long-term fairness constraint to allow fair participation among clients that may have important data while having a bad channel condition. The results show that the proposed scheduling policies achieve a higher accuracy than FedAvg with random scheduling and that different policies are recommended for different goals. To reduce the number of communication rounds, a scheduling using update importance metrics is recommended; whereas to shorten the total time in a communication round, a schedule using latency metrics is recommended.

In [7], the authors analyze the significance of the local models and the quality of the channels over the wireless multiple access channel from the users to the PS as user scheduling metrics. The main idea is to share the limited wireless resources with the users that have significant contributions to the model, rather than all the users. As a result, users with more significant updates can have more resources and can transmit their updates more accurately. On the other hand, users with very bad channels may not be able to communicate their updates accurately unless they are allocated a relatively significant portion of the resources; it is irrelevant whether the updates are significant or not. It is

shown numerically in [7] that considering both these metrics in user scheduling results in a better performance than considering each metric individually. The authors extend this result by deriving a convergence rate in [8] that corroborates the experimental results.

4.4.9 Energy harvesting and power transfer

One promising solution to overcome the energy limitations in IoT is energy harvesting, which allows devices to harvest radio frequency energy when communicating with a PS [38]. In FL over wireless IoT, the downlink transmission of the aggregated model parameters from the PS to the IoT devices could be used to provide energy to the devices. Hence, the use of energy harvesting for IoT devices with FL would be a perfect combination. However, how to allow the devices to harvest sufficient energy to train a FL model while not substantially increasing the communication round time is largely an open question. The use of energy harvesting for FL is highly novel and to the best of our knowledge, there are only three works in the literature [62], [148], [185].

In [185], the authors consider an FL application in which a wireless network uses power-beacons to transfer radio frequency energy. The key components of the work are the distributed gradient estimation, local-computation optimization, and optimal learning-wireless power transfer tradeoff. The distributed gradient estimation is related to the convergence of the FL method based on the mini-batch size of devices, number of active devices, and computation-outage probability. The computation-outage is an event in which a device does not harvest more energy than is necessary to transmit, thus not being able to transmit. The local-computation optimization aims to minimize the local gradient deviation present in the expected convergence rate expression, whose solution is accomplished through the optimization of the mini-batch size and processor clock frequency. Then, the authors derive an optimal learning-wireless power transfer tradeoff, which shows that a higher density of power beacons improves the learning convergence and the local gradient deviation. Moreover, it provides scaling laws of the convergence rate with respect to the transferred energy and the devices' computational capacities.

The authors in [148] analyze a multi-antenna PS using the simultaneous wireless information and power transfer technology for IoT devices. The sce-

nario considers FL simultaneously training a learning model while communicating with a PS (see Figure 4.4). The authors consider the use of FedProx [93], a recent generalization of FL that allows to optimize the number of local iterations at each device, while guaranteeing convergence to (non-)convex learning tasks. The work aims to minimize the number of communication rounds and communication round time while optimizing the number of local iterations, the time to transmit/receive, and to harvest a percentage of the total energy spent at each round and device. From the energy harvesting literature [157], maximum ratio transmission (MRT) beamforming is better at harvesting energy than zero-forcing (ZF) beamforming while ZF is better at providing higher rates than MRT due to the interference cancellation. Hence, it is non-trivial to choose the beamforming method due to a possible increase in the communication round time if the devices do not have sufficient energy to harvest or do not have sufficient rates to transmit the model parameters. For this reason, the authors consider MRT and ZF, and analyze which method is more suitable for energy harvesting within FL. The results indicate that the test accuracy using either MRT or ZF with the optimization of the local number of iterations outperforms a solution without such optimization. Moreover, it shows that MRT vastly outperforms ZF in terms of minimum communication round time for all percentages of the energy harvesting required.

4.4.10 Noisy downlink

Although the PS typically has access to more resources than the edge users, it is essential to consider imperfect transmission over wireless networks, where the PS shares the global model with the users for local training. In this case, users may not receive the global model available at the PS accurately, and the analysis of the convergence behavior of FL should account for a noisy version of the global model at the users. Digital transmission of the model over a bandwidth-limited noisy downlink leads to a relatively coarse estimation of the global model at the devices since the model vector has a high empirical variance, and quantizing the model itself does not provide an accurate estimate. Therefore, it is suggested in [26] to project the model vector linearly using a random matrix before quantizing it. This random linear projection spreads the information in the model vector more evenly across its dimensions, and leads to a smaller empirical variance. Then, the PS quantizes the projected model and

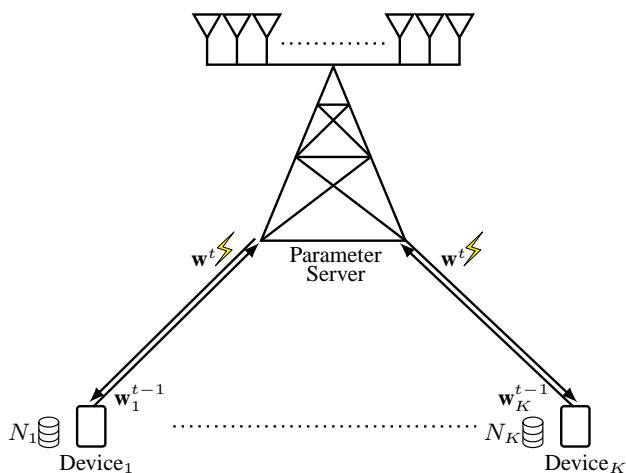


Figure 4.4: An example of a network employing FL and simultaneous wireless information and power transfer with K devices. The devices send the local model \mathbf{w}_k^{t-1} in the uplink in one time slot, while the edge server sends the global model \mathbf{w}^t and energy in the downlink on the subsequent time slot.

broadcasts the quantized vector over the downlink, where the users recover the actual model from the quantized vector having knowledge about the random matrix employed at the PS. In a follow-up work, [155] suggests to compress the model itself while accumulating and compensating the quantization error. This may lead to a coarse estimate of the model in the users if the downlink capacity is not large enough, in which case the model is compressed with a relatively low quantization level.

It is shown in [104] that the global model update, with respect to the last global model estimate available at the devices, has significantly less empirical variance than the global model. As a result, quantizing the global model update provides a more accurate estimate rather than quantizing the global model itself for the same quantization level. The authors in [104] introduce quantizing the global model updates at the PS with respect to the last model available at the users. This approach provides a significant improvement over the ones introduced in [26], [155], which is due to the availability of a more accurate estimate of the global model at the users. This approach is extended in [105] by considering broadcasting different global model descriptions to different users based on the broadcasting capacity region such that the users

with better capacities receive a more accurate estimate of the global model. This introduces a new user scheduling metric, which is based on the downlink capacity, through which at each iteration only the devices with relatively good channels, i.e., better global model estimates, can be selected to participate in the training.

It is worth highlighting that, as studied in [12], analog transmission of the model from the PS allows different devices to receive different noisy copies of the model, where less noisy devices receive a better version of the model. As such, devices with more accurate estimates of the model can compensate the lack of accurate estimates of the model at noisier devices when averaging the local models transmitted over uplink. This may lead to performance improvement compared to digital transmission of the global model from the PS [5], [12].

4.4.11 Federated meta-learning

Within certain narrow fields of ML, state-of-the-art systems are in parity with or even beyond human capabilities, such as playing the games of Chess and Go [20]. However, to reach such capabilities, state-of-the-art ML systems require significantly more exposure to data than a human. For instance, the training process of AlphaGo included approximately 600 billion moves of Go to train the value network [90]. If a human plays for 8 hours every day of their life, spending an average of 10 seconds per move, it would take the human more than 500,000 years to play 600 billion moves.

To address this efficiency gap, the field of meta learning was born [51]. In meta learning, the goal is to train a parametrized algorithm that, in turn, is used to train ML models, i.e., the parametrized algorithm is learning to learn (meta learning). Practically, the fundamental difference between meta learning and standard ML can be expressed as the division of testing and training cases. In standard ML, the dataset is divided into training data and testing data, where the training data is used to train the model and the testing data is used to evaluate its performance. In meta learning, there is instead a collection of tasks, which are divided into training tasks and testing tasks, where the tasks are generally non-overlapping, e.g., one task might be to classify different animals and another to classify plants. The idea is that the training tasks are used to train the parametrized learning algorithm, which learns to detect common

structures among the non-overlapping tasks. The testing tasks are then used to evaluate how well the learning algorithm trains ML models on the previously unseen tasks using just a few data samples.

In the space of DML, Federated Meta Learning (FML) is a recently proposed framework for achieving fast learning with distributed data [31]. The FML framework leverages the data of multiple devices to train the parametrized learning algorithm. This algorithm can then be used by the participating devices to train an ML model but more importantly, new devices can be given the parametrized learning algorithm upon joining the network so that they can quickly train an ML model using just a few data points. The underlying assumption here is that the devices carry data for a similar class of tasks, e.g., image classification, but with non-overlapping tasks within that class. In [183], the FML framework is brought into the wireless setting. First, the authors claim that the uniform selection of devices in each round (which is part of vanilla FML) leads to slow convergence rates. Then, they propose a non-uniform device selection scheme that maximizes a lower bound on the convergence speed of FML. In the same paper, the model is also extended to a joint user device selection and RRM problem. The paper contains both theoretical insights in terms of convergence bounds and numerical results that reveal strictly lower losses for the proposed system compared to a greedy and random RRM baseline.

5

Open Problems

The current literature on Wireless for ML has demonstrated that many critical metrics can be substantially improved by tailoring wireless network protocols to support ML, including latency, classification accuracy, energy consumption, and spectrum efficiency. However, the literature is still young and there are fundamental problems that remain unsolved. In this section, we give brief insights into these open problems to inspire future research.

5.1 Over-the-air Computation

CoMAC for ML is an exciting area of research since it offers a radically new way to think about wireless protocol design. However, the divergence from digital communications poses challenges of incompatibility with standard hardware and lack of prior experience. There is a need for careful investigation of assumptions in the theory and extensive testing in practice. If these challenges are overcome, substantial rewards await in the form of massively improved spectral efficiency, approximately proportional to the number of participating devices.

5.1.1 Digital over-the-air computation

As explained in Section 3.3.7, a recent study proposed a digital CoMAC protocol, based on one-bit quantization of gradient elements and BPSK modulation [196]. The proposal carries great importance for the practical implementation of CoMAC since it is compatible with the digital wireless transceivers we are using today. However, there are two potential issues with the scheme that should be investigated further.

First, BPSK demands more precise synchronization than comparable analog schemes. For example, [56] showed that analog CoMAC can be achieved with just coarse block-synchronization by encoding its real-valued message in the transmit power of a series of random signal pulses. In contrast, the BPSK-based scheme is dependent upon constructive and destructive interference of phase modulated signals to represent the transmission of "+1" and "-1". Such a scheme requires very precise alignment of the analog waveforms, which may be unreasonably difficult or expensive to achieve in practice [56].

Second, the restriction of using one-bit quantization of the gradient elements could pose problems. The numerical study in [196] found that the classification accuracy of one-bit quantization was comparable to analog communication, but this could easily change depending on the properties of the wireless network. As learning bounds on over-the-air FL demonstrate, noisy estimation of the local models slows down convergence and harms the final accuracy of the model [138], and the combination of quantization noise and channel noise can yield undesirable results.

5.1.2 Channel state information

As we have seen in Section 3.1 and [1], [43], [57], the CSI acquisition effort is greater for over-the-air computation than for digital communications. Multiple solutions have been devised to solve this issue, such as blind estimation using either MIMO or IRSs. However, there are still open questions related to channel estimation (CE) that remain unaddressed. In particular, the current literature assumes the availability of perfect CSI, which allows for perfect inversion of the channel. In reality, noisy CSI will lead to distorted sums. Instead of the channel inversion in Eq. (3.2), the received vector will be

$$\sum_{k=1}^K h_k z_k + v = \sum_{k=1}^K \frac{w_k h_k}{\hat{h}_k} + v, \quad (5.1)$$

where the estimated channels \hat{h}_k do not cancel out h_k . To understand the effect of imperfect CSI on learning performance, this needs to be studied. Additionally, the performance comparisons of CoMAC and digital communications have not considered the cost of CSI acquisition, which could be a non-negligible difference due to the increased channel estimation effort.

5.1.3 Security

A fundamental consequence of CoMAC, is that it is impossible to see who is transmitting model updates in the uplink. This can be seen as a blessing or a curse. The upside is that user privacy is guaranteed, stopping the potential for model inversion at the PS [52]. The downside is that it opens up for potential adversaries to corrupt the training process. Because of the inherent anonymity of CoMAC, it is easy for an adversary to send malicious model updates and harm training. This process is known as model poisoning and has received some attention from the FL community [13], [18], [53]. However, the defense strategies proposed in the literature depend on detecting anomalies in individual model updates, which is impossible for CoMAC. Hence, there is a need to find new strategies against model poisoning that work without seeing individual model updates. One possible countermeasure is the consideration of coded computing, but so far there is only one paper that would be applicable to CoMAC [147]. Another idea is briefly mentioned in [199] where all legitimate devices are assigned a common secret spreading code. Consequently, the PS can exclude devices that are not using the secret code. However, despite these initial steps, the security problem of over-the-air FL is far from solved.

5.1.4 Self-aware power control

The power control schemes developed for CoMAC are all reliant upon an assumption of random messages being transmitted by the devices [28], [29], [103], [184]. Such an assumption is made to reflect that the transmitting devices are unaware of the messages to be sent by other devices in the network. However, for mathematical simplicity, these schemes are not only assuming that other devices' messages are unknown but also the message of the transmitting device itself. In practice, each device of course knows the message it is about to transmit, therefore there is room to improve the power control by taking this information into account. Since these schemes use analog modulation,

the strength of the transmitted signal depends on the value being sent, and therefore the knowledge of this value should change the optimal transmission power.

5.2 Digital Communications

Today's digital communication systems are optimized for communication metrics such as data rate, packet error rate, latency, or fairness. These metrics are in some way beneficial for the goals of ML but are not completely aligned. Instead, digital Wireless for ML systems should optimize metrics such as classification accuracy, data importance, or training time. In the current Wireless for ML literature, we have seen that customized retransmission and RRM protocols generate significantly better ML models than generic communication protocols. However, since machine learning performance is difficult to predict ahead of training, it is not clear what the correct objective of these protocols should be, leaving us with proxies for classification accuracy, such as data importance, user participation, or bounds on the learning loss. A deeper theory of these objectives and the interplay between communication and learning is needed.

5.2.1 Data-importance metrics

In most Wireless for ML scenarios, the acquisition of data from user devices is the bottleneck of training. Therefore, the selection of which data points to collect or which devices to schedule is of critical importance to efficiently train an ML model. In much of the current literature [75], [99], [100], [128], [145], this selection is based on data-importance metrics from the field of Active Learning. The original problem studied in Active Learning was that of labeling data samples but there are important differences between the problem of labeling data samples and communicating them, which opens up new research directions. Specifically, we have listed two such differences below:

- In most Wireless for ML scenarios, the labels are available at the user devices. By using importance metrics from Active Learning as-is, potentially valuable information (the labels) is completely unutilized. This calls for the investigation of new importance metrics that incorporate the label;

- In DML, the devices do not communicate data samples but local models, model updates, or gradients. However, for the sake of device scheduling, we are still interested in the importance of the updates. In one paper, the local loss was proposed as a measure of gradient-importance [55] but no more work has been done in this direction. This measure could potentially be used to improve RRM for ML and other metrics for gradient-importance could be developed.

5.2.2 Data-importance staleness

In several importance-aware RRM schemes, the data-importance is not updated in every communication round. For instance in [55], the data importance is measured on a user basis and is calculated locally during training to be transmitted in conjunction with the local model on the uplink. However, only a subset of users is selected for any given round, leaving the PS with a mix of old and fresh data importance measurements. As the global model is trained, the importance of a user's data could change substantially. This calls for further studies on the effect of data importance staleness on learning convergence, and eventually solutions to combat this effect.

5.2.3 Channel uncertainty

Despite the strong progress on developing RRM schemes for FL, there are still fundamental questions that are unanswered. One example is the impact of channel uncertainty on the learning convergence. In practical systems, the RRM decisions will always be based on an imperfect estimate of the wireless channel and the impact of this uncertainty on these systems is still unexplored. Despite affecting the RRM decision, channel uncertainty will also have an impact on the packet error rates, which will thus worsen the optimality gap [35]. A recent study [161] has taken a first step to address imperfect CSI but more work is needed.

5.2.4 Energy harvesting for federated learning

With the increasing use of IoT devices for monitoring applications, the importance of energy harvesting for FL is quickly increasing. The works we discussed [62], [148], [185] are the first attempts to analyze this emerging

field, but substantial work is still necessary. Specifically, the impact of its application with bandwidth limited transmissions, such as narrowband IoT which limits the transmission rate for the IoT devices. Moreover, the impact of CSI errors in the process also needs to be considered given that the errors will impact the learning accuracy and may imply the need for retransmissions. If retransmissions are needed, this may also be beneficial for the energy harvesting of the devices, but will impact the ultimate convergence time of the process. Hence, there is a tradeoff in terms of retransmissions, in case of CSI errors, energy harvesting, and learning accuracy.

5.3 Problems Relevant to Analog and Digital Communications

An important missing piece of analytical performance evaluation of FL over wireless networks is its gap to centralized learning, where the entire data is available at a single server carrying out all the processing. FL over wireless networks suffers from unreliable communications between the nodes in addition to the various heterogeneity aspects that exist with the FL framework. This gap should capture the impact of various factors that exist with the FL framework due to its distributed nature and communications over noisy channels. It would be particularly interesting to analyze the impact of noisy communications on the performance gap to centralized learning.

6

Applications

The term Wireless for ML is meant to capture any wireless technology tailored to solve a machine learning problem, including model training, data collection, and inference. However, the current Wireless for ML literature is almost exclusively focused on supervised learning using a distributed data set. Therefore, the work we have surveyed in this monograph applies to any application that falls within that domain, given that the data-collecting devices are connected via a wireless link. There are already a number of such applications envisioned or used in practice, such as Vehicular Internet of Things [45], FL for wireless [116], environmental monitoring [122], mobile keyboard prediction [63], and Industrial IoT [115]. Besides the current applications, Wireless for ML argues for the creation of foundations of an infrastructure for DML. Such an infrastructure will be able to support many upcoming applications that we cannot envision today. In this section we expose a few current applications to discuss the challenges they pose and how Wireless for ML addresses those challenges.

6.1 Smart City

The future smart cities critically depend on the reliable monitoring of large civil infrastructures such as roads, tunnels, bridges, water networks, renewable energy sources, or electric grids. The denser we can measure relevant

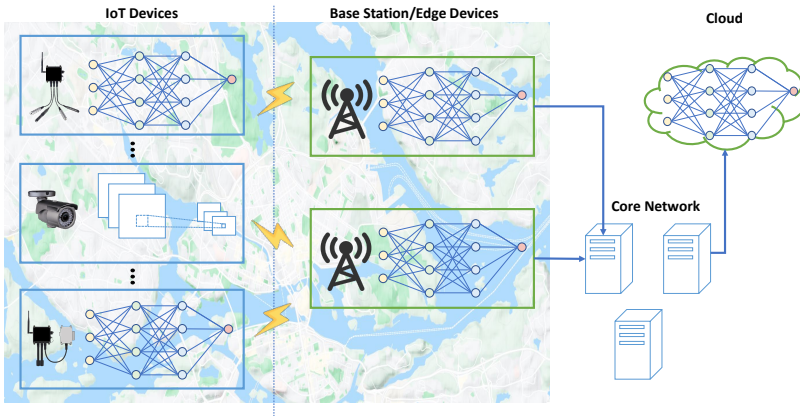


Figure 6.1: An example of wireless IoT for ML monitoring in smart cities, including IoT devices for water monitoring, security surveillance, and mobile monitoring.

information in space and time, the higher is the potential to perform an accurate monitoring. Recently, IoT is becoming instrumental to performing such fine-grained monitoring and is opening the potential for several new monitoring services. Although IoT devices can collect a large amount of data, it is challenging to have sustainable, secure, and reliable monitoring services. To overcome such challenges, a key promising solution is the use of ML over the IoT devices in a distributed manner across the wireless network, as illustrated in Figure 6.1.

Using data-driven and model-based solutions to perform reliable data analysis, it is possible to establish a methodology for scalable, resource-efficient learning and decision making under physical, communication, and security constraints. With the increase in the computation capacity of sensors, it is now possible to consider a scenario in which the IoT devices perform part of the learning and/or prediction tasks locally, rather than offloading to the cloud or edge server. Using DML across the wireless network, the IoT devices may reduce the need to transmit a large amount of data to the network, alleviate the storage and energy consumption due to less intensive transmission needs, and enhance privacy by not transmitting the raw data over the network. For example, the authors in [44] propose model compression for IoT devices monitoring water conditions in Sweden. The proposed model compression shows a degradation of 2.5% in test accuracy while saving 96% in transmissions compared to a scheme that sends all raw data.

Many Smart City IoT nodes will be placed in inaccessible or remote locations, such as chimneys, water pipes, lakes, and underground. As such, there is a large cost associated with performing maintenance on these devices, including charging or replacing the battery. The results from Section 4.4.2 suggest that RRM for energy-efficient learning can significantly prolong the battery life of such devices. While Section 4.4.9 suggests that energy harvesting can be leveraged to completely compensate for the consumed energy by increasing the communication round time. Therefore, the use of Wireless for ML can help to learn and predict relevant phenomena in critical infrastructures of smart cities, such as water leakage in water distribution networks and structural problems in the road infrastructure.

6.2 Vehicular Communication

To enable intelligent transportation systems, such as autonomous driving and advanced driver assistance systems, it is necessary to integrate vehicular communications and machine learning. Vehicular communications provide communications between vehicles, pedestrians, road infrastructure, and the Internet, and has severe requirements in terms of low latency, high reliability and high rates [45]. Due to the advantages of FL in terms of distributed computation, communication efficiency, and privacy by not sending raw data, its use in vehicular communications has started to gain momentum (see Figure 6.2).

Some studies have recently considered FL methods in learning tasks at vehicles, such as collision avoidance and traffic sign recognition, which can be considered as FL in vehicular applications but without tailoring wireless methods for ML. Specifically, the authors in [46] investigate FL applications for vehicular communications, including autonomous driving, road safety prediction, and vehicular object detection, and highlight some of the challenges and research directions for FL in vehicular communications. Conversely, FL methods have been applied to resource management problems in vehicular communications, such as power control. For instance, the authors in [134] address wireless resource management problems in vehicular communications by using FL to estimate the tail distribution of the network-wide queue lengths.

However, we are interested in this survey in the joint design of vehicular communications and FL, or RRM in vehicular communications for FL, in

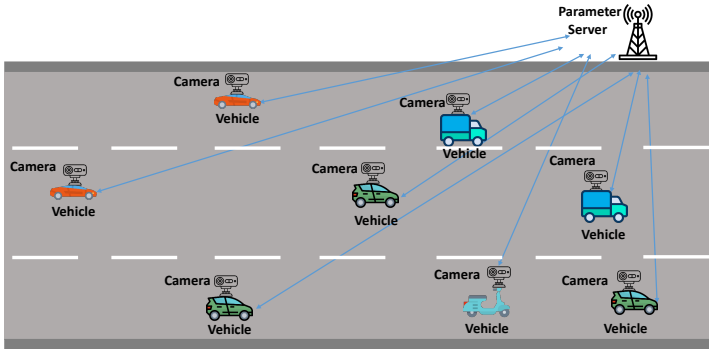


Figure 6.2: An example of vehicular networks using distributed machine learning for training. The vehicles exchange their learning models with a parameter server towards a global common goal, such as traffic sign recognition.

which a PS and vehicles jointly optimize their learning goals together with the communication requirements. These three directions are illustrated in a recent survey [45], in which the authors discuss mainly the communication and learning aspects, while briefly mentioning the challenges of joint learning and communication of FL in vehicular communications.

To the best of our knowledge, there is only one study that fits our criteria [187]. The authors in [187] consider the problem of learning and optimizing their autonomous controller design, which allows the vehicle to execute near real-time decisions, in the presence of wireless uncertainties and environmental dynamics. To this end, this work proposes a dynamic federated proximal algorithm to account for the varying participation of vehicles due to mobility and wireless channels. To improve the convergence of the proposed FL algorithm, the authors design an incentive mechanism for device participation using contract theory. The incentive mechanism acts as a device participation and importance RRM, such as the ones in Sections 4.4.1 and 4.4.8, by taking into account the data quality and devising a power allocation mechanism to maximize the convergence gain between two consecutive rounds. The results show substantial improvements in the convergence speed compared to the other FL algorithms, such as FedAvg, and baselines of their own proposed FL algorithm using maximum and random power allocations.

The work in [187] jointly analyzed some of the control and learning challenges, but the communication challenges are still open. Specifically, the

impact of severe requirements in terms of low latency, high reliability, and high rates in a joint communication and learning approach needs to be considered. Moreover, the impact of quick channel variations need to be analyzed together with the learning convergence of the FL method. Therefore, research for this application is still quite open and there are many challenges ahead.

6.3 Augmented and Virtual Reality

For augmented and virtual reality (AR) and (VR) services provided by wireless networks, any sudden drop in the data rate or increase in the delay can negatively affect the quality of experience (QoE) of VR users. Although 5G networks support operation at high frequency bands as well as flexible frame structure to minimize latency, the performance of communication links at high frequencies is highly prone to blockage thus reducing the QoE of VR users.

One key application of using FL for improving QoE of wireless VR users is presented in [33]. In the considered model, each BS serves several VR users. Each user will transmit tracking information to the BS. Then, the BS will generate VR images according to the received tracking information and transmit the generated VR images to the VR users over millimeter wave frequency, which can be seen in Figure 6.3. Since VR images are transmitted over millimeter wave links, user movement such as mobility and orientation will introduce blockages to the millimeter wave transmission links thus decreasing the QoE of VR users.

The goal of [33] is to minimize the breaks in presence (BIP) of all VR users via optimizing user association. Since user association depends on the user mobility patterns and orientation, it is necessary to design a novel learning method to analyze the mobility patterns and orientation of each VR user. Meanwhile, since user association changes over time, each user may connect to different BSs at different time slots and hence each BS can collect partial information related to user mobility patterns and orientation. Thus, traditional centralized learning algorithms that are implemented by a given BS cannot predict the entire VR user's locations and orientations without knowing the user's data collected by other BSs. To minimize the BIP of all users, an echo state network (ESN) based FL algorithm is designed, which enables the BSs to collaboratively generate a global ESN model to predict the whole set of

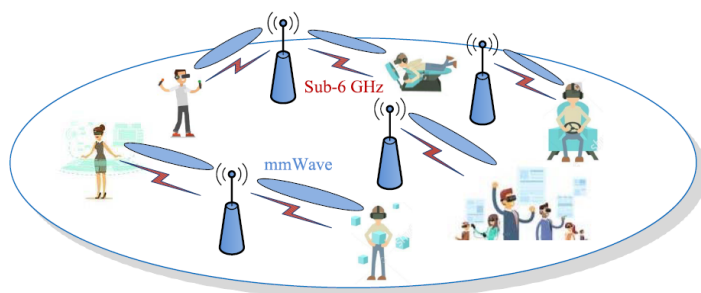


Figure 6.3: The architecture of a wireless virtual reality network. Virtual reality applications impose stringent throughput requirements. Therefore, millimeter wave communication is employed in the downlink.

locations and orientations for each user without transmitting the collected data to other BSs. Meanwhile, different from traditional FL algorithms that need to transmit the entire FL model, ESN based FL only needs to transmit the parameters of the output layer which can significantly reduce the size of data transmitted over wireless links thus improving convergence speed. In many envisioned VR and AR applications, co-located users share the same virtual world, for example in Smart Campus [177] and the Metaverse [117]. When many co-located users share an ML task, over-the-air FL offers radical communication-efficiency improvements over orthogonal communications, as discussed in Section 3. Therefore, Wireless for ML can assist in meeting the heavy communication demands imposed by AR and VR.

6.4 Edge Caching

Caching of popular content at the network edge has been introduced as a promising approach to push the network traffic closer to the edge and reduce data traffic on backhaul networks [36], [108]. Popular content is stored close to the edge terminals, at small BSs, APs, or edge devices, proactively, such that it can be accessed more easily by the edge users. This is particularly appealing for applications with stringent delay and bandwidth requirements. One of the challenges in edge caching is determining popularity of the content which is stored in the cache memories. Static and dynamic models have been introduced to capture the content popularity, where static models do not consider the time varying nature of the real-time content. On the other hand, dynamic models

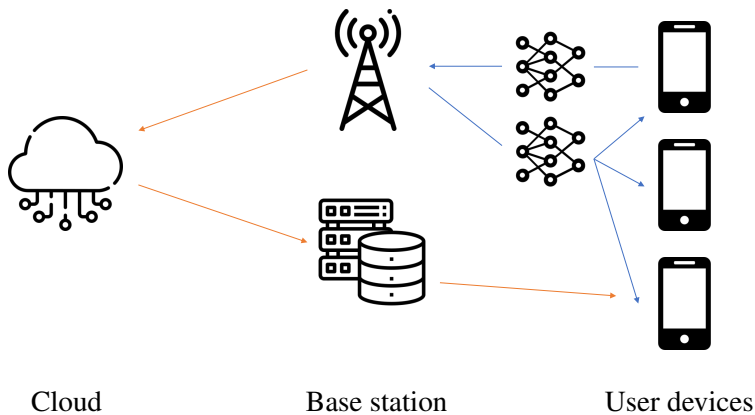


Figure 6.4: Illustration of federated learning for content popularity prediction. The blue lines illustrate how the user devices are collaboratively training the prediction model, using the base station as a PS. The red lines illustrate the communication of content. The base station requests content based on the output of the popularity prediction model.

require accessing data for content differentiation. This is not desirable in wireless systems since sharing data with other nodes may violate the privacy of users.

A distributed ML framework is a perfect fit to learn content popularity for edge caching by utilizing processing capabilities of edge devices. In this approach, local data at the users can be used to train a global model that is shared with all the users in order to learn the content popularity; see Figure 6.4. Therefore, the entire data across the network is used to determine the popularity of the content while data never leaves the users. The popular content is then stored close to the users to reduce the network traffic. For example, in augmented reality local data at the users can be used to learn popular elements, and the information about these elements can be cached proactively close to the users to reduce the latency and improve users experience. Furthermore, in an autonomous driving example, information about the traffic, which can be learned collaboratively using the data collected by different vehicles, can be pre-fetched into the roadside units.

Since the BS is often both the arbiter of RRM decisions and the host of the cache, it is natural to consider RRM tailored to learn content popularity, as discussed in Section 4.4. Such dedicated wireless methods could improve

the communication efficiency of training the content prediction model as well as reduce training and communication energy costs. Since trends in popular content changes regularly at a moment's notice, the prediction model should be retrained continuously, which further emphasizes the importance of communication and energy efficiency.

6.5 Unmanned Aerial Vehicles

The low-altitude airspace of contemporary cities is generally empty or dominated by urban wildlife. In the upcoming decades, this underutilized real estate is predicted to be populated by search-and-rescue drones, delivery vehicles, and aerial BSs [24], [65], [123]. These applications are enabled by the Unmanned Aerial Vehicle (UAV) technology, which provides cheap, easy to deploy, and highly maneuverable drones. However, there are many communication challenges associated with flying devices. First, there are stringent energy constraints as the weight of the battery increases the cost of flying. Secondly, the UAV air-to-ground channel is more susceptible to fading, path loss, and delay spread because of the 3D movement of the vehicles [24]. Finally, UAVs are never completely still, generating continual fast-fading.

One interesting use-case of UAVs are the deployment of flying BSs, especially in geographical zones lacking cellular infrastructure or as temporary deployment to increase cellular capacity during large events. Unlike a traditional BS, these would be able to dynamically adjust their location to improve channel quality. The prediction of the correct location is a challenging problem that depends on the propagation environment, the number of users, and their mobility patterns. FL is a natural choice for training such a prediction model using the distributed data collected by the UAV BS and mobile devices [24]; see Figure 6.5. In this case, the training data is channel state information collected by the UAV BSs. As such the data distribution changes quickly, and it is important to retrain continuously, which calls for efficient wireless protocols. Since the channels are changing quickly, the communication method must offer low latency to cope with the short channel coherence time. The over-the-air computation methods discussed in Section 3 can offer low latencies that scale inversely with the number of users, which is ideal for a flying BS deployed to a large event. Additionally, the blind methods discussed in Section 3.4.1 offer CSI-free over-the-air computation which is helpful when the channels are changing quickly.

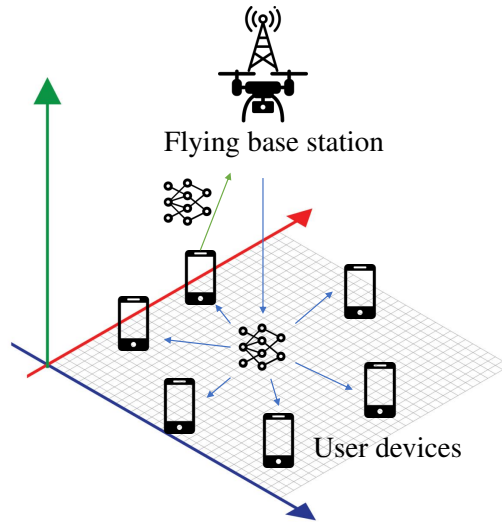


Figure 6.5: Illustration of federated learning for flying base stations. The goal is to predict the optimal 3D flying base station location that provides maximum channel gains to the user devices. Unlike a traditional base station, the flying base station can adjust its location dynamically.

Another critical application of UAVs is search and rescue missions at disaster locations. In these missions, the terrain is often unknown, since disasters such as floods, explosions, and earthquakes can change the known map completely [65]. The time to locate victims is critical since survival is often heavily dependent on quick retrieval. Unfortunately, the cellular infrastructure easily gets destroyed by the disaster, leaving rescue workers in an unknown environment, with strict time constraints, and without connectivity. UAVs could be helpful in these scenarios to quickly set up multi-hop ad-hoc networks as a replacement for the damaged cellular infrastructure and to map out the environment. However, the highly mobile environment results in uncertain channel conditions that make routing difficult. A possible solution is an ML-based model to predict the channels of potential next-hop nodes [24]. The inference of such models would be used to dynamically update the UAVs routing tables. Additionally, by training with the rescue team's devices, the UAVs can predict areas of poor coverage and adjust their locations to compensate.

Conclusions

Given the continuous growth in the numbers of IoT and mobile devices, the demand for ML over wireless networks is expected to grow significantly. However, traditional communication protocols have been shown to be greatly inefficient for carrying ML related data, creating a demand for new wireless solutions. In this survey, we have reviewed the most important contributions in this area, specifically focusing on analog over-the-air computation and digital RRM for DML.

Analog over-the-air computation offers the most radical improvements in communication efficiency, exhibiting a throughput improvement approximately proportional to the number of participating devices. However, for contemporary communication, digital transmission is the de-facto standard. Within digital RRM for DML, significant performance improvements are achieved by considering data-importance and tailored RRM protocols for FL. However, this field is still in its infancy and several fundamental problems remain. For analog over-the-air computation, the main concerns are with integration into contemporary wireless infrastructure and functionality in dynamic wireless environments. Within digital RRM for DML, there are still many open questions relating to data-importance, such as choice of metrics and staleness of importance updates.

It is highly relevant to find answers to these open questions, since efficient Wireless for ML solutions could have profound effects on society, which we have demonstrated by discussing five application areas: Smart City, Vehicular Communication, Virtual Reality, Edge Caching, and Unmanned Aerial Vehicles. The development of wireless methods specifically for ML is a fertile area of research that could provide significant benefits in terms of energy efficiency, spectral efficiency, and latency.

References

- [1] O. Abari, H. Rahul, and D. Katabi, “Over-the-Air Function Computation in Sensor Networks,” *arXiv abs/1612.02307*, 2016.
- [2] O. Abari, H. Rahul, D. Katabi, and M. Pant, “Airshare: Distributed Coherent Transmission Made Seamless,” in *Proceedings of the 2015 IEEE Conference on Computer Communications (INFOCOM)*, IEEE, pp. 1742–1750, 2015.
- [3] S. Abdulrahman, H. Tout, H. Ould-Slimane, A. Mourad, C. Talhi, and M. Guizani, “A Survey on Federated Learning: The Journey from Centralized to Distributed On-Site Learning and Beyond,” *IEEE Internet of Things Journal*, vol. 8, no. 7, 2020, pp. 5476–5497.
- [4] J.-H. Ahn, O. Simeone, and J. Kang, “Wireless Federated Distillation for Distributed Edge Learning with Heterogeneous Data,” in *Proceedings of the 2019 IEEE 30th Annual International Symposium on Personal, Indoor and Mobile Radio Communications (PIMRC)*, IEEE, pp. 1–6, 2019.
- [5] J.-H. Ahn, O. Simeone, and J. Kang, “Cooperative Learning via Federated Distillation over Fading Channels,” in *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 8856–8860, Barcelona, Spain, 2020.
- [6] M. M. Amiri, T. M. Duman, D. Gündüz, S. R. Kulkarni, and H. V. Poor, “Blind Federated Edge Learning,” *IEEE Transactions on Wireless Communications*, 2021.

- [7] M. M. Amiri, D. Gündüz, S. R. Kulkarni, and H. V. Poor, "Update Aware Device Scheduling for Federated Learning at the Wireless Edge," in *Proceedings of the IEEE International Symposium on Information Theory*, pp. 2598–2603, Los Angeles, CA, USA, 2020.
- [8] M. M. Amiri, D. Gündüz, S. R. Kulkarni, and H. V. Poor, "Convergence of Update Aware Device Scheduling for Federated Learning at the Wireless Edge," *IEEE Transactions on Wireless Communications*, vol. 20, no. 6, 2021, pp. 3643–3658.
- [9] M. M. Amiri and D. Gunduz, "Machine Learning at the Wireless Edge: Distributed Stochastic Gradient Descent Over-the-Air," *IEEE Transactions on Signal Processing*, vol. 68, 2020, pp. 2155–2169.
- [10] M. M. Amiri and D. Gündüz, "Over-the-Air Machine Learning at the Wireless Edge," in *Proceedings of the IEEE 20th International Workshop on Signal Processing Advances in Wireless Communications (SPAWC)*, IEEE, pp. 1–5, 2019.
- [11] M. M. Amiri and D. Gündüz, "Federated Learning over Wireless Fading Channels," *IEEE Transactions on Wireless Communications*, vol. 19, no. 5, 2020, pp. 3546–3557.
- [12] M. M. Amiri, D. Gündüz, S. R. Kulkarni, and H. V. Poor, "Convergence of Federated Learning over a Noisy Downlink," *IEEE Transactions on Wireless Communications*, 2021.
- [13] E. Bagdasaryan, A. Veit, Y. Hua, D. Estrin, and V. Shmatikov, "How to Backdoor Federated Learning," in *Proceedings of the International Conference on Artificial Intelligence and Statistics*, PMLR, pp. 2938–2948, 2020.
- [14] T. Ben-Nun and T. Hoefler, "Demystifying Parallel and Distributed Deep Learning: An In-Depth Concurrency Analysis," *ACM Computing Surveys*, vol. 52, no. 4, 2019.
- [15] M. Bennis, M. Debbah, and H. V. Poor, "Ultrareliable and Low-Latency Wireless Communication: Tail, Risk, and Scale," *Proceedings of the IEEE*, vol. 106, no. 10, 2018, pp. 1834–1853.
- [16] J. Bernstein, Y.-X. Wang, K. Azizzadenesheli, and A. Anandkumar, "signSGD: Compressed Optimisation for Non-Convex Problems," in *Proceedings of the International Conference on Machine Learning*, PMLR, pp. 560–569, 2018.

- [17] T. Bertin-Mahieux, D. P. Ellis, B. Whitman, and P. Lamere, "The Million Song Dataset," in *Proceedings of the 11th International Conference on Music Information Retrieval (ISMIR)*, 2011.
- [18] A. N. Bhagoji, S. Chakraborty, P. Mittal, and S. Calo, "Analyzing Federated Learning Through an Adversarial Lens," in *Proceedings of the International Conference on Machine Learning*, PMLR, pp. 634–643, 2019.
- [19] P. Blanchard, E. M. El Mhamdi, R. Guerraoui, and J. Stainer, "Machine Learning with Adversaries: Byzantine Tolerant Gradient Descent," *Advances in Neural Information Processing Systems*, vol. 30, 2017.
- [20] P. Bory, "Deep New: The Shifting Narratives of Artificial intelligence from Deep Blue to AlphaGo," *Convergence: The International Journal of Research into New Media Technologies*, vol. 25, no. 4, 2019, pp. 627–642.
- [21] L. Bottou, F. E. Curtis, and J. Nocedal, "Optimization Methods for Large-Scale Machine Learning," *SIAM Review*, vol. 60, no. 2, 2018, pp. 223–311.
- [22] S. Boyd, N. Parikh, E. Chu, B. Peleato, J. Eckstein, *et al.*, "Distributed Optimization and Statistical Learning via the Alternating Direction Method of Multipliers," *Foundations and Trends® in Machine Learning*, vol. 3, no. 1, 2011, pp. 1–122.
- [23] P. Branco, L. Torgo, and R. P. Ribeiro, "A Survey of Predictive Modeling on Imbalanced Domains," *ACM Computing Surveys (CSUR)*, vol. 49, no. 2, 2016, pp. 1–50.
- [24] B. Brik, A. Ksentini, and M. Bouaziz, "Federated Learning for UAVs-enabled Wireless Networks: Use Cases, Challenges, and Open Problems," *IEEE Access*, vol. 8, 2020, pp. 53 841–53 849.
- [25] R. W. Broadley, J. Klenk, S. B. Thies, L. P. Kenney, and M. H. Granat, "Methods for the Real-World Evaluation of Fall Detection Technology: A Scoping Review," *Sensors*, vol. 18, no. 7, 2018, p. 2060.
- [26] S. Caldas, J. Konecny, H. B. McMahan, and A. Talwalkar, "Expanding the Reach of Federated Learning by Reducing Client Resource Requirements," *arXiv abs/1812.07210*, 2019.
- [27] X. Cao, G. Zhu, J. Xu, and S. Cui, "Transmission Power Control for Over-the-Air Federated Averaging at Network Edge," *IEEE Journal on Selected Areas in Communications*, 2022, pp. 1571–1586.

- [28] X. Cao, G. Zhu, J. Xu, and S. Cui, "Optimized Power Control for Over-the-Air Federated Edge Learning," in *Proceedings of the 2021 IEEE International Conference on Communications*, IEEE, pp. 1–6.
- [29] X. Cao, G. Zhu, J. Xu, and K. Huang, "Optimized Power Control for Over-the-Air Computation in Fading Channels," *IEEE Transactions on Wireless Communications*, vol. 19, no. 11, 2020, pp. 7498–7513.
- [30] N. Carlini, C. Liu, Ú. Erlingsson, J. Kos, and D. Song, "The Secret Sharer: Evaluating and Testing Unintended Memorization in Neural Networks," in *Proceedings of the 28th {USENIX} Security Symposium ({USENIX} Security 19)*, pp. 267–284, 2019.
- [31] F. Chen, M. Luo, Z. Dong, Z. Li, and X. He, "Federated Meta-Learning with Fast Convergence and Efficient Communication," *abs/1802.07876*, 2018.
- [32] M. Chen, H. V. Poor, W. Saad, and S. Cui, "Convergence Time Optimization for Federated Learning Over Wireless Networks," *IEEE Transactions on Wireless Communications*, vol. 20, no. 4, 2021, pp. 2457–2471.
- [33] M. Chen, O. Semiari, W. Saad, X. Liu, and C. Yin, "Federated Echo State Learning for Minimizing Breaks in Presence in Wireless Virtual Reality Networks," *IEEE Transactions on Wireless Communications*, vol. 19, no. 1, 2020, pp. 177–191.
- [34] M. Chen, N. Shlezinger, H. V. Poor, Y. C. Eldar, and S. Cui, "Communication-efficient Federated Learning," *Proceedings of the National Academy of Sciences*, vol. 118, no. 17, 2021.
- [35] M. Chen, Z. Yang, W. Saad, C. Yin, H. V. Poor, and S. Cui, "A Joint Learning and Communications Framework for Federated Learning over Wireless Networks," *IEEE Transactions on Wireless Communications*, vol. 20, no. 1, 2020, pp. 269–283.
- [36] W. Chen and H. V. Poor, *Edge Caching for Mobile Networks*. London, UK: IET Press, 2021.
- [37] Y. Cheng, D. Wang, P. Zhou, and T. Zhang, "A Survey of Model Compression and Acceleration for Deep Neural Networks," *arXiv abs/1710.09282*, 2017.

- [38] B. Clerckx, K. Huang, L. R. Varshney, S. Ulukus, and M.-S. Alouini, “Wireless Power Transfer for Future Networks: Signal Processing, Machine Learning, Computing, and Sensing,” *arXiv abs/2101.04810*, 2021.
- [39] R. Collobert, J. Weston, L. Bottou, M. Karlen, K. Kavukcuoglu, and P. Kuksa, “Natural Language Processing (almost) from Scratch,” *Journal of Machine Learning Research*, vol. 12, 2011, pp. 2493–2537.
- [40] S. Deng, H. Zhao, W. Fang, J. Yin, S. Dustdar, and A. Y. Zo5a, “Edge Intelligence: the Confluence of Edge Computing and Artificial Intelligence,” *IEEE Internet of Things Journal*, 2020.
- [41] J. Ding, R. Calderbank, and V. Tarokh, “Gradient Information for Representation and Modeling,” *Advances in Neural Information Processing Systems*, vol. 32, 2019, pp. 2396–2405.
- [42] C. T. Dinh, N. H. Tran, M. N. H. Nguyen, C. S. Hong, W. Bao, A. Y. Zo5a, and V. Gramoli, “Federated Learning Over Wireless Networks: Convergence Analysis and Resource Allocation,” *IEEE/ACM Transactions on Networking*, vol. 29, no. 1, 2021, pp. 398–409.
- [43] J. Dong, Y. Shi, and Z. Ding, “Blind over-the-air Computation and Data Fusion via Provable Wirtinger Flow,” *IEEE Transactions on Signal Processing*, vol. 68, 2020, pp. 1136–1151.
- [44] R. Du, S. Magnusson, and C. Fischione, “The Internet of Things as a Deep Neural Network,” *IEEE Communications Magazine*, vol. 58, no. 9, 2020, pp. 20–25.
- [45] Z. Du, C. Wu, T. Yoshinaga, K. .-. A. Yau, Y. Ji, and J. Li, “Federated Learning for Vehicular Internet of Things: Recent Advances and Open Issues,” *IEEE Open Journal of the Computer Society*, 2020, pp. 45–61.
- [46] A. M. Elbir, B. Soner, and S. Coleri, “Federated Learning in Vehicular Networks,” *arXiv abs/2006.01412*, 2020, URL: <http://arxiv.org/abs/2006.01412>.
- [47] A. Elgabli, J. Park, C. B. Issaid, and M. Bennis, “Harnessing Wireless Channels for Scalable and Privacy-Preserving Federated Learning,” *IEEE Transactions on Communications*, 2021.
- [48] K. Elkhailil, A. Hasan, J. Ding, S. Farsiu, and V. Tarokh, “Fisher Auto-Encoders,” in *International Conference on Artificial Intelligence and Statistics*, PMLR, pp. 352–360, 2021.

- [49] *Ericsson Mobility Report*, 2019, URL: <https://www.ericsson.com/en/press-releases/2019/6/ericsson-mobility-report-5g-uptake-even-faster-than-expected>.
- [50] D. Fan, X. Yuan, and Y.-J. A. Zhang, “Temporal-Structure-Assisted Gradient Aggregation for Over-the-Air Federated Edge Learning,” *arXiv abs/2103.02270*, 2021.
- [51] C. Finn, P. Abbeel, and S. Levine, “Model-Agnostic Meta-Learning for Fast Adaptation of Deep Networks,” in *International conference on machine learning*, PMLR, pp. 1126–1135, 2017.
- [52] M. Fredrikson, S. Jha, and T. Ristenpart, “Model Inversion Attacks That Exploit Confidence Information and Basic Countermeasures,” in *Proceedings of the 22nd ACM SIGSAC Conference on Computer and Communications Security*, pp. 1322–1333, 2015.
- [53] C. Fung, C. J. Yoon, and I. Beschastnikh, “Mitigating Sybils in Federated Learning Poisoning,” *arXiv abs/1808.04866*, 2018.
- [54] T. Gafni, N. Shlezinger, K. Cohen, Y. C. Eldar, and H. V. Poor, “Federated Learning: A Signal Processing Perspective,” *arXiv abs/2103.17150*, 2021.
- [55] J. Goetz, K. Malik, D. Bui, S. Moon, H. Liu, and A. Kumar, “Active Federated Learning,” *arXiv abs/1909.12641*, 2019.
- [56] M. Goldenbaum and S. Stanczak, “Robust Analog Function Computation via Wireless Multiple-Access Channels,” *IEEE Transactions on Communications*, vol. 61, no. 9, 2013, pp. 3863–3877.
- [57] M. Goldenbaum and S. Stanczak, “On the Channel Estimation Effort for Analog Computation over Wireless Multiple-Access Channels,” *IEEE Wireless Communications Letters*, vol. 3, no. 3, 2014, pp. 261–264.
- [58] A. J. Goldsmith and S.-G. Chua, “Adaptive Coded Modulation for Fading Channels,” *IEEE Transactions on Communications*, vol. 46, no. 5, 1998, pp. 595–602.
- [59] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. MIT Press, 2016.
- [60] W. Guan, H. Zhang, and V. C. Leung, “Customized Slicing for 6G: Enforcing Artificial Intelligence on Resource Management,” *IEEE Network*, 2021, pp. 264–271.

- [61] S. Ha, J. Zhang, O. Simeone, and J. Kang, “Coded Federated Computing in Wireless Networks With Straggling Devices and Imperfect CSI,” in *Proceedings of the IEEE International Symposium on Information Theory (ISIT)*, IEEE, pp. 2649–2653, 2019.
- [62] R. Hamdi, M. Chen, A. B. Said, M. Qaraqe, and H. V. Poor, “Federated Learning Over Energy Harvesting Wireless Networks,” *IEEE Internet of Things Journal*, vol. 9, no. 1, 2022, pp. 93–103.
- [63] A. Hard, K. Rao, R. Mathews, S. Ramaswamy, F. Beaufays, S. Augenstein, H. Eichner, C. Kiddon, and D. Ramage, “Federated Learning for Mobile Keyboard Prediction,” *arXiv abs/1811.03604*, 2018.
- [64] T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*, ser. Springer series in statistics. Springer, 2009.
- [65] S. Hayat, E. Yanmaz, and R. Muzaffar, “Survey on Unmanned Aerial Vehicle Networks for Civil Applications: A Communications Viewpoint,” *IEEE Communications Surveys & Tutorials*, vol. 18, no. 4, 2016, pp. 2624–2661.
- [66] K. He, X. Zhang, S. Ren, and J. Sun, “Deep Residual Learning for Image Recognition,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 770–778, 2016.
- [67] H. Hellström, V. Fodor, and C. Fischione, “Over-the-Air Federated Learning with Retransmissions,” in *Proceedings of the IEEE 22nd International Workshop on Signal Processing Advances in Wireless Communications (SPAWC)*, IEEE, pp. 1–5, 2021.
- [68] H. Hellström, V. Fodor, and C. Fischione, “Over-the-Air Federated Learning with Retransmissions (Extended Version),” *arXiv abs/2111.10267*, 2021.
- [69] G. Hinton, O. Vinyals, and J. Dean, “Distilling the Knowledge in a Neural Network,” *arXiv abs/1503.02531*, 2015.
- [70] Y. Hu, M. Chen, M. Chen, Z. Yang, M. Shikh-Bahaei, H. V. Poor, and S. Cui, “Energy Minimization for Federated Learning with IRS-Assisted Over-the-Air Computation,” in *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Toronto, ON, Canada, 2021.

- [71] K. Huang, G. Zhu, C. You, J. Zhang, Y. Du, and D. Liu, "Communication, Computing, and Learning on the Edge," in *Proceedings of the 2018 IEEE International Conference on Communication Systems (ICCS)*, IEEE, pp. 268–273, 2018.
- [72] S.-J. Huang and Z.-H. Zhou, "Active Query Driven by Uncertainty and Diversity for Incremental Multi-Label Learning," in *2013 IEEE 13th International Conference on Data Mining*, IEEE, pp. 1079–1084, 2014.
- [73] F. Hussain, S. A. Hassan, R. Hussain, and E. Hossain, "Machine Learning for Resource Management in Cellular and IoT Networks: Potentials, Current Solutions, and Open Challenges," *IEEE Communications Surveys & Tutorials*, vol. 22, no. 2, 2020, pp. 1251–1275.
- [74] A. Imteaj, U. Thakker, S. Wang, J. Li, and M. H. Amini, "A Survey on Federated Learning for Resource-Constrained IoT Devices," *IEEE Internet of Things Journal*, 2021, pp. 1–24.
- [75] Y. Inagaki, R. Shinkuma, T. Sato, and E. Oki, "Prioritization of Mobile IoT Data Transmission Based on Data Importance Extracted From Machine Learning Model," *IEEE Access*, vol. 7, 2019, pp. 93 611–93 620.
- [76] S. R. Islam, M. Zeng, O. A. Dobre, and K.-S. Kwak, "Nonorthogonal Multiple Access (NOMA): How It Meets 5G and Beyond," *Wiley 5G Ref: The Essential 5G Reference Online*, 2019, pp. 1–28.
- [77] M. Jaggi, V. Smith, M. Takác, J. Terhorst, S. Krishnan, T. Hofmann, and M. I. Jordan, "Communication-efficient Distributed Dual Coordinate Ascent," in *Advances in Neural Information Processing Systems*, pp. 3068–3076, 2014.
- [78] Y.-S. Jeon, M. M. Amiri, and N. Lee, "Communication-Efficient Federated Learning over MIMO Multiple Access Channels," 2022, under review.
- [79] Y.-S. Jeon, M. M. Amiri, J. Li, and H. V. Poor, "A Compressive Sensing Approach for Federated Learning over Massive MIMO Communication Systems," *IEEE Transactions on Wireless Communications*, vol. 20, no. 3, 2021, pp. 1990–2004.

- [80] E. Jeong, S. Oh, H. Kim, J. Park, M. Bennis, and S.-L. Kim, “Communication-Efficient On-Device Machine Learning: Federated Distillation and Augmentation Under Non-i.i.d. Private Data,” *arXiv abs/1811.11479*, 2018.
- [81] R. Jiang and S. Zhou, “Cluster-Based Cooperative Digital Over-the-Air Aggregation for Wireless Federated Edge Learning,” in *Proceedings of the 2020 IEEE/CIC International Conference on Communications in China (ICCC)*, IEEE, pp. 887–892, 2020.
- [82] T. Jiang and Y. Shi, “Over-the-Air Computation via Intelligent Reflecting Surfaces,” in *Proceedings of the 2019 IEEE Global Communications Conference (GLOBECOM)*, IEEE, pp. 1–6, 2019.
- [83] R. Jin, X. He, and H. Dai, “Communication Efficient Federated Learning with Energy Awareness over Wireless Networks,” *IEEE Transactions on Wireless Communications*, 2022, Early Access.
- [84] P. Kairouz, H. B. McMahan, B. Avent, A. Bellet, M. Bennis, A. N. Bhagoji, K. Bonawitz, Z. Charles, G. Cormode, R. Cummings, R. G. L. D’Oliveira, S. E. Rouayheb, D. Evans, J. Gardner, Z. A. Garrett, A. Gascón, B. Ghazi, P. B. Gibbons, M. Gruteser, Z. Harchaoui, C. He, L. He, Z. Huo, B. Hutchinson, J. Hsu, M. Jaggi, T. Javidi, G. Joshi, M. Khodak, J. Konečný, A. Korolova, F. Koushanfar, O. Koyejo, T. Lepoint, Y. Liu, P. Mittal, M. Mohri, R. Nock, A. Özgür, R. Pagh, M. Raykova, H. Qi, D. Ramage, R. Raskar, D. X. Song, W. Song, S. U. Stich, Z. Sun, A. T. Suresh, F. Tramèr, P. Vepakomma, J. Wang, L. Xiong, Z. Xu, Q. Yang, F. X. Yu, H. Yu, and S. Zhao, “Advances and Open Problems in Federated Learning,” *Foundations and Trends® in Machine Learning*, vol. 14, no. 1-2, 2021, pp. 1–210.
- [85] Y. Koda, K. Yamamoto, T. Nishio, and M. Morikura, “Differentially Private Aircomp Federated Learning with Power Adaptation Harnessing Receiver Noise,” in *Proceedings of the IEEE Global Communications Conference*, pp. 1–6, Taipei, Taiwan, 2021.
- [86] J. Konečný, H. B. McMahan, D. Ramage, and P. Richtárik, “Federated Optimization: Distributed Machine Learning for On-Device Intelligence,” *arXiv abs/1610.02527*, 2016.
- [87] A. Krizhevsky, *Learning Multiple Layers of Features from Tiny Images*, 2009, URL: <https://www.cs.toronto.edu/~kriz/learning-features-2009-TR.pdf>.

- [88] B. Lakshminarayanan, A. Pritzel, and C. Blundell, “Simple and Scalable Predictive Uncertainty Estimation Using Deep Ensembles,” *Advances in Neural Information Processing Systems*, vol. 30, 2017.
- [89] Y. LeCun, L. Bottou, Y. Bengio, P. Haffner, *et al.*, “Gradient-Based Learning Applied to Document Recognition,” *Proceedings of the IEEE*, vol. 86, no. 11, 1998, pp. 2278–2324.
- [90] C.-S. Lee, M.-H. Wang, S.-J. Yen, T.-H. Wei, I.-C. Wu, P.-C. Chou, C.-H. Chou, M.-W. Wang, and T.-H. Yan, “Human vs. Computer Go: Review and Prospect [Discussion Forum],” *IEEE Computational intelligence magazine*, vol. 11, no. 3, 2016, pp. 67–72.
- [91] J. Leng, Z. Lin, M. Ding, P. Wang, D. Smith, and B. Vucetic, “Client Scheduling in Wireless Federated Learning Based on Channel and Learning Qualities,” *IEEE Wireless Communications Letters*, 2022.
- [92] T. Li, A. K. Sahu, A. Talwalkar, and V. Smith, “Federated Learning: Challenges, Methods, and Future Directions,” *IEEE Signal Processing Magazine*, vol. 37, no. 3, 2020, pp. 50–60.
- [93] T. Li, A. K. Sahu, M. Zaheer, M. Sanjabi, A. Talwalkar, and V. Smith, “Federated Optimization in Heterogeneous Networks,” in *Proceedings of Machine Learning and Systems*, 2020, pp. 429–450.
- [94] T. Li, M. Sanjabi, A. Beirami, and V. Smith, “Fair Resource Allocation in Federated Learning,” *arXiv abs/1905.10497*, 2019.
- [95] Y. Lin, S. Han, H. Mao, Y. Wang, and W. J. Dally, “Deep Gradient Compression: Reducing the Communication Bandwidth for Distributed Training,” *arXiv abs/1712.01887*, 2017.
- [96] D. Liu and O. Simeone, “Privacy for Free: Wireless Federated Learning via Uncoded Transmission With Adaptive Power Control,” *IEEE Journal on Selected Areas in Communications*, vol. 39, no. 1, 2021, pp. 170–185.
- [97] D. Liu and O. Simeone, “Channel-driven Monte Carlo Sampling for Bayesian Distributed Learning in Wireless Data Centers,” *IEEE Journal on Selected Areas in Communications*, vol. 40, no. 2, 2022, pp. 562–577.
- [98] D. Liu, G. Zhu, Q. Zeng, J. Zhang, and K. Huang, “Wireless Data Acquisition for Edge Learning: Data-Importance Aware Retransmission,” *IEEE Transactions on Wireless Communications*, vol. 20, no. 1, 2020, pp. 406–420.

- [99] D. Liu, G. Zhu, J. Zhang, and K. Huang, "Wireless Data Acquisition for Edge Learning: Importance-Aware Retransmission," in *Proceedings of the 2019 IEEE 20th International Workshop on Signal Processing Advances in Wireless Communications (SPAWC)*, IEEE, pp. 1–5, 2019.
- [100] D. Liu, G. Zhu, J. Zhang, and K. Huang, "Data-Importance Aware User Scheduling for Communication-efficient Edge Machine Learning," *IEEE Transactions on Cognitive Communications and Networking*, vol. 7, no. 1, 2020, pp. 265–278.
- [101] H. Liu, X. Yuan, and Y.-J. A. Zhang, "CSIT-Free Federated Edge Learning via Reconfigurable Intelligent Surface," *arXiv abs/1905.10497*, 2021.
- [102] H. Liu, X. Yuan, and Y.-J. A. Zhang, "Reconfigurable Intelligent Surface Enabled Federated Learning: A Unified Communication-Learning Design Approach," *IEEE Transactions on Wireless Communications*, 2021, pp. 7595–7609.
- [103] W. Liu, X. Zang, Y. Li, and B. Vucetic, "Over-the-Air Computation Systems: Optimization, Analysis and Scaling Laws," *IEEE Transactions on Wireless Communications*, vol. 19, no. 8, 2020, pp. 5488–5502.
- [104] M. M. Amiri, D. Gündüz, S. R. Kulkarni, and H. V. Poor, "Federated Learning With Quantized Global Model Updates," *arXiv abs/2006.10672*, 2020.
- [105] M. M. Amiri, S. R. Kulkarni, and H. V. Poor, "Federated Learning With Downlink Device Selection," in *Proceedings of the IEEE International Workshop on Signal Processing Advances in Wireless Communications*, Lucca, Italy, 2021.
- [106] M. M. Amiri, T. M. Duman, and D. Gündüz, "Collaborative Machine Learning at the Wireless Edge with Blind Transmitters," in *Proceedings of the IEEE Global Conference on Signal and Information Processing (GlobalSIP)*, pp. 1–5, Ottawa, ON, Canada, 2019.
- [107] C. Ma, J. Konečný, M. Jaggi, V. Smith, M. I. Jordan, P. Richtárik, and M. Takáč, "Distributed Optimization with Arbitrary Local Solvers," *Optimization Methods and Software*, vol. 32, no. 4, 2017, pp. 813–848.

- [108] M. A. Maddah-Ali and U. Niesen, “Fundamental Limits of Caching,” *IEEE Transactions on Information Theory*, vol. 60, no. 5, 2014, pp. 2856–2867.
- [109] Q. Mao, F. Hu, and Q. Hao, “Deep Learning for Intelligent Wireless Networks: A Comprehensive Survey,” *IEEE Communications Surveys & Tutorials*, vol. 20, no. 4, 2018, pp. 2595–2621.
- [110] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas, “Communication-efficient Learning of Deep Networks from Decentralized Data,” in *Artificial intelligence and Statistics*, PMLR, pp. 1273–1282, 2017.
- [111] L. Melis, C. Song, E. D. Cristofaro, and V. Shmatikov, “Exploiting Unintended Feature Leakage in Collaborative Learning,” in *Proceedings of the IEEE Symposium on Security and Privacy*, pp. 691–706, San Francisco, CA, USA, 2019.
- [112] B. Nazer and M. Gastpar, “Computation over Multiple-Access Channels,” *IEEE Transactions on Information Theory*, vol. 53, no. 10, 2007, pp. 3498–3516.
- [113] A. Nedić, A. Olshevsky, and M. G. Rabbat, “Network Topology and Communication-Computation Tradeoffs in Decentralized Optimization,” *Proceedings of the IEEE*, vol. 106, no. 5, 2018, pp. 953–976.
- [114] H. Q. Ngo, A. Ashikhmin, H. Yang, E. G. Larsson, and T. L. Marzetta, “Cell-free Massive MIMO Versus Small Cells,” *IEEE Transactions on Wireless Communications*, vol. 16, no. 3, 2017, pp. 1834–1850.
- [115] D. C. Nguyen, M. Ding, P. N. Pathirana, A. Seneviratne, J. Li, D. Niyato, and H. V. Poor, “Federated Learning for Industrial Internet of Things in Future Industries,” *arXiv abs/2105.14659*, 2021.
- [116] S. Niknam, H. S. Dhillon, and J. H. Reed, “Federated Learning for Wireless Communications: Motivation, Opportunities, and Challenges,” *IEEE Communications Magazine*, vol. 58, no. 6, 2020, pp. 46–51.
- [117] H. Ning, H. Wang, Y. Lin, W. Wang, S. Dhelim, F. Farha, J. Ding, and M. Daneshmand, “A Survey on Metaverse: the State-of-the-art, Technologies, Applications, and Challenges,” *arXiv abs/2111.09673*, 2021.

- [118] T. Nishio and R. Yonetani, "Client Selection for Federated Learning with Heterogeneous Resources in Mobile Edge," in *Proceedings of the 2019 IEEE International Conference on Communications (ICC)*, IEEE, pp. 1–7, 2019.
- [119] S. Oh, J. Park, E. Jeong, H. Kim, M. Bennis, and S.-L. Kim, "Mix2FLD: Downlink Federated Learning After Uplink Federated Distillation With Two-Way Mixup," *IEEE Communications Letters*, vol. 24, no. 10, 2020, pp. 2211–2215.
- [120] J. Park, S. Samarakoon, M. Bennis, and M. Debbah, "Wireless Network Intelligence at the Edge," *Proceedings of the IEEE*, vol. 107, no. 11, 2019, pp. 2204–2239.
- [121] J. Park, S. Wang, A. Elgabli, S. Oh, E. Jeong, H. Cha, H. Kim, S.-L. Kim, and M. Bennis, "Distilling on-Device Intelligence at the Network Edge," *arXiv abs/1908.05895*, 2019.
- [122] S. Park, S. Jung, H. Lee, J. Kim, and J.-H. Kim, "Large-Scale Water Quality Prediction Using Federated Sensing and Learning: A Case Study with Real-World Sensing Big-Data," *Sensors*, vol. 21, no. 4, 2021, p. 1462.
- [123] M. Półka, S. Ptak, and Ł. Kuziora, "The Use of UAV's for Search and Rescue Operations," *Procedia Engineering*, vol. 192, 2017, pp. 748–752.
- [124] S. Prakash, H. Hashemi, Y. Wang, M. Annavaram, and S. Avestimehr, "Byzantine-Resilient Federated Learning with Heterogeneous Data Distribution," *arXiv abs/2010.07541*, 2020.
- [125] S. Pu and A. Nedić, "Distributed Stochastic Gradient Tracking Methods," *Mathematical Programming*, vol. 187, no. 1, 2020, pp. 409–457.
- [126] Z. Qin, G. Y. Li, and H. Ye, "Federated Learning and Wireless Communications," *IEEE Wireless Communications*, vol. 28, no. 5, 2021.
- [127] A. Reisizadeh, A. Mokhtari, H. Hassani, A. Jadbabaie, and R. Pedarsani, "FedPAQ: A Communication-Efficient Federated Learning Method with Periodic Averaging and Quantization," in *Proceedings of the International Conference on Artificial Intelligence and Statistics*, PMLR, pp. 2021–2031, 2020.

- [128] J. Ren, Y. He, D. Wen, G. Yu, K. Huang, and D. Guo, "Scheduling for Cellular Federated Edge Learning With Importance and Channel Awareness," *IEEE Transactions on Wireless Communications*, vol. 19, no. 11, 2020, pp. 7690–7703.
- [129] J. Ren, G. Yu, and G. Ding, "Accelerating DNN Training in Wireless Federated Edge Learning Systems," *IEEE Journal on Selected Areas in Communications*, vol. 39, no. 1, 2020, pp. 219–232.
- [130] P. Rost, C. Mannweiler, D. S. Michalopoulos, C. Sartori, V. Sciancalepore, N. Sastry, O. Holland, S. Tayade, B. Han, D. Bega, *et al.*, "Network Slicing to Enable Scalability and Flexibility in 5G Mobile Networks," *IEEE Communications Magazine*, vol. 55, no. 5, 2017, pp. 72–79.
- [131] W. Saad, M. Bennis, and M. Chen, "A Vision of 6G Wireless Systems: Applications, Trends, Technologies, and Open Research Problems," *IEEE Network*, vol. 34, no. 3, 2019, pp. 134–142.
- [132] Y. Saito, Y. Kishiyama, A. Benjebbour, T. Nakamura, A. Li, and K. Higuchi, "Non-Orthogonal Multiple Access (NOMA) for Cellular Future Radio Access," in *Proceedings of the 2013 IEEE 77th Vehicular Technology Conference (VTC Spring)*, IEEE, pp. 1–5, 2013.
- [133] M. Salehi and E. Hossain, "Federated Learning in Unreliable and Resource-Constrained Cellular Wireless Networks," *IEEE Transactions on Communications*, 2021.
- [134] S. Samarakoon, M. Bennis, W. Saad, and M. Debbah, "Distributed Federated Learning for Ultra-Reliable Low-Latency Vehicular Communications," *IEEE Transactions on Communications*, vol. 68, no. 2, 2019, pp. 1146–1159.
- [135] F. Seide, H. Fu, J. Droppo, G. Li, and D. Yu, "1-bit Stochastic Gradient Descent and its Application to Data-Parallel Distributed Training of Speech DNNs," in *Proceedings of the 15th Annual Conference of the International Speech Communication Association*, 2014.
- [136] M. Seif, R. Tandon, and M. Li, "Wireless Federated Learning with Local Differential Privacy," in *Proceedings of the IEEE International Symposium on Information Theory*, pp. 2604–2609, 2020.

- [137] T. Sery and K. Cohen, "A Sequential Gradient-Based Multiple Access for Distributed Learning over Fading Channels," in *Proceedings of the 2019 57th Annual Allerton Conference on Communication, Control, and Computing (Allerton)*, IEEE, pp. 303–307, 2019.
- [138] T. Sery and K. Cohen, "On Analog Gradient Descent Learning over Multiple Access Fading Channels," *IEEE Transactions on Signal Processing*, vol. 68, 2020, pp. 2897–2911.
- [139] B. Settles, "Active Learning Literature Survey," University of Wisconsin-Madison Department of Computer Sciences, Tech. Rep., 2009.
- [140] B. Settles, "Active Learning," *Synthesis Lectures on Artificial Intelligence and Machine Learning*, vol. 6, no. 1, 2012, pp. 1–114.
- [141] O. Shamir, N. Srebro, and T. Zhang, "Communication-Efficient Distributed Optimization Using an Approximate Newton-Type Method," in *Proceedings of the International Conference on Machine Learning*, PMLR, pp. 1000–1008, 2014.
- [142] W. Shi, S. Zhou, and Z. Niu, "Device Scheduling with Fast Convergence for Wireless Federated Learning," in *Proceedings of the 2020 IEEE International Conference on Communications (ICC)*, IEEE, pp. 1–6, 2020.
- [143] Y. Shi, Y. Zhou, and Y. Shi, "Over-the-Air Decentralized Federated Learning," in *Proceedings of the 2021 IEEE International Symposium on Information Theory (ISIT)*, IEEE, pp. 455–460, 2021.
- [144] Y. Shi, K. Yang, T. Jiang, J. Zhang, and K. B. Letaief, "Communication-efficient Edge AI: Algorithms and Systems," *IEEE Communications Surveys & Tutorials*, vol. 22, no. 4, 2020, pp. 2167–2191.
- [145] R. Shinkuma and T. Nishio, "Data Assessment and Prioritization in Mobile Networks for Real-Time Prediction of Spatial Information with Machine Learning," in *Proceedings of the IEEE First International Workshop on Network Meets Intelligent Computations (NMIC)*, IEEE, pp. 1–6, 2019.
- [146] N. Shlezinger, M. Chen, Y. C. Eldar, H. V. Poor, and S. Cui, "UVE-QFed: Universal Vector Quantization for Federated Learning," *IEEE Transactions on Signal Processing*, vol. 69, 2020, pp. 500–514.
- [147] H. Sifaou and G. Y. Li, "Robust Federated Learning via Over-The-Air Computation," *arXiv abs/2111.01221*, 2021.

- [148] J. M. B. da Silva Jr., K. Ntougias, I. Krikidis, G. Fodor, and C. Fischione, “Simultaneous Wireless Information and Power Transfer for Federated Learning,” *arXiv abs/2104.12749*, 2021.
- [149] K. Simonyan and A. Zisserman, “Very Deep Convolutional Networks for Large-Scale Image Recognition,” *arXiv abs/1409.1556*, 2014.
- [150] J. So, B. Güler, and A. S. Avestimehr, “CodedPrivateML: A Fast and Privacy-Preserving Framework for Distributed Machine Learning,” *IEEE Journal on Selected Areas in Information Theory*, vol. 2, no. 1, 2021, pp. 441–451.
- [151] E. C. Strinati, S. Barbarossa, J. L. Gonzalez-Jimenez, D. Ktenas, N. Cassiau, L. Maret, and C. Dehos, “6G: The Next Frontier: From Holographic Messaging to Artificial Intelligence Using Subterahertz and Visible Light Communication,” *IEEE Vehicular Technology Magazine*, vol. 14, no. 3, 2019, pp. 42–50.
- [152] C. Sun, A. Shrivastava, S. Singh, and A. Gupta, “Revisiting Unreasonable Effectiveness of Data in Deep Learning Era,” in *Proceedings of the IEEE international conference on computer vision*, pp. 843–852, 2017.
- [153] Y. Sun, S. Zhou, and D. Gündüz, “Energy-Aware Analog Aggregation for Federated Learning with Redundant Data,” in *Proceedings of the 2020 IEEE International Conference on Communications (ICC)*, IEEE, pp. 1–7, 2020.
- [154] S. Sundhar Ram, A. Nedić, and V. V. Veeravalli, “Distributed Stochastic Subgradient Projection Algorithms for Convex Optimization,” *Journal of optimization theory and applications*, vol. 147, no. 3, 2010, pp. 516–545.
- [155] H. Tang, X. Lian, C. Yu, T. Zhang, and J. Liu, “DoubleSqueeze: Parallel Stochastic Gradient Descent with Double-Pass Error-Compensated Compression,” in *Proceedings of the International Conference on Machine Learning*, Long Beach, CA, 2019.
- [156] *Timing Advance (TA) in LTE*, 2010, URL: <http://4g5gworld.com/blog/timing-advance-ta-lte>.
- [157] S. Timotheou, I. Krikidis, G. Zheng, and B. Ottersten, “Beamforming for MISO Interference Channels with QoS and RF Energy Transfer,” *IEEE Transactions on Wireless Communications*, vol. 13, no. 5, 2014, pp. 2646–2658.

- [158] N. H. Tran, W. Bao, A. Zo5a, M. N. Nguyen, and C. S. Hong, "Federated Learning over Wireless Networks: Optimization Model Design and Analysis," in *Proceedings of the IEEE INFOCOM 2019 Conference on Computer Communications*, IEEE, pp. 1387–1395, 2019.
- [159] A. Vempaty, L. Tong, and P. K. Varshney, "Distributed Inference with Byzantine Data: State-of-the-Art Review on Data Falsification Attacks," *IEEE Signal Processing Magazine*, vol. 30, no. 5, 2013, pp. 65–75.
- [160] T. T. Vu, D. T. Ngo, N. H. Tran, H. Q. Ngo, M. N. Dao, and R. H. Middleton, "Cell-Free Massive MIMO for Wireless Federated Learning," *IEEE Transactions on Wireless Communications*, vol. 19, no. 10, 2020, pp. 6377–6392.
- [161] M. M. Wadu, S. Samarakoon, and M. Bennis, "Federated Learning Under Channel Uncertainty: Joint Client Scheduling and Resource Allocation," in *Proceedings of the 2020 IEEE Wireless Communications and Networking Conference (WCNC)*, IEEE, pp. 1–6, 2020.
- [162] J. Wang, C. Jiang, H. Zhang, Y. Ren, K.-C. Chen, and L. Hanzo, "Thirty Years of Machine Learning: The Road to Pareto-Optimal Wireless Networks," *IEEE Communications Surveys & Tutorials*, vol. 22, no. 3, 2020, pp. 1472–1514.
- [163] S. Wang, Y.-C. Wu, M. Xia, R. Wang, and H. V. Poor, "Machine Intelligence at the Edge With Learning Centric Power Allocation," *IEEE Transactions on Wireless Communications*, vol. 19, no. 11, 2020, pp. 7293–7308.
- [164] X. Wang, Y. Han, V. C. Leung, D. Niyato, X. Yan, and X. Chen, "Convergence of Edge Computing and Deep Learning: A Comprehensive Survey," *IEEE Communications Surveys & Tutorials*, vol. 22, no. 2, 2020, pp. 869–904.
- [165] Z. Wang, J. Qiu, Y. Zhou, Y. Shi, L. Fu, W. Chen, and K. B. Letaief, "Federated Learning via Intelligent Reflecting Surface," *IEEE Transactions on Wireless Communications*, 2021, pp. 808–822.
- [166] K. Wei, J. Li, M. Ding, C. Ma, H. H. Yang, F. Farokhi, S. Jin, T. Q. Quek, and H. V. Poor, "Federated Learning with Differential Privacy: Algorithms and Performance Analysis," *IEEE Transactions on Information Forensics and Security*, vol. 15, 2020, pp. 3454–3469.

- [167] D. Wen, X. Li, Q. Zeng, J. Ren, and K. Huang, "An Overview of Data-Importance Aware Radio Resource Management for Edge Machine Learning," *Journal of Communications and Information Networks*, vol. 4, no. 4, 2019, pp. 1–14.
- [168] A. G. Wilson and P. Izmailov, "Bayesian Deep Learning and a Probabilistic Perspective of Generalization," *Advances in neural information processing systems*, vol. 33, 2020, pp. 4697–4708.
- [169] T. Wu, F. Wu, J.-M. Redoute, and M. R. Yuce, "An Autonomous Wireless Body Area Network Implementation Towards IoT Connected Healthcare Applications," *IEEE Access*, vol. 5, 2017, pp. 11 413–11 422.
- [170] Q. Xia, W. Ye, Z. Tao, J. Wu, and Q. Li, "A Survey of Federated Learning for Edge Computing: Research Problems and Solutions," *High-Confidence Computing*, 2021, p. 100 008.
- [171] W. Xia, W. Wen, K. -. Wong, T. Q. S. Quek, J. Zhang, and H. Zhu, "Federated-Learning-Based Client Scheduling for Low-Latency Wireless Communications," *IEEE Wireless Communications*, vol. 28, no. 2, 2021, pp. 32–38.
- [172] R. Xin, S. Kar, and U. A. Khan, "Decentralized Stochastic Optimization and Machine Learning: A Unified Variance-Reduction Framework for Robust Performance and Fast Convergence," *IEEE Signal Processing Magazine*, vol. 37, no. 3, 2020, pp. 102–113.
- [173] H. Xing, O. Simeone, and S. Bi, "Decentralized Federated Learning via SGD over Wireless D2D Networks," in *Proceedings of the 2020 IEEE 21st International Workshop on Signal Processing Advances in Wireless Communications (SPAWC)*, IEEE, pp. 1–5, 2020.
- [174] J. Xu and H. Wang, "Client Selection and Bandwidth Allocation in Wireless Federated Learning Networks: A Long-Term Perspective," *IEEE Transactions on Wireless Communications*, vol. 20, no. 2, 2021, pp. 1188–1200.
- [175] J. Xu, H. Wang, and L. Chen, "Bandwidth Allocation for Multiple Federated Learning Services in Wireless Edge Networks," *arXiv abs/2101.03627*, 2021.

- [176] P. Xue, P. Gong, J. H. Park, D. Park, and D. K. Kim, “Max-Min Fairness Based Radio Resource Management in Fourth Generation Heterogeneous Networks,” in *Proceedings of the 9th International Symposium on Communications and Information Technology*, IEEE, pp. 208–213, 2009.
- [177] P. Yagol, F. Ramos, S. Trilles, J. Torres-Sospedra, and F. J. Perales, “New Trends in Using Augmented Reality Apps for Smart City Contexts,” *ISPRS International Journal of Geo-Information*, vol. 7, no. 12, 2018, p. 478.
- [178] H. H. Yang, A. Arafa, T. Q. S. Quek, and H. V. Poor, “Age-Based Scheduling Policy for Federated Learning in Mobile Edge Networks,” in *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2020.
- [179] K. Yang, T. Jiang, Y. Shi, and Z. Ding, “Federated Learning via Over-the-Air Computation,” *IEEE Transactions on Wireless Communications*, vol. 19, no. 3, 2020, pp. 2022–2035.
- [180] Q. Yang, Y. Liu, Y. Cheng, Y. Kang, T. Chen, and H. Yu, *Federated Learning*, ser. Synthesis Lectures on Artificial Intelligence and Machine Learning. Morgan & Claypool Publishers, 2019.
- [181] Z. Yang, M. Chen, W. Saad, C. S. Hong, and M. Shikh-Bahaei, “Energy Efficient Federated Learning Over Wireless Communication Networks,” *IEEE Transactions on Wireless Communications*, vol. 20, no. 3, 2021, pp. 1935–1949.
- [182] W. Yu and J. M. Cioffi, “On Constant Power Water-Filling,” in *IEEE International Conference on Communications. Conference Record (ICC)*, IEEE, vol. 6, pp. 1665–1669, 2002.
- [183] S. Yue, J. Ren, J. Xin, D. Zhang, Y. Zhang, and W. Zhuang, “Efficient Federated Meta-Learning over Multi-Access Wireless Networks,” *IEEE Journal on Selected Areas in Communications*, vol. 40, 2022, pp. 1556–1570.
- [184] X. Zang, W. Liu, Y. Li, and B. Vucetic, “Over-the-Air Computation Systems: Optimal Design with Sum-Power Constraint,” *IEEE Wireless Communications Letters*, vol. 9, no. 9, 2020, pp. 1524–1528.
- [185] Q. Zeng, Y. Du, and K. Huang, “Wirelessly Powered Federated Edge Learning: Optimal Tradeoffs Between Convergence and Power Transfer,” *arXiv abs/2102.12357*, 2021.

- [186] Q. Zeng, Y. Du, K. Huang, and K. K. Leung, "Energy-Efficient Radio Resource Allocation for Federated Edge Learning," in *Proceedings of the 2020 IEEE International Conference on Communications Workshops (ICC Workshops)*, IEEE, pp. 1–6, 2020.
- [187] T. Zeng, O. Semiari, M. Chen, W. Saad, and M. Bennis, "Federated Learning on the Road: Autonomous Controller Design for Connected and Autonomous Vehicles," *arXiv abs/2102.03401*, 2021.
- [188] H. Zhang, M. Cisse, Y. N. Dauphin, and D. Lopez-Paz, "mixup: Beyond Empirical Risk Minimization," *arXiv abs/1710.09412*, 2017.
- [189] N. Zhang and M. Tao, "Gradient Statistics Aware Power Control for Over-the-Air Federated Learning," *IEEE Transactions on Wireless Communications*, 2021, pp. 5115–5128.
- [190] S. Zhang, S. C. Liew, and P. P. Lam, "Hot Topic: Physical-Layer Network Coding," in *Proceedings of the 12th Annual International Conference on Mobile Computing and Networking*, ACM, pp. 358–365, 2006.
- [191] Z. Zhang, Y. Xiao, Z. Ma, M. Xiao, Z. Ding, X. Lei, G. K. Karagiannis, and P. Fan, "6G Wireless Networks: Vision, Requirements, Architecture, and Key Technologies," *IEEE Vehicular Technology Magazine*, vol. 14, no. 3, 2019, pp. 28–41.
- [192] J. Zhao, "A Survey of Intelligent Reflecting Surfaces (IRSs): Towards 6G Wireless Communication Networks," *arXiv abs/1907.04789*, 2019.
- [193] Y. Zhao, M. Li, L. Lai, N. Suda, D. Civin, and V. Chandra, "Federated Learning With non-IID Data," *arXiv abs/1806.00582*, 2018.
- [194] Z. Zhou, X. Chen, E. Li, L. Zeng, K. Luo, and J. Zhang, "Edge Intelligence: Paving the Last Mile of Artificial Intelligence with Edge Computing," *Proceedings of the IEEE*, vol. 107, no. 8, 2019, pp. 1738–1762.
- [195] G. Zhu, L. Chen, and K. Huang, "Over-the-Air Computation in MIMO Multi-Access Channels: Beamforming and Channel Feedback," *CoRR*, vol. abs/1803.11129, 2018.
- [196] G. Zhu, Y. Du, D. Gündüz, and K. Huang, "One-Bit Over-the-Air Aggregation for Communication-Efficient Federated Edge Learning: Design and Convergence Analysis," *IEEE Transactions on Wireless Communications*, vol. 20, no. 3, 2020, pp. 2120–2135.

- [197] G. Zhu, D. Liu, Y. Du, C. You, J. Zhang, and K. Huang, “Towards an Intelligent Edge: Wireless Communication Meets Machine Learning,” *arXiv abs/1809.00343*, 2018.
- [198] G. Zhu, D. Liu, Y. Du, C. You, J. Zhang, and K. Huang, “Toward an Intelligent Edge: Wireless Communication Meets Machine Learning,” *IEEE Communications Magazine*, vol. 58, no. 1, 2020, pp. 19–25.
- [199] G. Zhu, Y. Wang, and K. Huang, “Broadband Analog Aggregation for Low-Latency Federated Edge Learning (extended version),” *arXiv abs/1812.11494*, 2018.
- [200] G. Zhu, Y. Wang, and K. Huang, “Broadband Analog Aggregation for Low-Latency Federated Edge Learning,” *IEEE Transactions on Wireless Communications*, vol. 19, no. 1, 2020, pp. 491–506.
- [201] G. Zhu, J. Xu, K. Huang, and S. Cui, “Over-the-Air Computing for Wireless Data Aggregation in Massive IoT,” *IEEE Wireless Communications*, vol. 28, no. 4, 2021, pp. 57–65.
- [202] M. Zhu and S. Gupta, “To Prune, or not to Prune: Exploring the Efficacy of Pruning for Model Compression,” *arXiv abs/1710.01878*, 2017.
- [203] M. Zinkevich, M. Weimer, L. Li, and A. J. Smola, “Parallelized Stochastic Gradient Descent,” in *Advances in Neural Information Processing Systems*, pp. 2595–2603, 2010.