

# Federated Inference with Reliable Uncertainty Quantification over Wireless Channels via Conformal Prediction

Meiyi Zhu (Member), IEEE, Matteo Zecchin (Student Member), IEEE, Sangwoo Park (Member), IEEE, Chai Guo (Senior Member), IEEE, Thuy Nam Fang (Senior Member), IEEE, Osvaldo Simeone (Fellow), IEEE

### Abstract

Consider a setting in which a device and a server share a shared channel. The device wishes to make an inference on a new input given in input  $x$ . Devices have devices that have previously learned data used for training and for communicating to the server over a common wireless channel. If the devices have no access to the new input, can communication from device to the server enhance the quality of the inference decision at the server? Recent work has introduced federated learning [1]. Recent work [2] introduced GP which leverages devices-to-server (GP), which leverages devices-to-server communication to improve the reliability of the server's decision. With federated GP, devices communicate to the server information about the loss accrued by the shared information model on the data calculated by the server to refine its information to calibrate a and the server leverages or lists it for guaranteed talkback in the decision interval, with respect that it is guaranteed to provide the worker's answer with pre-defined level of reliability level. Previous work [3] assumed a single user for the server. In this paper, we study federated GP in the setting of multiple users.

We introduce a novel protocol, termed *wireless federated conformal prediction* (WFCP), which builds upon the work of M. Zhu, C. Guo and C. Feng was supported by the Fundamental Research Funds for the Central Universities (No.2221XD-A01-1) and by the Beijing Natural Science Foundation (J2202043). The work of M. Zhu was also supported by the BUPT Excellent Ph.D. Students Foundation (2020-4). The work of M. Zechin was also supported by the Beijing Natural Science Foundation (2020-4). The work of M. Zechin and Q. Simeone was supported by the European Union's Horizon Europe project (ZENTRIG ID 101896379). The support of Simeone was also supported by an Open Fellowship of the EPSRC (EP/W024101/3) and by Project REASON, an UK Government funded EPSRC project under the Future Open Networks Research Challenge (FONRC) sponsored by the Future Department of Science Innovation and Technology (DSIT).

Meiyi Zhu, Caili Guo and Chunyan Feng are with the Beijing Key Laboratory of Network System Architecture and Omnipotent Engineering School, School of Informatics, Command and Computer Engineering, Beijing Jiaotong University, Beijing, China. (e-mail: diannuo@bjtu.edu.cn; guocaili@bjtu.edu.cn; cyfeng@bjtu.edu.cn).

Matteo Zecchin, Sangwoo Park, and Osvaldo Simeone are with the King's Communications, Learning & Information Processing (KCLIP) lab, Department of Engineering, King's College London, London WC2R 2LS, UK. (e-mail: matteo.zecchin@kcl.ac.uk; sangwoo.park@kcl.ac.uk; osvaldo.simeone@kcl.ac.uk).

In this paper, we study federated inference for type-based multiple access in the wireless setting. WFCP is proposed to provide reliability guarantees for federated inference and prediction (WFCP), which provides reliability guarantees. It is type-based (TBM) and demonstrates the significant advantages of WFCP against digital implementations of existing federated CP schemes, especially in regimes with limited communication resources and/or large number of devices.

Conformal prediction, federated inference, wireless communications, type-based multiple access.  
Index Terms

Conformal prediction, federated inference, wireless communications, type-based multiple access.

**I. INTRODUCTION**  
*Federation* is a data processing paradigm whereby distributed devices with local, possibly private, data sets collaborate for the purpose of carrying out a shared information processing task without the direct exchange of the local data sets. The main exemplar of federated data

Federation is a data processing paradigm whereby distributed devices with local, possibly private, data sets collaborate for the purpose of carrying out a shared information processing task without the direct exchange of the local data sets. The main exemplar of federated data processing is *federated learning*, which addresses the task of training a machine learning model. Federated learning has been widely studied in recent years, with research activities ranging from theoretical analyses [?], [?] to the design of communication protocols [?], [?] and to testbeds [?], processing is federated learning, which addresses the task of training a machine learning [?]. This paper focuses on a different federated data processing task, namely *federated inference*, model. Federated learning has been widely studied in recent years, with research activities with the goal of leveraging collaboration across devices to ensure *reliable decision-making*, ranging from theoretical analyses [?], [?] to the design of communication protocols [?], [?]

and to testbeds [?], [?]. This paper focuses on a different federated data processing task, **A. Federated Reliable Inference**

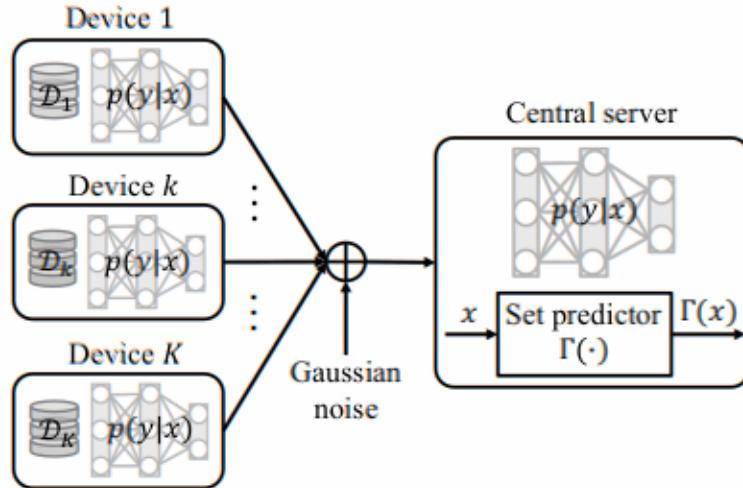
namely federated inference, with the goal of leveraging collaboration across devices to ensure reliable decision-making.

As illustrated in Fig. ??, we study a setting in which devices and a server share a pre-trained machine learning model. The model may have been obtained through a previous phase of federated learning, or it may have been downloaded from a repository of existing models **A. Federated Reliable Inference**

trained in any other arbitrary manner. The server wishes to make an inference on a new input

As illustrated in Fig. ??, we study a setting in which devices and a server share a pre-given the model. Specifically, given an input, it wishes to produce an *interval*, or *set*, of possible trained machine learning model. The model may have been obtained through a previous output values that is guaranteed to contain the correct answer with a pre-defined target *reliability* phase of federated learning, or it may have been downloaded from a repository of existing level. Devices have access to data, previously not used for training, and can communicate to the models trained in any other arbitrary manner. The server wishes to make an inference on server over wireless channels. The devices do not have access to the new input. a new input given the model. Specifically, given an input, it wishes to produce an interval,

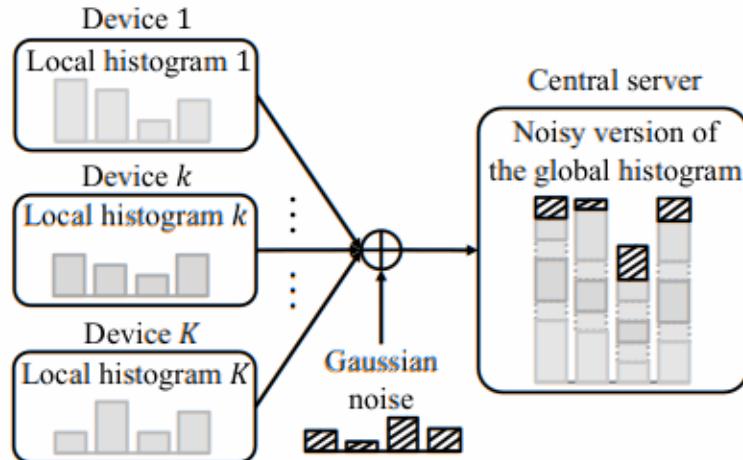
For this setting, recent work has introduced *federated conformal prediction* (CP), which or set, of possible output values that is guaranteed to contain the correct answer with leverages devices-to-server communication to support reliable decision-making at the server [?]. a pre-defined target reliability level. Devices have access to data, previously not used for With federated CP, devices communicate to the server information about the performance accrued training, and can communicate to the server over wireless channels. The devices do not by the shared pre-trained model on the local data. Intuitively, this information provides a yardstick have access to the new input.



**Fig. 1.** Illustration of the virile label fedated in infected people blend medical study preprint aimed at learning a model  $p(y|x)$  available and used. The server wishes to make a reliable prediction input  $x$ , which is not available at the doctor's office. Following the prediction rule, the prediction takes  $\Gamma(x)$  of the label space  $\mathcal{V}$ . The goal is to predict the label  $y$  with probability  $p(y|x)$  against the predicted with probability  $p(\Gamma(x)|x)$ . The label weight probability is zero smaller than  $p(y|x)$ . To get this reliability leveled  $k$ -confidence (??), we must endow the label with information about shared cell death. This information is then used as they secrete to update this information) but the nested calibration prediction to update sharing the label  $p(y|x)$  to condition (??) calibrate the prediction  $\Gamma(x)$ , ensuring the reliability condition (??).

with which the system can work has a probability of failure value of the order of 0.01, for which given a single instance if the model fails to support reliable decision-making at 90% of the data points fed to the device, the device server may safely exclude from the path one interval per output value as unreliable. The model assigns a loss margin of 1%, which is it wishes to give, to this 90% reliability level. In other words, the device server informs the vehicle of the reliability of each value it outputs in relation to the given input. For instance, if the model obtains a loss no higher than 0.01, it is assumed noise 90% of the time it is in operation, where the device can be considered safe up to the specified level of risk [2]. The proposed values are quantiles of confidence (QCs) which are assigned, respectively, to the device's output values. In this paper, we therefore consider the calibration of the device in a two-tier setting.

Previous work [?] assumed noise-free communication, whereby devices can communicate a single real number to the server. Specifically, reference [?] proposed a quantile-of-quantile (QQ) scheme, referred to as FedCP-QQ1 whereby each device computes FedCP and communicates a predetermined quantile of the local losses. In this paper, we study for the first time



**Fig. 2.** TBBM A neighbor histogram of the global histogram of received data over available devices. All devices, which are assigned to each histogram bin, are used to transmit simultaneously their initial local histograms. The individual channel histograms are allocated to wireless channels proportionally to the corresponding bin probability. This way, the corresponding bin probability is the sum of the global histogram thanks to the superposition of the global histogram thanks to the superposition of the signals received for each orthogonal codeword.

the 90-percentile of the losses obtained by the pre-trained model across *all* devices. However, the quantile-of-quantiles targeted by FedCP-QQ provides a generally inaccurate estimate of the overall quantile, particularly when the number of devices is large. Furthermore, a direct implementation of FedCP-QQ [23] on wireless quantiles, though requiring the transmission of quantized local quantiles, would require the transmission of FedCP-QQ of quantized local quantiles. In fact, requiring a bandwidth that increases proportionally to the number of available devices.

In this paper, however, the quantile-of-quantiles targeted by FedCP-QQ performs much better (WEGP), which addresses the shortcomings by building a hybrid-based multiple access of devices (TBMA) [3]. Furthermore, a novel quantile implementation of TBMA-QQ [multiple wireless channels that would accommodate the aggregation of quantiles rather than quantiles, requiring a hybrid channel, that increased proportionally with the number of available devices] available across the devices at the server. To explain this, suppose that each device has its own quantized data wireless channel, a generally different histogram (WECP), which in Fig. 2 shows the goings on the server by building on the aggregated histograms across all devices [2], i.e., the histogram of the data available at all TBMA is without multiple messages separately estimate the histograms of all devices.

To accomplish this objective in TBMA, each histogram bin is assigned an orthogonal code word. Devices divide their transmission in a manner that number of orthogonal codewords in

data with a given energy generally is allocated histogram bins corresponding to the bin width of the larger is to best of my knowledge. A low level of noise as it varies with the histogram bin width due to the availability of available devices, without the inverse separation in the global histogram of all devices the superposition of this objective function. In TBM iteration, histogram bin is assigned an orthogonal code word. By adopting TBM, we provide the communication protocol of the proposed WFGP scheme allows the server to estimate the histogram of the losses accrued by the preexisting model histograms all the devices. The estimate of the update to the histogram bins is achieved by transmitting important values by dividing them into  $\Delta$  bins. The number of bins and position of the bins are communicated.

By adopting the FBM challenge in the design of WFGPs, how proposed WFGPs try to handle the situation that it implied the integral(s) of the losses were contained by the pre-trained model has addressed the liability level. This is mainly based on the fact that desired global of the global importance of the contributions, and denoted by its TBM, is also independently noisy (see Figure 23). The quantifications this problem by proposing a novel quantitative method that is proved to guarantee reliability.

The main technical challenge in the design of WFCP is how to ensure reliability – that is, the condition that the predicted interval/set at the server contains the true output with the desired probability level. This is challenging as it requires the server to have access to the estimated CP probabilities and the mean of TBM losses, and hence of its quantile, is inherently noisy (see Fig. ??). WFCP addresses this problem by proposing a federated CP-QQ estimator that provides a guarantee of reliability applied CP in federated settings, aiming to provide distribution-free, set-valued predictions with reliability guarantees. In [?], each device calculates a quantile of its local losses, and the server aggregates these quantiles from all devices to form an average. However, applying a brief review of related work by focusing on guaranteed reliability CP procedures, the limitation of CP-QQ [?] was proposed whereby a QQ estimator is used in lieu of an average of quantiles. Prior to establishing a formal reliability of CP-QQs from federated CP, reference [?] initially applies CP in generalized loss settings with statistical distributions across the predictions with the reliability. In [?], the authors proposed an approach that quantifies that the predictor is well calibrated with respect to a specific structural distribution of the deviance. In contrast, applying CP without the notion of quantile was proposed to apply reliability to unbiased estimation methods [?]. FedCP-QQ [?] was proposed to approximate quantile QQs to implement an estimate of the quantile of reliability, guaranteeing a reliable prediction with a confidence interval. The quality of the estimation error. In parallel, reference [?] studied a related setting with label

dissemination of GRs has been defined by generalizing to the weighted type with optimization scheme proposed in [8] for decentralized CP. In that devices by means of joint channel and power beamforming techniques that minimize total power cost for the whole applied with distributed approach to jointly of distributed units from distributed nodes. According to unlike when setting is considered, this paper proposed previous apply distributed approach to estimation which has [9] parallel, proposed a new scheme in [10]. In this iterative, the inquiry function of the channel and liability guarantees are only proved under existing failure probability of quality of service influence of reliability, confirming that studied in this paper is with added distribution knowledge gap by investigating the problem of weightless feature of GR and by focusing on the impact of channel estimation quality in particularity, noting that model calibration be unlikely [11], as [12] has in [13] two of the following reliability guarantees will be additional assessed through the quality of the quantile estimates. Distributed devices in According [14], we focus on statistically homogeneous data across devices work, which assumes one TBMA. The pioneering paper [15], [16] introduced the TBMA where [17] is orthogonal, codebooks and assigned to different measurement values across multiple devices and a variant of the manner is likely to be tested in a CP is designed to accomplish single path estimation irrespective the joint TBMA is applied. This paper multi-path knowledge gap by calves it up by leveraging the orthogonal-cell wise of the TBMA GR inter-cell non-orthogonal frequency of the strategy under centralized and decentralized decoding setting can be found in [18], [19], [20] developed a target for orthogonal strategy of TBMA without added overhead in random quality of radio resources. It is shown that assuring further more sparse retransmission [8]. We focus on statistically homogeneous passing between devices designed for a single-cell [21] and a multi-cell fog-radio access network [22] respectively. Ref TBMA. [23] proposed papers [24] and [25] design of TBMA protocol whereby orthogonal information is bottleneck principle different adopted user criteria to jointly to optimize and enhance and of the network has been limited to estimator unknown device source and accomplish single path. Note of this paper is provide an [26] in TBMA as the key problem of testing TBMA for reliable parameters estimation in-cell set-up by leveraging in-cell orthogonal TBMA and inter-cell non-orthogonal frequency reuse.

### D. Contributions and Organization

In this paper, we introduce a WFCP of the TBMA wireless protocol for the implementation of federated scenario. The main contributions of this paper are summarized as follows:

- We first review conventional centralized GR which handles that all data is available at the access network. Then we propose a digital communication framework for the design of TBMA protocols of the early federated GR is bottlenecked. The QoS [27] is adopted as the metric which will jointly

optimize the link selection for WFCP. To this end, we assume that under the known orthogonal channel statistics available in the paper, a time division insights into the FedMA-based quantile estimation over the correctly received quantiles.

#### D. Contributions and Organization

In this paper, we introduce WFCP, the first wireless protocol for the implementation of federated CP. The main contributions of this paper are summarized as follows:

- We provide a rigorous analysis of the reliability performance of WFCP, proving that it can achieve any target reliability level. This is achieved by providing the details of this achievable implementation design [2]. Furthermore, we introduce a digital quantization levels framework in order to validate the scheme.
- Supplying simulation results to demonstrate the advantage of the proposed WFCP scheme over existing, inhomogeneous, especially baseline schemes of WFCP communication, resources and data among the devices.

The remainder of this paper is organized as follows. In Sec. ??, we describe the setting and define local quantile. Sec. ?? presents the general framework of conventional CP, and introduces a quantile threshold that accounts for the presence of channel noise.

We propose WFCP, a novel protocol based on TBMA that hinges on a carefully selected quantile threshold that accounts for the presence of channel noise. We provide a rigorous analysis of the reliability performance of WFCP, also providing design guidelines and proof of reliability. Sec. ?? evaluates the performance of WFCP as compared to benchmarks via experiments, validating the effectiveness of WFCP. Sec. ?? summarizes this paper and points to directions for future work.

existing strategies, especially in the presence of limited communication resources and/or large number of devices.

## II. SETTING AND PROBLEM DEFINITION

**A. Setting** The remainder of this paper is organized as follows. In Sec. ??, we describe the setting and define the problem. Sec. ?? presents the general framework of conventional CP, and introduces a quantile threshold that accounts for the presence of channel noise. We consider a wireless federated inference scenario in which a set of  $K$  devices and a central server communicate over a multiple access channel. A pre-trained machine learning model is available at both server and devices side. This model may have been previously trained using federated learning [?]. As in [?], we focus on the problem of reliable collaborative, or federated, inference. Sec. ?? using a fixed model along with communication between devices and server. In this setting, communication is not used to optimize the machine learning model, as points to directions for future work.

## II. Setting and Problem Definition

A. **Setting** communication devices can help the server *calibrate* its decision, enhancing the server's estimate of the corresponding uncertainty.

We consider a wireless federated inference scenario in which a set of  $K$  devices and a central server communicate over a multiple access channel. A pre-trained machine learning model is available at both server and devices side. This model may have been previously trained using federated learning [?]. As in [?], we focus on the problem of reliable collaborative, interpreted as a measure of the *confidence* that the model has in label  $y$  being the correct one. Conventionally, a decision  $y^*$  is obtained by selecting the label on which the model has maximum confidence, i.e.,

$$\text{decision produced by the server on a new input } x = \arg \max_{y \in \mathcal{Y}} p(y|x). \quad (1)$$

The corresponding confidence level  $p(y|x)$  produced by the model should ideally provide an indication of the true accuracy of the decision  $y$ . However, it is well known that machine learning models tend to be overconfident [?], [?], and hence systems using the decision  $y$

elements in set  $\mathcal{Y} = \{1, 2, \dots, C\}$ . Given any input  $x \in \mathcal{X}$ , the predictive model produces a conditional probability distribution  $p(y|x)$  over the labels  $y \in \mathcal{Y}$ . Probability  $p(y|x)$  is typically interpreted as a measure of the confidence that the model has in label  $y$  being the correct one. Conventionally, a decision  $y^*$  is obtained by selecting the label on which the model has maximum confidence, i.e.,

In this paper, we are interested in producing decisions that provide trustworthy measures of uncertainty. To this end, following the CP framework [?], our goal is to ensure that, based on communication with the devices, the server outputs

*set-valued* decisions for any new input  $x$  with *formal reliability guarantees*.

$$(1)$$

To explain, we define a mapping from an input  $x$  to a subset of the label space as  $\Gamma(x) \subseteq \mathcal{Y}$ . The corresponding confidence level  $p(y^*|x)$  produced by the model should ideally provide an indication of the true accuracy of the decision  $y^*$ . However, it is well known that machine learning models tend to be overconfident [?], [?], and hence systems using the decision  $y^*$  produced by the model cannot trust the confidence level  $p(y^*|x)$  to provide a reliable measure of the reliability of the decision. For a target reliability level  $1 - \alpha$ , with  $\alpha \in [0, 1]$ , a set decision is said to be *reliable* if the set contains the true label  $y$  with probability at least  $1 - \alpha$ , i.e., if the inequality

In this paper, we are interested in producing decisions that provide trustworthy measures of uncertainty. To this end, following the CP framework [?], our goal is to ensure that, based on communication with the devices, the server outputs

$$\Pr(y \in \Gamma(x)) \geq 1 - \alpha \quad (2)$$

The probability  $\Pr(\cdot)$  is measured with respect to the randomness of data generation and computation, i.e., we decisions for new inputs with formal reliability guarantees.

Before detailing the role of mapping functions, we observe that the reliability requirement (2) can be trivially met by choosing as a set decision the set of all possible labels, i.e., the model output distribution  $p(y|x)$ . Along with information received from the devices, the server aims at producing a set-valued decision  $\Gamma(x)$  with reliability of the set decision on the basis of the

average size of its prediction. A set decision is said to be reliable if the predictor, which is defined by label  $y$  with probability at least  $1 - \alpha$ , i.e., if the inequality

$$\Pr(y \in \Gamma(x)) \geq 1 - \alpha \quad (2)$$

where  $|\cdot|$  represents the cardinality of the argument set. The expectation in (??) is evaluated holds. The probability in (??) is evaluated with respect to the randomness of data generation with respect to the randomness of both data and communications, as in (??). and communications, as we discuss in the next subsections.

Before detailing the role of communications, we observe that the reliability requirement **B. Data Model**

(??) can be trivially met by choosing as a set decision the set of all possible labels, i.e., As mentioned in the previous subsection, we assume that the model  $p(y|x)$  is pre-trained  $\Gamma(x) = \mathcal{Y}$ , irrespective of the input  $x$ . This set predictor, while reliable, would be completely using an arbitrary training technique and an arbitrary data set. Accordingly, we do not concern uninformative. It is hence also important to evaluate the performance of the set decision ourselves with the training data and with the training process in this paper. That said, the CP on the basis of the average size of its prediction. This is known as the inefficiency of the procedure requires a data set – distinct from training data – that is used to calibrate model predictor, which is defined as

$p(y|x)$  so as to obtain a reliable set-valued prediction in the sense of condition (??) for some  $\mathbb{E} [|\Gamma(x)|]$ , **B. Data Model**

input  $x$ . In conventional CP, such data set, known as *calibration data set*, is directly available at where  $|\cdot|$  represents the cardinality of the argument set. The expectation in (??) is evaluated with respect to the randomness of both data and communications, as in (??).

instead, calibration data sets are only present at the devices. By communicating information about such data to the server, the devices can facilitate the implementation of CP-based mechanisms

**B. Data Model** to produce reliable estimates  $\Gamma(x)$  that satisfy the inequality (??) for a desired reliability level

$1 - \alpha$ . We will explain how CP works in the next section.

We assume that there are a total of  $N$  *calibration data points*, denoted as  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$ , which are equally split across all  $K$  devices. Accordingly, each device  $k$  stores  $N_k = N/K$  calibration data points, denoted as  $\mathcal{D}_k = \{(x_{i,k}, y_{i,k})\}_{i=1}^{N_k}$ . The union of all disjoint sets of data points  $\mathcal{D}_k$  across all  $K$  devices recovers the overall calibration data set, i.e.,  $\mathcal{D} = \mathcal{D}_1 \cup \dots \cup \mathcal{D}_K$ . Following the standard machine learning model, all calibration data points in data set  $\mathcal{D}$  are assumed to be generated i.i.d. from some unknown distribution  $p(x, y)$ .

In this paper, as in [2], we assume that, instead calibration data sets are only present at the devices. By communicating information about such data to the server, the devices can facilitate the implementation of CP-based mechanisms to produce reliable estimates  $\Gamma(x)$  that satisfy the inequality (??) for a desired reliability level  $1 - \alpha$ . We will explain how CP works in the next section.

**C. Communication Model** We assume that there are a total of  $N$  calibration data points, denoted as  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$ , which prior knowledge is iterated in the device CPA [2] or [3], each device has storage  $M_d$  channels  $K$  calibration devices point set denoted as  $\mathcal{D}_d = \{(x_{i,d}, y_{i,d})\}_{i=1}^{N_d}$ . In contrast to this disjoint sets of data

points. Challenging all  $K$  devices which the devices are all located on a data set  $\mathcal{D}$ , no By channels. Following the standard machine learning model, all reliable federated data points in data set  $\mathcal{D}$  are specified to be Gaussian distributed with the second distribution probability access channel using the channel uses not yet assigned. Each channel uses  $(x_k, y_k)$  is also assigned from a Gaussian distribution  $\mathcal{N}(\mu, \sigma^2)$  independently from the actual fading channel. As explained below, we will consider two types of protocols, namely orthogonal-access systems in which each device uses distinct subsets of channel uses, and non-Orthogonal protocols in which devices are simultaneously active on all channel uses.

We assume average per symbol power constraint  $P$  for each device  $k$  and the channel between symbol powers of the channel assigned to device  $k$  is denoted as  $N_0$ . In contrast, signal-to-noise ratio (SNR) is the more challenging scenario in which devices are connected to the server via noisy channels.  $\text{SNR} = \frac{P}{N_0}$ .

Accordingly, we will refer to this problem as wireless reliable federated inference.

*1) Orthogonal Multiple Access:* Time-division multiple access (TDMA) is a conventional orthogonal multiple access scheme that assigns distinct subsets of the  $T$  channel uses to the  $K$  devices. In this paper, we focus on equal allocations whereby all devices are assigned  $[T/K]$  symbol, and we focus on Gaussian-noise channels [?]. In practice, fading channels can also be approximated by Gaussian-noise channels via pre-equalization, as often done in studies involving federated learning [?], [?], [?], [?]. As we will detail below, we consider two types of protocols, namely orthogonal-access systems in which each device uses distinct subsets of channel uses, and non-orthogonal protocols in which devices are simultaneously active on all channel uses.

We assume an average per-symbol power constraint  $P$  for each device  $k$ . Accordingly, denoting the per symbol power of the channel noise as  $N_0$ , we define the signal-to-noise ratio (SNR) as

$$\text{SNR} = \frac{P}{N_0}.$$

*1) Orthogonal Multiple Access:* Time-division multiple access (TDMA) is a conventional orthogonal multiple access scheme that assigns distinct subsets of the  $T$  channel uses to the  $K$  devices. In this paper, we focus on equal allocations whereby all devices are assigned  $[T/K]$  channel uses.

**III. BACKGROUND ON CONFORMAL PREDICTION**

In this section, we provide a brief primer on CP in order to set the necessary background required by benchmarks and proposed schemes for the problem of wireless reliable federated inference.

inference, described in all the previous sections. The hypothesis testing and the presentation also include discussions about the impact of quantization on the performance of CP, which is not covered in standard references on CP. Note that the federated multiple access interest as this will be considered. CP applies to a centralized scenario where a server holding all the available calibration data, which will be assumed throughout this section. In orthogonal protocols, all devices transmit concurrently in each symbol period  $t$ . Accordingly, the signal received at the server in period  $t$  can be written as

#### A. Validation-Based Conformal Prediction

$$y_t = \sum_{k=1}^K x_{k,t} + z_t, \quad (6)$$

We focus on a practical variant of CP, known as split, inductive, or validation-based CP, that operates on a pre-defined model  $p(y|x)$  [1], [2], [3]. Given a new input  $x$ , the goal of CP is to produce a set predictor  $\Gamma(x) \in \mathcal{Y}$  with the property of satisfying the reliability condition (??) for some pre-determined target reliability level  $1 - \alpha$ .

To this end, the server is given access not only to the model  $p(y|x)$  and to a test input  $x$ , but also to a *calibration data set*  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$  consisting of  $N$  data points. As in the previous section, the  $N$  calibration data points and the test data pair  $(x, y)$  are assumed to be i.i.d. according to an unknown distribution  $p(x, y)$ . The probability in (??) is evaluated with respect to the joint distribution of calibration and test data. In the centralized setting under study here, the server builds the set predictor  $\Gamma(x)$  using the test input  $x$ , the calibration data, and the model  $p(y|x)$ . Note that the true label  $y$  is not known at the server, since it is the subject of the inference process.

To this end, we introduce the *nonconformity (NC) score function*

#### A. Validation-Based Conformal Prediction

$$s(x, y) = 1 - p(y|x). \quad (7)$$

We focus on a practical variant of CP, known as split, inductive, or validation-based CP, that operates on a pre-trained model  $p(y|x)$  [1], [2], [3]. Given a new input  $x$ , the goal of CP is to produce a set predictor  $\Gamma(x) \subset \mathcal{Y}$  with the property of satisfying the reliability condition (??) for some pre-determined target reliability level  $1 - \alpha$ . To this end, the server is given access not only to the model  $p(y|x)$  and to a test input  $x$ , but also to a calibration data set  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$  consisting of  $N$  data points. As in the previous section, the  $N$  calibration data points and the test data pair  $(x, y)$  are assumed to be i.i.d. according to an unknown distribution  $p(x, y)$ . The probability in (??) is evaluated with respect to the joint distribution of calibration and test data. Note that the upper bound is set to 1 without loss of generality since one can always re-scale a bounded NC score to fit in the range (??).

CP includes in the predicted set  $\Gamma(x)$  all labels  $y \in \mathcal{Y}$  with a NC score smaller than a given threshold  $s_\alpha$ , i.e.,

In the centralized setting under study here, the server builds the set predictor  $\Gamma(x)$  using the test input  $x$ , the calibration data, and the model  $p(y|x)$ . Note that the true label  $y$  is not known at the server, since it is the subject of the inference process.

$$\Gamma_\alpha(x) = \{y \in \mathcal{Y} : s(x, y) \leq s_\alpha\}. \quad (9)$$

As we discuss next, in the threshold noise deformed NC score function reliability level  $1 - \alpha$  in the reliability constraint (??). Dependence on parameter  $\alpha$  is accordingly added in the subscript of the set predictor  $\Gamma_\alpha(x)$ .  
 $s(x, y) = 1 - p(y|x)$ . (7)

The NC score is a measure of the loss of the model  $p(y|x)$  on the data point  $(x, y)$ . In fact, **Evaluation in the Threshold** the model assigns a low probability to example  $(x, y)$ . Other NC scores are also possible [?], [?], [?] and the methodology developed in this paper applies more broadly to any scores as long as they are non-negative and upper bounded. To evaluate the threshold  $s_\alpha$ , the server computes the NC scores (7) for all the  $N$  calibration data points, obtaining the collection  $\mathcal{S} \triangleq \{s(x_i, y_i)\}_{i=1}^N$  of NC scores. Note that multiple data points may have the same NC score, which is accordingly counted multiple times. Then, the server sets the threshold  $s_\alpha$  to be approximately equal to the  $\lceil(1 - \alpha)N\rceil$ -th smallest NC score in the set  $\mathcal{S}$  (counting possible repetitions). Intuitively, as the reliability level  $1 - \alpha$  increases, so does the threshold  $s_\alpha$ , ensuring that the predicted set (??) includes a larger number of labels.

To formalize the operation of CP, let us introduce a function that given a set  $\mathcal{S}$  produces the  $\lceil(1 - \alpha)(N + 1)\rceil$ -th smallest value in the set. Note that the smallest value is evaluated with respect to a set with cardinality  $N + 1$  and not  $N$ , as required by CP (see, e.g., [?]). For

any given collection of real numbers  $\mathcal{S} = \{s_1, \dots, s_N\}$  with possible repetitions, we denote as  $s_{(1)} \leq s_{(2)} \dots \leq s_{(N)}$  the sorted values in ascending order. Ties are broken arbitrarily. Then, the in the reliability constraint (??). Dependence on parameter  $\alpha$  is accordingly added in the desired function is defined as

subscript of the set predictor  $\Gamma_\alpha(x)$ .

$$Q_{1-\alpha}(\mathcal{S}) \triangleq \begin{cases} s_{(\lceil(1-\alpha)(N+1)\rceil)} & \text{if } \alpha \geq 1/(N+1), \\ 1 & \text{otherwise,} \end{cases} \quad (10)$$

**B. Evaluation of the Threshold**

where  $\lceil \cdot \rceil$  denotes the ceiling operation. Accordingly, function  $Q_{\text{NC}}(\mathcal{S})$  returns the  $\lceil(1 - \alpha)(N + 1)\rceil$ -th smallest value in the set as long as  $\lceil(1 - \alpha)(N + 1)\rceil \leq N$ , or equivalently  $\alpha \geq 1/(N + 1)$ , while returning the maximum value 1 otherwise. The value  $Q_{1-\alpha}(\mathcal{S})$  can also be interpreted as the  $(1 - \alpha)(N + 1)/N$ -quantile of the empirical distribution of the entries of set  $\mathcal{S}$ . In fact, the  $(1 - \alpha)(N + 1)/N$ -quantile of the empirical distribution is, by definition, the smallest number in the set  $\mathcal{S}$  that is at least as large as a fraction  $(1 - \alpha)(N + 1)/N$  of the elements in  $\mathcal{S}$ .

With function  $Q_{1-\alpha}(\cdot)$ , the CP set predictor (??) can be succinctly expressed as  
 $\Gamma_\alpha(x) = \{y \in \mathcal{Y} : s(x, y) \leq s_\alpha\} \triangleq Q_{1-\alpha}(\mathcal{S})$ . (11)

with respect to a set with cardinality  $N + 1$ , and not  $N$ , as required by CP (see, e.g., As mentioned, it can be proved that the prediction set  $\Gamma_\alpha(x)$  in (??) satisfies the reliability [?]). For any given collection of real numbers  $\mathcal{S} = \{s_1, \dots, s_N\}$  with possible repetitions, condition (??), irrespective of the accuracy of the underlying model  $p(y|x)$  and of the ground-

truth distribution  $p^*(x, y)$  of the data [?], [?].

We Quantized Conformal Prediction the sorted values in ascending order. Ties are broken arbitrarily. Then, the desired function is defined as

As discussed in the previous section, in this paper, we are concerned with decentralized settings in which calibration data is not available at the server. In such a setting, communication between devices holding the calibration data and the server is limited by the available transmission resources. As a step in the direction of accounting for limitations arising from finite communication capacity, in this subsection, we discuss a centralized CP setting in which NC scores used to evaluate the threshold  $s_\alpha^{\text{CP}}$  in (??) are constrained to take a discrete finite set of values.

$\alpha \geq 1/(N+1)$ ; while returning the maximum value 1 otherwise.

To this end, we adopt a uniform scalar quantizer in which the range  $[0, 1]$  of possible values for the NC score, by assumption (??), is divided into  $M$  equal intervals  $[S_0, S_1], [S_1, S_2], \dots, [S_{M-1}, S_M]$  empirical distribution of the entries of set  $\mathcal{S}$ . In fact, the  $(1-\alpha)(N+1)/N$ -quantile of the with  $S_0 = 0$  and  $S_M = 1$ . Given an input NC score  $x \in [0, 1]$ , the quantized output  $q(x)$  equals the upper value  $S_m$  of the interval  $(S_{m-1}, S_m]$  containing  $x$ . Accordingly, the quantization function large as a fraction  $(1-\alpha)(N+1)/N$  of the elements in  $\mathcal{S}$ .

is defined as

With function  $Q_{1-\alpha}(\cdot)$ , the CP set predictor (??) can be succinctly expressed as

$$q(x) \triangleq \begin{cases} S_1 & x \in [S_0, S_1], \\ S_m & x \in [S_{m-1}, S_m] \text{ for } m = Q_{1-\alpha}(\mathcal{S}), \\ S_M & x \in [S_M, 1]. \end{cases} \quad (12)$$

As  $\Delta$  is small enough that the server has access to the data, quantized NC scores  $s_\alpha^{\text{CP}}$  satisfies the reliability condition (??). Following the CP perspective, we define the set of labels  $y$  underling model  $p(y|x)$  and of the ground-truth distribution  $p^*(x, y)$  of the data [?], [?],  $\Gamma_\alpha^q(x) = \{y \in \mathcal{Y} : q(s(x, y)) \leq s_\alpha^{\text{CP}} \triangleq Q_{1-\alpha}(\mathcal{S}^q)\}$ , (13)

that is, as the set of labels  $y \in \mathcal{Y}$  whose quantized NC scores  $q(s(x, y))$  are no larger than the

$(1-\alpha)(N+1)$ -th smallest NC score,  $Q_{1-\alpha}(\mathcal{S}^q)$ , in the calibration set.

As discussed in the previous section, in this paper, we are concerned with decentralized

settings in which calibration data is not available at the server. In such a setting, communication between devices holding the calibration data and the server is limited by the available transmission resources. As a step in the direction of accounting for limitations arising from finite communication capacity, in this subsection, we discuss a centralized CP setting in which NC scores used to evaluate the threshold  $s_\alpha^{\text{CP}}$  in (??) are constrained to

take a discrete finite set of values.

#### D. Quantized Conformal Prediction via Empirical Quantiles

To this end, we adopt a uniform scalar quantizer in which the range  $[0, 1]$  of possible values for the NC score, by assumption (??), is divided into  $M$  equal intervals set predictor (??) can be expressed in terms of the empirical distribution of the quantized NC  $[S_0, S_1], [S_1, S_2], \dots, [S_{M-1}, S_M]$  with  $S_0 = 0$  and  $S_M = 1$ . Given an input NC score  $x \in [0, 1]$ , scores in set  $\mathcal{S}^q$ . More precisely, it can be evaluated, approximately, as the  $(1-\alpha)$ -quantile of the empirical distribution. This fact will be instrumental in the design of the proposed federated inference protocol in Sec. ??.

the Top  $\alpha$ -quantile output defined equals the upper threshold  $S_m$  off quantized NC scores  $S_q$  due to quantizing the set of input  $x$  to the NC quantization  $S^q$ , function is defined as

$$q(x) \triangleq \begin{cases} S_1 & p_m \in [S_1, S_2] \\ \frac{1}{N} \sum_{i=1}^N \mathbb{1}_{\{m_i = m\}} & \text{for } m = 2, \dots, M. \end{cases} \quad (14)$$

where  $m_i$  is the index of the quantized  $i$ -th NC score, i.e.,  $S_{m_i} = q(s(x_i, y_i))$ . We collect all  $M$  fractions into the vector

Suppose now that the server has access to the set of quantized NC scores  $\mathcal{S}^q \triangleq \{q(s(x_i, y_i))\}_{i=1}^N$ . Following the CP procedure, we define the set predictor as

$$\mathbf{p} = [p_1, \dots, p_M]^T, \quad (15)$$

which satisfies the equality  $\sum_{m=1}^M p_m = q(s(x, y))$  we have  $\mathbf{p}^q \stackrel{\text{CP}}{\leftarrow} Q_{1-\alpha}(\mathcal{S}^q)$  discussed in Sec. IV. This relies on the evaluation of the  $[(1-\alpha)(N+1)]$ -th smallest element in the set  $\mathcal{S}^q$ , i.e.,  $Q_{1-\alpha}(\mathcal{S}^q)$ . To evaluate this quantity, that is, as the set of labels  $y \in \mathcal{Y}$  whose quantized NC scores  $q(s(x, y))$  are no larger than the  $[(1-\alpha)(N+1)]$ -th smallest NC score,  $Q_{1-\alpha}(\mathcal{S}^q)$ , in the calibration set.

Since any function of the input-output pair  $(x, y)$  is a valid NC score, so is the quantized value  $q(s(x, y))$ . Therefore  $\mathbf{p}^q = \frac{N}{N+1} \left[ 0, \frac{1}{N+1}, \dots, \frac{1}{N+1} \right]^T$  satisfies the reliability condition (??). However, one should generally expect that, due to information loss caused by quantization, the size of the predicted set  $\Gamma_\alpha^q(x)$  is generally larger than that of the quantization level  $S_{m_\alpha(\mathbf{p}^q)}$ , where the index  $m_\alpha(\mathbf{p}^q)$  is obtained by evaluating the  $(1-\alpha)$ -quantile predicted set  $\Gamma_\alpha^q(x)$  obtained from the original NC score function  $s(x, y)$ .

With this definition, the  $[(1-\alpha)(N+1)]$ -th smallest element in set  $\mathcal{S}^q$  can be obtained as the predicted set  $\Gamma_\alpha^q(x)$  obtained from the original NC score function  $s(x, y)$ .

$$\text{D. Quantized Conformal Prediction via Empirical Quantiles} \quad p_m^+ \geq 1 - \alpha. \quad (17)$$

In this subsection, we make the observation that the threshold  $s^{q-\text{CP}} = Q_{1-\alpha}(\mathcal{S}^q)$  used in the set predictor (??) can be expressed in terms of the empirical distribution of the

While the conventional CP scheme reviewed in the previous section assumes that predictive quantized NC scores in set  $\mathcal{S}^q$ . More precisely, it can be evaluated, approximately, as the model and calibration data are both present at the server, in the wireless reliable federated  $(1-\alpha)$ -quantile of the empirical distribution. This fact will be instrumental in the design inference setting as explained in Sec. ??, calibration data are only available at the devices. In of the proposed federated inference protocol in Sec. ??.

this section, we first review the FedCP-QQ scheme proposed in [?], which addresses this problem

To elaborate, let us define as  $p_m \in [0, 1]$  the fraction of quantized NC scores equal to  $S_m$  by assuming noiseless links from devices to server that can support the noiseless transmission in the set of quantized NC scores  $\mathcal{S}^q$ , i.e.,

of a single real number from each device. Then, as a benchmark, we describe a direct digital wireless implementation of FedCP-QQ that accounts for the presence of noisy channels between devices and server.

where  $m_i$  is the index of the quantized  $i$ -th NC score, i.e.,  $S_{m_i} = q(s(x_i, y_i))$ . We collect all

#### A. Federated Conformal Prediction with Noiseless Communications

The FedCP-QQ scheme introduced in [?] is based on the quantile-of-quantiles (QQ) operation. Accordingly, as the detail in next section, it sets two probabilities  $\alpha$  and  $1 - \alpha$  to identify target quantiles on behalf of clients (divided into  $M$  groups) and servers respectively in the set  $\mathcal{S}^q$ , i.e.,  $Q_{1-\alpha}(\mathcal{S}^q)$ . To evaluate this

entity device modifies its periodic NC contribution  $S_k$  of the quantized  $i=1 \dots N_d$  NCs based on this adding on  $\alpha_d$  NC scores ( $N_d$  it can inputs NCs (one equal  $y_d=1$ ) in  $N_d$  maximum  $Q$  value  $(S_k)$ . This yields policy empirical distribution communicated noiselessly to the server.

The server collects all the  $p^+$ -quantiles  $\mathcal{Q}_{1-\alpha_d}^{1:K} \left[ 0, \dots, Q_{1-\alpha_d} \left( \frac{1}{N+1} \right), \dots, Q_{1-\alpha_d} \left( S_K \right) \right]^T$  from the (16) devices, and it evaluates the  $(1-\alpha_d)(K+1)/K$ -quantile of the  $K$  quantiles, i.e., With this definition, the  $\lceil (1-\alpha_d)(N+1) \rceil$ -th smallest element in set  $\mathcal{S}^q$  can be obtained as the quantization level  $S_{m_\alpha(p^+)}$ , where  $m_\alpha \triangleq \text{the index } \mathcal{Q}_{1-\alpha_d}^{1:K}(\mathbf{p}^+)$  is obtained by evaluating (18)  $(1-\alpha)$ -quantile of the empirical distribution  $\mathbf{p}^+$ , i.e.,

The set predictor of the FedCP-QQ scheme is constructed using the obtained threshold as

$$m_\alpha(\mathbf{p}^+) = \min_{\Gamma_{\alpha_d, \alpha_s}^{\text{QQ}}(x)} \left\{ m \in \{1, \dots, M\} : p_1^+ + \dots + p_m^+ \geq 1 - \alpha \right\}. \quad (17)$$

The pair of miscoverage levels  $(\alpha_d, \alpha_s)$  must be selected in order to satisfy the coverage condition (??). To this end, reference [?] proved the following result. The previous section assumes that predictive model and calibration data are both present at the server, in the wireless reliable federated inference setting as explained in Sec. ??, calibration data are only available at the devices. In this section, we first review the FedCP-QQ scheme proposed in [?], which addresses this problem by assuming noiseless links from devices to server that can support the noiseless transmission of a single real number from each device. Then, as a benchmark, we describe a direct digital wireless implementation of FedCP-QQ that accounts for the presence of noisy channels between devices and server.

**Theorem 1 (Theorem 3.2 [?]).** For any  $(\alpha_d, \alpha_s) \in [1/(N_d+1), 1] \times [1/(K+1), 1]$ , the coverage of the set predictor  $\Gamma_{\alpha_d, \alpha_s}^{\text{QQ}}(x)$  is lower bounded as

the devices. In this section, we first review the FedCP-QQ scheme proposed in [?], which

addresses this problem by assuming noiseless links from devices to server that can support

the noiseless transmission of a single real number from each device. Then, as a benchmark,

we describe a direct digital wireless implementation of FedCP-QQ that accounts for the presence of noisy channels between devices and server.

and the operation  $\sum_{I_{1,j}=n}^N$  stands for the cascade of summations that takes into account for all elements in  $I_{1,j}$  starting from  $n$  up to  $N$ , i.e.,  $\sum_{I_{1,j}=n}^N = \sum_{i_1=n}^N \sum_{i_2=n}^N \dots \sum_{i_j=n}^N$ .

#### A. Federated Conformal Prediction with Noiseless Communications

With this result, one can find a pair of miscoverage levels  $(\alpha_d, \alpha_s)$  that minimizes the lower bound  $M_{\alpha_d, \alpha_s}^{\text{QQ}}$  while satisfying the target coverage rate  $1 - \alpha$ . The optimization objective can be formulated as

Accordingly, as we detail next, it sets two probabilities  $\alpha_d$  and  $\alpha_s$  to identify target quantiles to be computed at devices and server, respectively.

$$\alpha_d, \alpha_s \in \arg \min_{\alpha_d, \alpha_s} \{M_{\alpha_d, \alpha_s}^{\text{QQ}} : M_{\alpha_d, \alpha_s}^{\text{QQ}} \leq 1 - \alpha\}. \quad (21)$$

Each device  $k$  has access to the local NC scores  $\mathcal{S}_k = \{s(x_{i,k}, y_{i,k})\}_{i=1}^{N_d}$ . Based on this If the solution of (??) is not unique, it is suggested to find the pairs with the largest value  $\alpha_s^*$  collection of NC scores, it computes the  $(1-\alpha_d)(N_d+1)/N_d$ -quantile  $Q_{1-\alpha_d}^q(\mathcal{S}_k)$ . This real and then choose among those the pair with the largest value  $\alpha_s^*$ . Efficient ways to address this positive number is then communicated noiselessly to the server.

problem are discussed in [?], which also covers the more general case in which devices have

The server collects all the quantiles  $\mathcal{Q}_{1-\alpha_d}^{1:K} = \{Q_{1-\alpha_d}^1(\mathcal{S}_1), \dots, Q_{1-\alpha_d}^K(\mathcal{S}_K)\}$  from the  $K$  different data set sizes.

devices, and it evaluates the  $(1-\alpha_d)(K+1)/K$ -quantile of the  $K$  quantiles, i.e.,

$$s_\alpha^{\text{QQ}} \triangleq Q_{1-\alpha_s}(\mathcal{Q}_{1-\alpha_d}^{1:K}). \quad (18)$$

**B. This is a Prediction of the FedCP-QQ scheme** is constructed using the obtained threshold as In this subsection, we propose a digital implementation of the FedCP-QQ scheme, which we refer to as Digital FedCP-QQ or DQQ for short. A direct implementation of the FedCP-QQ scheme requires every device  $k$  to quantize its local quantile  $Q_{1-\alpha_d}(S_k)$  in (??) before transmission in order to meet the capacity constraints on the shared noisy channel to the receiver. To this end, the device  $k$  applies the function  $q(\cdot)$  defined in (??) to quantize the local quantile  $Q_{1-\alpha_d}(S_k)$  in (??). Then, each device uses conventional digital communications to convey the quantized set predictor  $\Gamma_{\alpha_d, \alpha_s}(x)$  is lower bounded as

Specifically, to transmit the quantized local quantiles from  $K$  devices on the shared channel, we adopt a TDMA protocol whereby, as discussed in Sec. ??, the  $K$  devices are assigned with  $[K]$  channel uses (each  $d$ ). According to the probability of error for each device can be closely approximated [2], Theorem ?? provides for the cascade of summations that takes into account for all elements in  $I_{1,j}$  starting from  $i = Q\left(\frac{\lfloor T/K \rfloor C}{\sqrt{[T/K] V}} + n\right) = \sum_{i_1=n}^N \sum_{i_2=n}^N \cdots \sum_{i_j=n}^N$  in which this result, one can implement a misscoverage levels  $(\alpha_d, \alpha_s)$  that minimizes the Gaussian variable while satisfying the given target coverage rate  $1 - \alpha$ . The optimization objective can be formulated as

$$C = \frac{1}{2} \log(1 + \text{SNR}), \quad (23)$$

and the *channel dispersion*  $V_{\alpha_d, \alpha_s}$  is defined as  $\{M_{\alpha_d, \alpha_s} : M_{\alpha_d, \alpha_s} \geq 1 - \alpha\}$ .

If the solution of (??) is not unique, it is suggested to find the pairs with the largest value  $\alpha_d^*$ . Accordingly, with probability  $\alpha_d^*$ , the pairs with the largest value  $\alpha_d^*$ . As mentioned above, to address this problem, the DQ estimator [??] which is applied on the noise of quantiles that are devices have different set sizes. Note that the sets in (??) should now be evaluated by including only the correctly received quantiles from the devices.

**B. Digital Transmission Benchmark** satisfies the reliability condition (??) by Theorem 1, the impact of subsequent wireless channel transmissions that of introducing the FedCP-QQ of active devices and reinforce the a Digital of FedCP-QQ data DQQ for the short-term. This in turn FedCP-QQ increases the coverage predicted set size (??) quantize its local quantile  $Q_{1-\alpha_d}(S_k)$  in (??) before transmission in order to meet the capacity constraints on the shared noisy channel to the receiver. To this end, the device  $k$  applies the function  $q(\cdot)$  defined in (??).

In this section, we introduce the proposed Wireless Federated Conformal Prediction (WFCP) to quantize the local quantile into one of  $M$  levels. Then, each device uses conventional scheme. WFCP implements a novel combination of TBMA and over-the-air computing to communicate the empirical distribution of the quantized NC scores from the devices to the server.

## V. WIRELESS FEDERATED CONFORMAL PREDICTION

Via specifically computing it through the superposition property of the multiple access channel, served to obtain TBMA as proposed here by time division multiplexing. Distribution of the NC assigned across  $[T/K]$  devices. Based on this, the probability of error of the global quantile  $q$  is given by

Uniquely quantized in which is justified to ensure the coverage property (??). Unlike the existing FedCP-QQ scheme ( $\epsilon = Q\left(\frac{T/K + C}{\sqrt{T/K}V}\right)$ ) in the previous section, WFCP does not require devices to compute their local quantiles. This local computation, implemented by FedCP-QQ to reduce bandwidth requirements, generally results in a performance loss, since the QQ Gaussian variable, the capacity  $C$  is given by

$$C = \frac{1}{2} \log(1 + \text{SNR}), \quad (23)$$

could only be achieved by communicating separately all the quantized NC scores from devices to server. However, for a given transmission reliability level (??), this transmission would require

$$a \text{ number } T \text{ of channel uses that scale linearly with the number of calibration data points across all devices.} \quad V = \frac{\text{SNR}}{\text{SNR} + 2} \log^2 e, \quad (24)$$

Accordingly, with probability  $\epsilon$ , the transmission is unsuccessful. Assuming that the server can detect errors, the QQ estimator (??) can be applied on the subset of quantiles that are received correctly. Note that the bound in (??) should now be evaluated by including only the correctly received quantiles from the devices.

While the resulting set predictor satisfies the reliability condition (??) by Theorem 1, the impact of lost quantiles due to channel errors is that of reducing the number of active devices, and hence the amount of calibration data effectively accessible by the server. This, in turn, as it will be detailed, the main challenge in ensuring the reliability condition (??) is the determination of a suitable correction for the quantile estimated based on the noisy received signals. Finally, we provide some discussion on the trade-offs involved in the design choices, and in this section, condition (??) is the proposed WFCP's Federated Conformal Prediction (WFCP) scheme. WFCP implements a novel combination of TBMA and over-the-air computing to communicate the empirical distribution of the quantized NC scores from the devices to the server. Via over-the-air computing, thanks to the superposition property of the multiple access channel (??), the server obtains a noisy and unbiased estimate of the empirical distribution of the NC scores across all devices. Based on this estimate, the server computes an estimate of a global empirical quantile, which is judiciously selected to ensure the coverage property (??).

$$S_{m_{i,k}} = q(s(x_{i,k}, y_{i,k})). \quad (25)$$

Unlike the existing FedCP-QQ scheme reviewed in the previous section, WFCP does not require devices to compute their local quantiles. This local computation, implemented by

Further, note that in addition to the quantized NC scores, the presentation in Sec. ?? results in a performance dimensionless probability vector that collects the quantized NC scores to derive the global quantile required to implement CP on the overall calibration data set stored across all devices as reviewed in Sec. ??.

$$\mathbf{p}_k = [p_{1,k}, \dots, p_{M,k}]^T, \quad (26)$$

This result could only be achieved by communicating separately all the quantized NC scores from devices to server. However, for a given transmission reliability level (??), this transmission would require a number  $p_{m,k} = \frac{N_d}{N_d} \sum_{i=1}^{N_d} \mathbb{1}\{m_{i,k} = m\}$  of channel uses that scale linearly with the number of calibration data points across all devices.

To mitigate this loss, WFCP enables a direct estimate of the global quantile at the server without imposing bandwidth requirements that scale linearly with the number of devices.

As an intermediate goal, WFCP obtains an unbiased estimate of the global empirical distribution  $p_m$  in (??) of the quantized NC scores across all devices. To this end, we first note that we have the equality

we first detail the transmission protocol based on TBMA adopted by WFCP. Then, we describe the set predictor implemented by WFCP on the basis of the received baseband signal. As it will be detailed, the main challenge in ensuring the reliability and fairness (??) of the fraction of NC scores in the global quantization bin is equal to the corresponding fraction noisy across all devices. Finally, once such an estimate is available, WFCP can apply the procedure discussed in Sec. ?? in order to mitigate operation of centralized WFCP. As we will see, this requires a judicious adjustment of the threshold used in evaluating the predicted A set (??).

To WFCP and its related CPQ, the WFCP flow starts with the compilation of TBMA and over-the-air compiling of WFCP. At each device, the NC scores are batched by device [quantizes] separately each of the  $N_d$  codewords, using the uniform quantizer  $\mathbb{R}^M$  with  $M$  levels of Weiszfeld's The number  $M$  of than indexes of the quantization bin is produced for the  $N_d$  NC scores of symbols. We will discuss the choice of the number of quantization points in Sec. ??.

Assuming that each codeword is normalized to satisfy the energy constraint  $\|q_m(\mathbf{x}_{i,k}, \mathbf{y}_{i,k})\|$ , the orthogonality of the codewords (??), we have the equality  $\mathbf{C}^T \mathbf{C} = \mathbf{I}_{M \times M}$ . Each codeword  $\mathbf{c}_{i,k}$  is assigned to the  $m$ -th quantization bin. Furthermore, in a manner that mirrors the presentation in Sec. ??, we introduce an  $M$ -level  $S_{m,k}$  for  $m \in \{1, \dots, M\}$ , dimensional probability vector that collects the quantized NC scores at device  $k$  as

Accordingly, each device  $k$  transmits a superposition of the codewords  $\mathbf{c}_{m_{i,k}}$  that correspond to the  $N_d$  quantized NC score  $S_{m_{i,k}}$  for  $i = 1, \dots, N_d$ . We use the simplified notation  $\mathbf{c}_{i,k} = \mathbf{c}_{m_{i,k}}$ . To express the transmitted signal mathematically, let us denote as  $\mathbf{u}_{i,k} \in \mathbb{B}^{M \times 1}$  the one-hot vector

$$p_{m,k} = \frac{1}{N_d} \sum_{i=1}^{N_d} \mathbb{1}\{m_{i,k} = m\} \quad (27)$$

with the plone so tenth the fraction of a NC scores at the intermediate position. The scaled superposition of the codewords equality  $\sum_{m=1}^M p_{m,k}$  holds then be written as

As an intermediate goal, WFCP obtains an unbiased estimate of the global empirical distribution  $\mathbf{p}$  in (??) of the quantized NC scores across all devices. To this end, we first note that we have power control gain at device  $k$ , and we recall that  $\mathbf{p}_k$  in (??) is the empirical probability vector of the quantized NC scores at device  $k$ .

All devices transmit simultaneously on the shared Gaussian-noise channel (??). Accordingly, thanks to the superposition property of the multiple access channel, the received signal at the server is given by  $\mathbf{y} = \sum_{k=1}^K p_{m,k} \mathbf{c}_{i,k} + \mathbf{z}$ . Once such an estimate is available, WFCP can apply the procedure discussed in Sec. ?? in order to mimic the operation of centralized quantized CP. As we will see, this requires a judicious adjustment of the threshold used in evaluating the predicted set (??).

To obtain an estimate of distribution  $\mathbf{p}$ , WFCP leverages TBMA and over-the-air computing. With TBMA, the devices share a codebook  $\mathbf{C} = [\mathbf{c}_1, \dots, \mathbf{c}_M] \in \mathbb{R}^{M \times M}$  of  $M$  orthogonal codewords, where each codeword  $\mathbf{c}_m \in \mathbb{R}^{M \times 1}$  consists of  $M$  real symbols. The number  $M$

The server wishes to extract an estimate of the global empirical distribution  $\mathbf{p}$  of the quantized channel uses should not be larger than the available number  $T$  of symbols. We will discuss the choice of the number of quantization points in Sec. ??.

Assuming that each codeword is normalized to satisfy the energy constraint  $\|\mathbf{c}_m\|^2 = 1$ , by the orthogonality of

the codewords, we have the equality  $\mathbf{G}^T \mathbf{G} = \sum_{k=1}^K \gamma_k \mathbf{p}_k$ . Each codeword  $\mathbf{c}_m$  is assigned to the  $m$ -th quantization level  $S_m$  for  $m \in \{1, \dots, M\}$ .

From (??), we note that the server obtains a weighted sum of the local empirical distribution

Accordingly, each device  $k$  transmits a superposition of the codewords  $\mathbf{c}_{m_{i,k}}$  that correspond to the  $N_d$  quantized NC score  $S_{m_{i,k}}$  for  $i = 1, \dots, N_d$ . We use the simplified notation power control coefficients at the transmitters are set to be equal, i.e.,  $\gamma = \gamma_1 = \dots = \gamma_K$ . The  $\mathbf{c}_{i,k} = \mathbf{c}_{m_{i,k}}$ . To express the transmitted signal mathematically, let us denote as  $\mathbf{u}_{i,k} \in \mathbb{B}^{M \times 1}$  common scaling parameter  $\gamma$  must satisfy the average per-symbol transmit power constraint  $P$ . This constraint results in the inequalities

scaled superposition of the codewords transmitted by the device  $k$  can then be written as

$$\|\mathbf{x}_k\|^2 = \gamma^2 N_d^2 \|\mathbf{p}_k\|^2 \leq M_N P, \quad \text{for } k = 1, \dots, K \quad (32)$$

$$\mathbf{x}_k = \gamma_k \sum_{i=1}^{N_d} \mathbf{c}_{i,k} = \gamma_k \mathbf{C} \sum_{i=1}^{N_d} \mathbf{u}_{i,k} = \gamma_k N_d \mathbf{C} \mathbf{p}_k, \quad (29)$$

where  $\gamma_k > 0$  is a power control gain at device  $k$ , and we recall that  $\mathbf{p}_k$  in (??) is the empirical probability vector of the quantized NC scores at device  $k$ .

where  $H_2(\mathbf{p}_k)$  is the 2-Renyi entropy [?].

All devices transmit simultaneously on the shared Gaussian-noise channel (??). Accordingly, thanks to the superposition property of the multiple access channel, the received

of the empirical probability vector  $\mathbf{p}_k$ . Inequality (??) reflects the fact that a more concentrated empirical distribution, with a smaller 2-Rényi entropy, yields a stricter restriction on the choice of the transmit power. Given that, in general, no prior information is available on the distribution of the NC scores, we set the power scaling factor by considering the worst-case situation  $N_0$ , yielding the choice

$$\gamma = \frac{\sqrt{MP}}{N_d}. \quad (35)$$

### B. Estimate of the Global Empirical Distribution

With (??) in (??), the matched filtered received signal is given by

The server wishes to extract an estimate of the global empirical distribution  $\mathbf{p}$  of the quantized NC scores in  $\mathbf{w}$  from the perceived signal  $\mathbf{C}^T \mathbf{z}$ . To this end, matched filtering (36) applied by left-multiplying the received signal by matrix  $\mathbf{C}^T$ , yielding which is indeed a scaled and noisy version of the empirical probability distribution of all the  $N = KN_d$  NC scores from  $K$  devices:

$$\mathbf{w} = \mathbf{C}^T \mathbf{y} = N_d \sum_{k=1}^K \gamma_k \mathbf{p}_k + \mathbf{C}^T \mathbf{z}. \quad (31)$$

From (Predictor) note that the server obtains a weighted sum of the local empirical distribution vectors. In order to recover a noisy version of the global empirical distribution vector  $\mathbf{p}$ , the server recovers the scaled and power control coefficients at the transmitters are set to be equal, i.e.,  $\gamma_k = \gamma$ . The noisy version  $\mathbf{w}$  of the empirical distribution  $\mathbf{p}$  of all the calibration NC scores. In the ideal case of noiseless channels, i.e., with  $N_0 = 0$ , the server would have access to the empirical  $\mathbf{P}$  distribution  $\mathbf{p}$  of all the NC scores, and the quantized CP procedure in Sec. ?? could be directly applied to obtain a reliable set predictor

$$\|x_k\|^2 = \gamma N_d \|\mathbf{p}_k\|^2 \leq MP, \quad \text{for } k = 1, \dots, K \quad (32)$$

To address the availability of an estimate of vector  $\mathbf{p}$ , the proposed WFCP preprocesses or equivalently in the single inequality the scaled and noisy version of the global empirical distribution (??), and then it computes the  $(1 - \alpha_c)(N + 1)/N$ -quantile of the distribution at a corrected unreliability level  $\alpha_c$  (33), accounting for the presence of channel noise. We will demonstrate that, with a specific choice where  $H_2(\mathbf{p}_k)$  is the 2-Rényi entropy [?], of the corrected unreliability level  $\alpha_c$ , the proposed approach preserves the coverage property (??), with probability now taken over the channel noise.

First, in order to facilitate the estimate of the global empirical distribution of all the quantized NC scores (??), the server carries out two steps: (i) it rescales the vector  $\mathbf{w}$  by  $N/(\sqrt{MP}(N + 1))$ ; and (ii) it adds  $1/(N + 1)$  to the last entry of the received signal vector. Step (ii) amounts to the same operation carried out in (??) within the centralized quantized CP scheme reviewed in Sec. ??.

worst-case situation, yielding the choice

This preprocessing yields the  $M \times 1$  vector

$$\mathbf{v} = \frac{N}{\sqrt{MP}(N + 1)} \mathbf{w} + \left[ \frac{\gamma}{N_d}, \dots, \frac{1}{N + 1} \right]^T = \mathbf{p}^+ + \tilde{\mathbf{z}}, \quad (35)$$

While  $(\bar{p})$ , as defined in (??), has the empirical distribution given by the aggregated NC scores from the  $K$  devices as well as an additional NC score at the maximum quantization level  $S_M$ ; and  $\tilde{z} \triangleq N/(\sqrt{MPK(N+1)})\mathbf{C}^T\mathbf{z} \sim \mathcal{N}(0, \sigma^2\mathbf{I})$  is the effective noise vector with power

$$\sigma^2 = \frac{N^2}{MPK^2(N+1)^2}N_0 = \frac{N_d^2}{MSNR(N+1)^2} \approx \frac{1}{MSNRK^2}. \quad (36)$$

which is indeed a scaled and noisy version of the empirical probability distribution of all the  $N = KN_d$  NC scores from  $K$  devices.

Note that for a fixed number  $N$  of NC scores, the effective noise power is inversely proportional to the square of the number  $K$  of devices. This can be interpreted as a form of coherent gain due to the use of TBMA [?]. Furthermore, for a fixed number  $K$  of transmitting devices, the server recovers the scaled

smallest value of the effective noise power  $\sigma^2$  is obtained when every device sends exactly one and noisy version  $\mathbf{w}$  of the empirical distribution  $\mathbf{p}$  of all the calibration NC scores. In NC score, i.e., when  $N_d = 1$  and  $N = K$ .

Then, WFCP computes the index  $m_{\alpha_c}(\mathbf{v})$  corresponding to the  $(1 - \alpha_c)$ -“quantile” of the empirical distribution  $\mathbf{p}$  of all the NC scores, and the quantized CP procedure in Sec. ?? noisy distribution  $\mathbf{v}$ . Note that, since the vector  $\mathbf{v}$  is not normalized, and its entries may be even could be directly applied to obtain a reliable set predictor.

To address the availability of an estimate of vector  $\mathbf{p}$ , the proposed WFCP preprocesses define the set

the scaled and noisy version of the global empirical distribution (??), and then it computes the  $(1 - \alpha_c)(N\mathbf{v} + 1)/N$ -quantile of the distribution  $\mathbf{v}_m \geq \tilde{\mathbf{v}}$  uncorrected unreliability level  $\alpha_c < \alpha$ , accounting for the presence of channel noise. We will demonstrate that, with a specific choice of the corrected unreliability level  $\alpha_c$ , the proposed approach preserves the coverage

property (??) with probability now taken over the channel noise. Noting that there may be no value  $m \in \{1, \dots, M\}$  that satisfies the inequality in (??), we

define the index  $m_{\alpha_c}(\mathbf{v})$  as First in order to facilitate the estimate of the global empirical distribution of all the quantized NC scores (??) from the estimate (??), the server carries out two steps: (i) it rescales the vector (??) by  $M/N(\sqrt{MPK(N+1)})$ ; and (ii) it adds  $1/(N+1)$  to the last entry of the received signal vector. Step (ii) amounts to the same operation carried out in (??).

In the absence of noise, i.e. with  $\sigma^2 = 0$ , the index  $m_{\alpha_c}(\mathbf{v})$  corresponds to the index in (??) computed by the centralized quantized CP with  $\alpha_c = \alpha$ . This preprocessing yields the  $M \times 1$  vector

$$\mathbf{v} = \frac{N}{\sqrt{MPK(N+1)}}\mathbf{w} + \left[0, \dots, \frac{1}{N+1}\right]^T = \mathbf{p}^+ + \tilde{\mathbf{z}}, \quad (37)$$

$$\Gamma_{\alpha_c}^{\text{WFCP}}(x|\mathbf{v}) = \{y \in \mathcal{Y} : q(s(x, y)) \leq s_{\alpha_c}^{\text{WFCP}} \triangleq S_{m_{\alpha_c}(\mathbf{v})}\}. \quad (41)$$

where  $\mathbf{p}^+$ , as defined in (??), is the empirical distribution vector of the aggregated NC scores from the  $K$  devices as well as an additional NC score at the maximum quantization level  $S_M$ ; and  $\tilde{\mathbf{z}} \triangleq N/(\sqrt{MPK(N+1)})\mathbf{C}^T\mathbf{z} \sim \mathcal{N}(0, \sigma^2\mathbf{I})$  is the effective noise vector with power

$$\sigma^2 = \frac{N^2}{MPK^2(N+1)^2}N_0 = \frac{N_d^2}{MSNR(N+1)^2} \approx \frac{1}{MSNRK^2}. \quad (38)$$

**D. Optimization of the Number of Quantization Levels** the effective noise power is inversely proportional to the square of the number  $K$  of devices. This can be interpreted as a form of coherent gain due to the use of TBMA [?]. Furthermore, for a fixed number  $K$  of transmitting devices, the smallest value of  $M$  generally yields a more informative set predictor every device sends exactly one NC score, i.e., when  $N_d = 1$  and  $N = K$ .

Then, WFCP computes the index  $m_{\alpha_c}(\mathbf{v})$  corresponding to the  $(1 - \alpha_c)$  “quantile” of the noisy distribution  $\mathbf{v}$ . Note that, since the vector  $\mathbf{v}$  is not normalized, and its entries may be even negative, the index  $m_{\alpha_c}(\mathbf{v})$  does not correspond to a true quantile in general. To proceed, we define the set

$$\mathcal{M}_{\alpha_c}(\mathbf{v}) = \left\{ m \in \{1, \dots, M\} : v_1 + \dots + v_m \geq 1 - \alpha_c \right\}$$

In fact, with  $M < T$ , a simple repetition coding strategy stipulates that the devices repeat their transmissions  $R \triangleq \lfloor T/M \rfloor$  times, where  $R$  is the repetition rate. With this approach, the effective SNR, upon averaging the matched filter outputs (??), equals

$$\text{SNR}_{\text{rep}} = R \cdot \text{SNR} = \lfloor T/M \rfloor \text{SNR}.$$

Overall, the choice of the number  $M$  of quantization levels is a tension between improving the resolution of the CP set predictor, which would require increasing the value of  $M$ , and decreasing the effective noise power (??), which calls for a decrease in the value of  $M$ . In the absence of noise, i.e., with  $\sigma^2 = 0$ , the index  $m_{\alpha_c}(\mathbf{v})$  corresponds to the index in (??) computed by the centralized quantized CP with  $\alpha_c = \alpha$ .

#### E. Reliability Analysis

The index (??) is used to define the WFCP set predictor

WFCP satisfies the following reliability guarantee.

$$\Gamma_{\alpha_c}^{\text{WFCP}}(x|\mathbf{v}) = \{y \in \mathcal{Y} : q(s(x, y)) \leq s_{\alpha}^{\text{WFCP}} \triangleq S_{m_{\alpha_c}(\mathbf{v})}\}. \quad (41)$$

**Theorem 2.** Select the corrected unreliability level as

$$\alpha_c = \alpha - \frac{\sigma^2 M}{4\alpha} = \alpha - \frac{N_d^2}{4\alpha[T/M]\text{SNR}(N+1)^2}, \quad (43)$$

Then, the WFCP set predictor (??) satisfies the reliability guarantee

#### D. Optimization of the Number of Quantization Levels

$$\Pr(y \in \Gamma_{\alpha_c}^{\text{WFCP}}(x|\mathbf{v})) \geq 1 - \alpha, \quad (44)$$

The number of quantization levels,  $M$ , causes an increase in the number of channel uses where the average is taken with respect to the joint distribution of calibration and test data, as necessary for the server to recover vector  $\mathbf{v}$  in (??). On the flip side, in the ideal case of noiseless transmission, a larger value of  $M$  generally yields a more informative set predictor

thanks to the higher resolution of the NC scores.

As explained in Sec. ??, the system has access to  $T$  channel uses for transmission from devices to server. A possible choice of the number of quantization levels  $M$  would be to set Gaussian noise, in order to ensure the satisfaction of the reliability constraint (??). The correction

From  $\mathcal{D}^2 M \alpha$ , we increase  $\sigma^2 M \alpha$  with this effect to enhance channel utilization for optimal design. The fact that loss of quantization levels is proportional to the definition of  $\mathcal{D}$  implies the effective noise power with three SNR ( $\text{SNR}_{\text{eff}}$ ) effectiveness can be expressed as decreasing the SNR.

In fact, with  $M < \frac{\sigma^2 M}{4\alpha} = \frac{N_d^2}{4\alpha[T/M]\text{SNR}(N+1)^2} \approx \frac{N_d^2}{4\alpha T \text{SNR} K^2}$ , (45) repeat their transmissions  $R \triangleq \lfloor T/M \rfloor$  times, where  $R$  is the repetition rate. With this which is inversely proportional to the square of the number  $K$  of devices and to the number  $T$  of approach, the effective SNR, upon averaging the matched filter outputs (??), equals channel uses available for transmission. The dependence on  $K$  is particularly noteworthy: While  $\text{SNR}_{\text{rep}} = R \cdot \text{SNR} = \lfloor T/M \rfloor \text{SNR}$ , (42) conventional protocols like DQQ require communication resources in terms of channel uses  $T$  and over SNR, which one of the increasingly more quantization  $K$  increases. WFCP can benefit from the presence of resolution of the IC part of the  $K$ , which, in turn, decreases (32) the value of  $M/K^2$  and allows WFCP to offset a range of reliability (13) while the target becomes increasingly close to the true target  $1 - \alpha$ .

## E. Reliability Analysis VI. EXPERIMENTAL SETTINGS AND RESULTS

WFCP satisfies the following guarantee on the performance of the proposed WFCP via numerical results. We use as a benchmark the digital wireless implementation of FedCP-QQ [?], abbreviated Theorem 2. Select the corrected unreliability level as as DQQ, reviewed in Section ??.

$$\alpha_c = \alpha - \frac{\sigma^2 M}{4\alpha} = \alpha - \frac{N_d^2}{4\alpha[T/M]\text{SNR}(N+1)^2}. \quad (43)$$

*A. Setting*  
Then, the WFCP set predictor (??) satisfies the reliability guarantee

Following the experimental setting in reference [?], we use the CoronaHack data set, a public chest X-ray data set involving 5908 images classified using  $C = 3$  labels, namely normal, viral pneumonia, and bacterial pneumonia. We use  $N = 500$  data points for training predictive data, as well as the remaining 5408 image pairs are divided into  $N = 300$  points for calibration and  $N^{\text{te}} = 208$  points for testing. In the federated inference setup under study, each device holds

Proof: The proof is detailed in Appendix ??.  
 $N/K$  calibration data points, while only the server has access to the  $N^{\text{te}}$  test data points on Theorem ?? determines a correction for the threshold level  $\alpha_c < \alpha$  in the presence of which it wishes to generate reliable predictions. non-zero Gaussian-noise, in order to ensure the satisfaction of the reliability constraint (??).

The predictive model adopts the ResNet-18 architecture [?] with minor modifications. Specifically, the last layer is replaced by a linear layer with  $C = 3$  output neurons, followed by a softmax layer that outputs the conditional probability distribution  $p(y|x)$ . We train the model using the standard federated gradient descent protocol [?]. To this end, we divide the  $N^{\text{tr}}$  training examples evenly across all devices. Following federated stochastic gradient descent, the server collects and averages the local gradients from a subset of the devices that are evaluated based on

which is inversely proportional to the square of the parameter  $K$  of the stochastic gradient descent. Specifically, using a tilde-based entropy measure, the dependence while  $K$  is small is quadratic. While learning rate protocols like SGD or QPQs require the update step due to the increments of depicted in Fig. 2, another SNR-based method (e.g., SNG) deployed in both device and server. WFCP can benefit from this framework and since its performance increase is proportional to the noise reduction factor (ratio), it achieves a similar performance. WFCP also has a trainingability index that adapts to noisy channels. It is shown that the target can be directly accommodated within the proposed federated inference framework.

We adopt as performance measures the empirical coverage and empirical inefficiency, which are defined respectively as

$$\text{Empirical coverage} = \frac{1}{N^{\text{te}}} \sum_{i=1}^{N^{\text{te}}} \mathbf{1}(y_i \in \Gamma(x_i)) \quad (46)$$

and

#### A. Setting

$$\text{Empirical inefficiency} = \frac{1}{N^{\text{te}}} \sum_{i=1}^{N^{\text{te}}} |\Gamma(x_i)|. \quad (47)$$

Following the experimental setting in reference [?], we use the CoronaHack data set, a publicly available dataset containing 5908 images of the COVID-19 labels, averaged. Each example involves a bounding box of one of the 508 data points. We used 400 training points for training and 190 test points, while the remaining 508 data pairs are divided into  $N = 300$  points for calibration and  $N^{\text{te}} = 208$  points for testing. In the federated inference setup under study, each device holds  $N/K$  quantization levels, while only the server has access to the most by default. The performance of the proposed WFCP predictions is a function of the number of quantizable levels, the ResNet-18 number (Fig. 1[2]) of channels used in this study. Specifically, the softmax layer is replaced by a linear layer with optimal 3 output neurons, followed by a soft-improved effective SNR, requiring conditional probability distribution  $p(y|x)$ . We train the model using the standard federated gradient interpretation of WFCP, which simply divides the target training examples evenly across all devices. Following the consideration of the impact of channel noise, we collects and averages the local gradients from a subset of the devices that are Fig. 2. This empirical corresponds to the local training efficiency to update the model parameters and SNR stochastic gradient function. Specifically, first utilizing the loss function of WFCP, while adopting the coordinate optimality with diagonal (23) for all of 100 iterations per epoch. To obtain this goal, applying the thresholded target depicted in Fig. 1[2], the trained prediction model is also implemented in both the naive implementation of WFCP fails to meet the coverage requirements

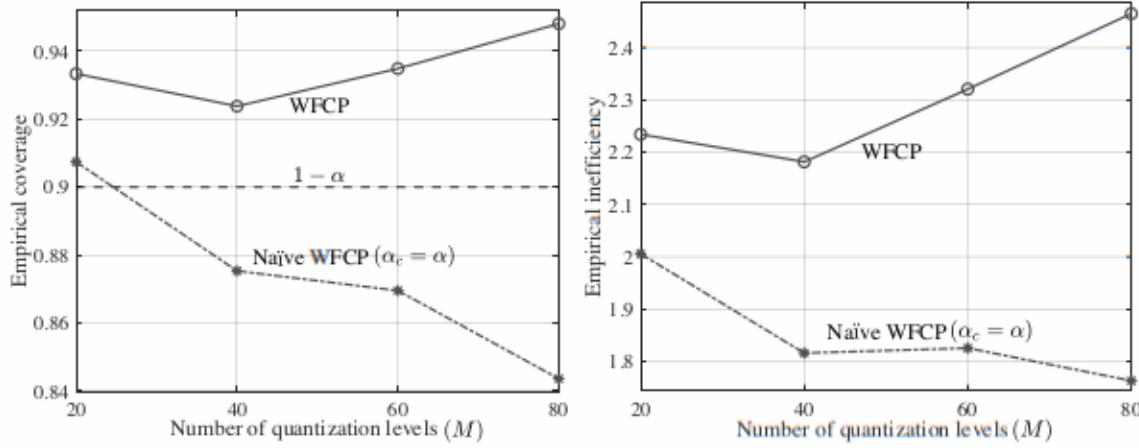


Fig. 3. Empirical coverage and empirical inefficiency of WFCP, WFCP vs. MFCP vs. DQQ versus the number of quantization levels with target unreliability level 10% of channel uses, number 20 of channel devices, and SNR  $K = 110$  dB devices, and SNR = -10 dB.

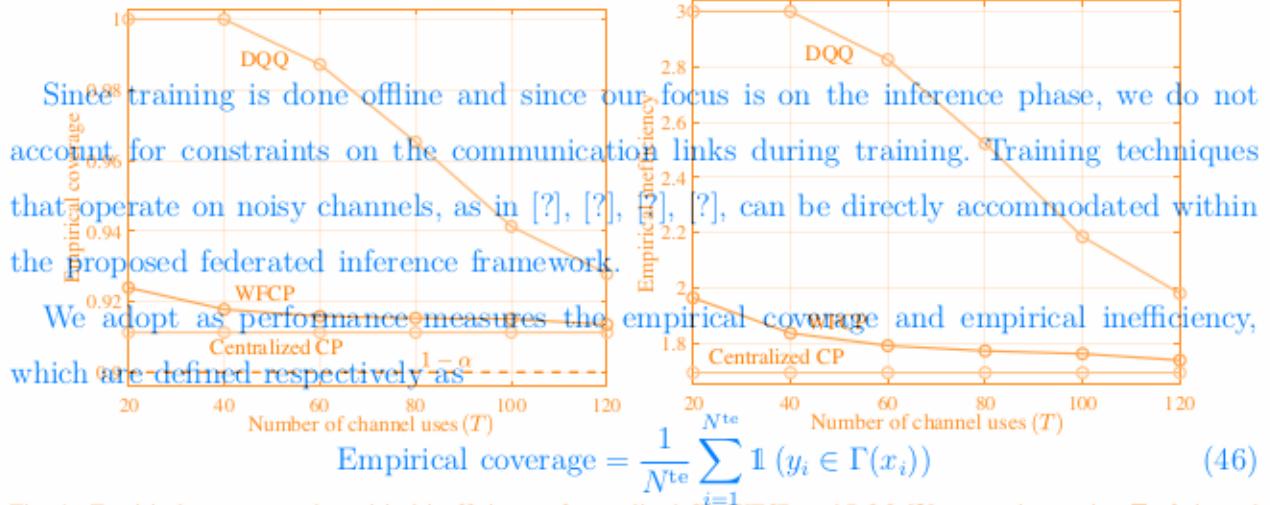


Fig. 4. Empirical coverage and empirical inefficiency of centralized CP, WFCP, and DQQ [?] versus the number  $T$  of channel uses available with target unreliability level  $\alpha = 0.1$ , number  $M = 20$  of quantization levels, number  $K = 20$  of devices, and SNR = 0 dB.

$$\text{Empirical inefficiency} = \frac{1}{N^{\text{te}}} \sum_{i=1}^{N^{\text{te}}} |\Gamma(x_i)|. \quad (47)$$

(?) was run independent 100 experiments to validate the above performance and obtain an average. Each experiment involves a random WFCP, of the 500 data points not used for training into N calibration and  $N^{\text{te}}$  test pairs, with smaller values causing a degraded performance due to an insufficient resolution and larger values being impaired by the smaller effective SNR.

### B. On the Choice of the Number of Quantization Levels

C. We start by focusing on the performance of the proposed WFCP scheme as a function of the number of quantization levels  $M$  for WFCP and DQQ [20] (Section 2). We start this study using empirical coverage and empirical inefficiency as a function of the number  $T$  of channel

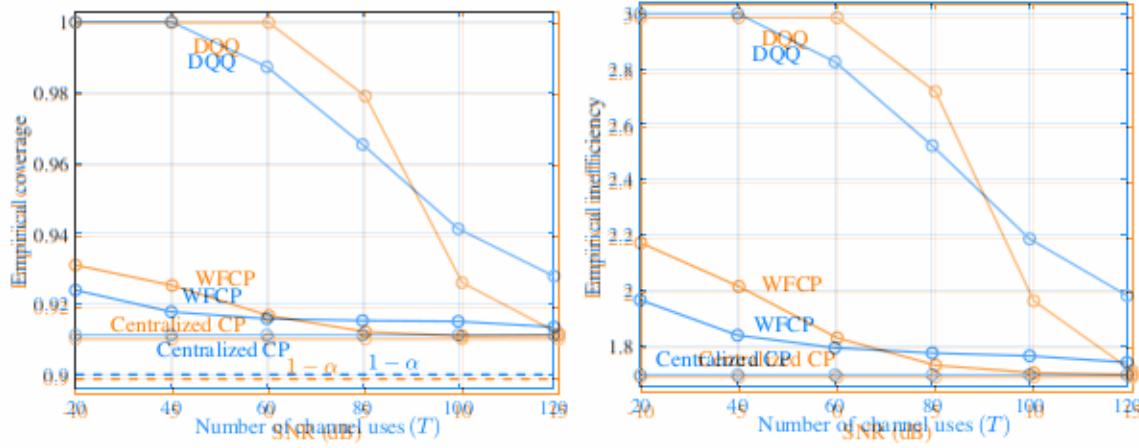


Fig. 5. Empirical coverage and empirical inefficiency of centralized CP, WFCP, and DQQ [7] versus SNR with target unreliability number  $T = 0.1$ , number  $M = 20$  of quantization levels, reliability level  $1 - \alpha_c = 0.1$  of channel number, and number  $M = 20$  of quantization levels, number  $K = 20$  of devices, and SNR = 0 dB.

uses for a fixed number  $M = 20$  of quantization levels. As seen in Fig. ??, as  $T$  increases, trade-off between improved effective SNR, requiring a smaller  $M$ , and a larger resolution, both methods maintain the target  $(1 - \alpha)$ -coverage, while offering a decreasing inefficiency, calling for a larger  $M$ . For reference, we also consider a naïve implementation of WFCP. This is because a larger  $T$  weakens the effect of channel noise by reducing the probability of which simply sets the target reliability level  $1 - \alpha_c$  in (??) to the true target  $1 - \alpha$  without error  $\epsilon$  in (??) for DQQ, and by improving the effective SNR in (??) for WFCP. The proposed considering the impact of channel noise.

WFCP consistently outperforms DQQ, yielding highly informative prediction sets, with efficiency

Fig. ?? shows empirical coverage and empirical inefficiency for  $\alpha = 0.1$ ,  $K = 10$  devices, improvements being particularly evident in the regime of limited communication resources with and SNR = -10 dB as a function of  $M$ . As a first observation, confirming Theorem ??, low number  $T$  of channel uses. As  $T$  grows sufficiently large, the performance of both schemes WFCP achieves the target coverage reliability condition (??) for all quantization levels approaches that of the centralized noiseless CP (Sec. ??).

To obtain this goal, applying the corrected target reliability level  $1 - \alpha_c$  in (??) is The performance gains of WFCP in the presence of limited communication resources are essential. In fact, as also seen in the figure, the naïve implementation of WFCP fails to further explore in Fig. ??, which evaluates the performance of WFCP and DQQ as a function meet the coverage requirements (??) as soon as  $M$  is sufficiently large, in which regime the of the SNR. As the SNR increases, the effective SNR in (??) improves along with a decrease in performance is more sensitive to the presence of channel noise. For WFCP, the optimal the correction term in (??), resulting in a more informative predicted set, which approaches the value of  $M$  in terms of inefficiency is observed to be around  $M = 40$ , with smaller values performance of the centralized CP. In a similar manner, as the SNR improves, the probability causing a degraded performance due to an insufficient resolution and larger values being of error  $\epsilon$  in (??) for DQQ decreases, thereby generating a smaller-sized predicted set, which impaired by the smaller effective SNR. approaches the performance of WFCP for SNR levels around 15 dB.

Fig. ?? evaluates the performance of WFCP and DQQ when varying the number of devices, C. Comparison between WFCP and DQQ

K. Note that the number  $N_d = 10$  of per-device calibration data points is kept fixed, so that, as

We now turn to comparing the performance of WFCP and DQQ (Sec. ??). We start by  $K$  increases, the total number of calibration data points increases. For DQQ, as the number of evaluating empirical coverage and empirical inefficiency as a function of the number  $T$  of devices increases, the inefficiency tends to increase. In fact, an increase in the number of devices channel uses for a fixed number  $M = 20$  of quantization levels. As seen in Fig. ??, as  $T$

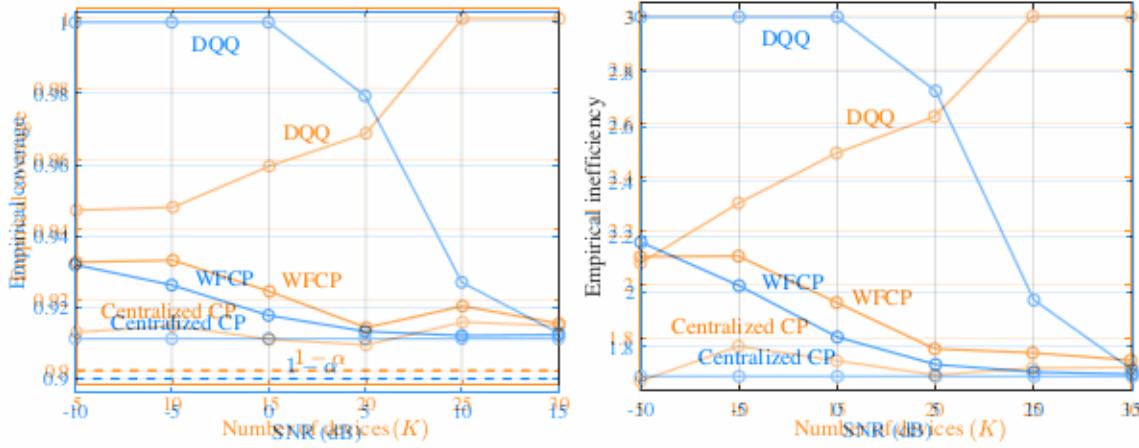


Fig. 6. Empirical coverage and inefficiency of centralized CP, WFCP, and DQQ [?] versus SNR and number of devices with  $N_d = 10$  per device calibration data points, target reliability level  $\alpha_c = 0.01$ , channel uses  $T = 40$ , quantization  $B = 20$  bits, number  $T = 60$  of channel uses, and SNR = 0 dB.

increases, both methods maintain the target  $(1 - \alpha_c)$ -coverage, while offering a decreasing leads to a higher error probability  $\epsilon$  in (??), which causes the average number of correctly received inefficiency. This is because a larger  $T$  weakens the effect of channel noise by reducing the local quantiles,  $K(1 - \epsilon)$ , to decrease probability of error  $\epsilon$  in (??) for DQQ, and by improving the effective SNR in (??) for

In stark contrast, WFCP is observed to reduce the average predicted set size as the number WFCP. The proposed WFCP consistently outperforms DQQ, yielding highly informative  $K$  of devices increases. Intuitively, this is due to the adoption of the TBMA protocol, which prediction sets, with efficiency improvements being particularly evident in the regime of allows the on-air combination of signals transmitted by all the devices. At a technical level, this limited communication resources with low number  $T$  of channel uses. As  $T$  grows sufficiently result is aligned with (??), which shows that the correction term is approximately independent large, the performance of both schemes approaches that of the centralized noiseless CP (Sec. of the number of calibration data per device and that it is inversely proportional to the square of ??).

The performance gains of WFCP in the presence of limited communication resources approaches the true level  $1 - \alpha_c$ , and the performance of WFCP approaches that of centralized CP. a function of the SNR. As the SNR increases, the effective SNR in (??) improves along with a decrease in the correction term in (??), resulting in a more informative predicted

## VII. CONCLUSIONS

set, which approaches the performance of the centralized CP. In a similar manner, as the

This paper has introduced wireless federated conformal prediction (WFCP), the first protocol SNR improves, the probability of error  $\epsilon$  in (??) for DQQ decreases, thereby generating for the deployment of federated inference via CP in shared noisy communication channels. Like a smaller-sized predicted set, which approaches the performance of WFCP for SNR levels conventional centralized CP and some of the existing federated extensions of CP for noiseless around 15 dB.

channels, WFCP provably provides formal guarantees of reliability, indicating that the predicted

Fig. ?? evaluates the performance of WFCP and DQQ when varying the number of set produced at the server contains the true output with any target probability. WFCP builds on devices,  $K$ . Note that the number  $N_d = 10$  of per-device calibration data points is kept type-based multiple access (TBMA), a communication protocol that allows the estimate of a fixed, so that, as  $K$  increases, the total number of calibration data points increases. For global histogram from distributed observations with a bandwidth that scales with the resolution

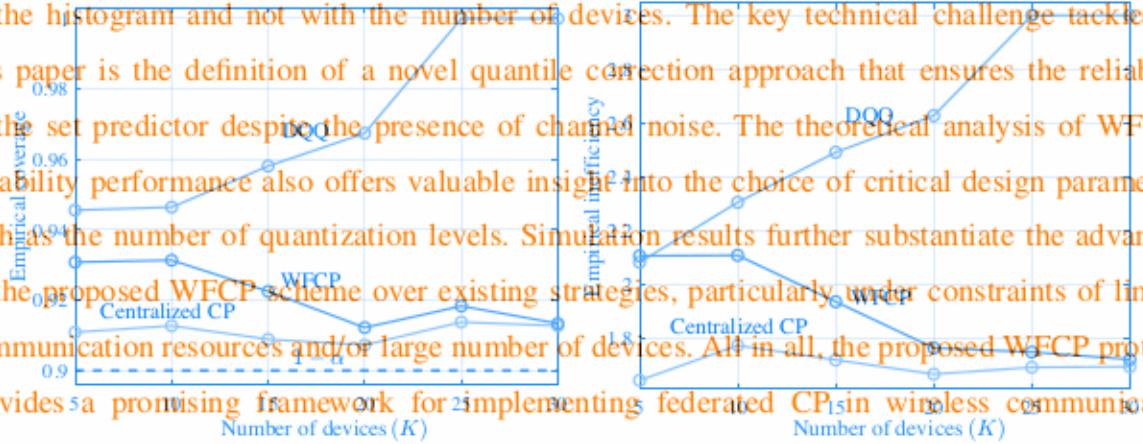
of the histogram and not with the number of devices. The key technical challenge tackled by this paper is the definition of a novel quantile correction approach that ensures the reliability of the set predictor despite the presence of channel noise. The theoretical analysis of WFCP's reliability performance also offers valuable insight into the choice of critical design parameters, such as the number of quantization levels. Simulation results further substantiate the advantage of the proposed WFCP scheme over existing strategies, particularly under constraints of limited communication resources and/or large number of devices. All in all, the proposed WFCP protocol provides a promising framework for implementing federated CP in wireless communication scenarios, thereby establishing a robust foundation for future exploration in this domain.

For future work, we suggest broadening the application of WFCP to encompass a wider range of practical scenarios. This could involve exploring the impact of heterogeneous data distributions across the devices, as done in [?], [?] for noiseless channels, as well as the performance degradation arising from imperfect channel state estimation in fading channels. DQO, as the number of devices increases, the inefficiency tends to increase. In fact, an increase in the number of devices leads to a higher error probability  $\epsilon$  in (??), which causes the average number of correctly received local quantiles,  $K(1 - \epsilon)$ , to decrease.

In stark contrast, WFCP is observed APPENDIX to reduce the average predicted set size as the number  $K$  of devices increases. Intuitively, this is due to the adoption of the TBMA protocol, which allows the on-air combination of signals transmitted by all the devices. At this section, we denote quantized calibration NC scores as  $s_i = q(s(x_i, y_i))$  for  $i = 1, \dots, N$ . A technical level, this result is aligned with (??), which shows that the correction term and quantized NC score for the true test pair  $(x, y)$  as  $s_{n+1} = q(s(x, y))$ . We also introduce two sets of  $N + 1$  NC scores. The first includes both calibration and test NC scores, i.e., is approximately independent of the number of calibration data per device and that it is inversely proportional to the square of the number of devices,  $K$ . Accordingly, as  $K$  grows, the corrected target reliability level  $1 - \alpha_c^{N+1}$  approaches the true level  $1 - \alpha$ , and the performance of WFCP approaches NC score  $s_{N+1}$  with the maximum NC score value  $S_M$ , i.e.,

$$\mathcal{S}^{\max} = \{s_i\}_{i=1}^N \cup \{S_M\}. \quad (49)$$

This paper has introduced wireless federated conformal prediction (WFCP), the first protocol for the deployment of federated inference via CP in shared noisy communication channels. Like the data point to which each NC score is assigned. Furthermore, we write as of CP the indices of the data point WFCP signed to each element for the bag  $\mathcal{S}^{\max}$ . We use the same notation for  $\mathcal{S}^{\max}$ . Based on these definitions, the set  $\mathcal{S}^{\max}$  is unambiguously identified by the bag  $\mathcal{S}^{\max}$  and by the synonym  $\pi(\mathcal{S})$  on type-based multiple access (TBMA), a communication



Finally, that all the bags estimate in a global histogram from a split (if used) in which each with an average depth that is a function of the NC resolution of the histogram quantization with the number of bins. With this definition, the key technical challenge is how to define a novel quantile correction approach that ensures the reliability of the set predictor despite the presence of channel noise. The theoretical analysis of WFCP's reliability performance also offers valuable insight into the choice of critical design parameters, such as the number of quantization levels. Simulation results further substantiate the advantage of the proposed WFCP scheme over existing strategies, particularly under constraints of limited communication resources and/or large number of devices. All in all, the proposed WFCP protocol provides a promising framework for implementing federated CP in wireless communication scenarios, thereby establishing a robust foundation for future exploration in this domain.

$$\Pr(y \in \Gamma_{\alpha_c}^{\text{WFCP}}(x|\mathbf{v})) = \mathbb{E}_{\mathcal{S}^*, \tilde{z}} [\mathbb{1}(y \in \Gamma_{\alpha_c}^{\text{WFCP}}(x|\mathbf{v}))]$$

For future work, we suggest broadening the application of WFCP to encompass a wider range of practical scenarios. This could involve exploring the impact of heterogeneous data distributions across the devices, as done in [?], or for noiseless channels, as well as the performance degradation arising from imperfect channel state estimation in fading channels. Another interesting direction for research is to devise a differentially private implementation of WFCP, potentially leveraging the idea of channel noise as a masking mechanism [?]. We have followed a standard trick of CP (see, e.g., Lemma 1 of [?]). This states that replacing any single value with the maximum value will never decrease the respective empirical quantile value. In the last equality, we have used the law of iterated expectations, which allows us to apply the expectations in the sequence as explained next.

In this section, we denote quantized calibration NC scores as  $s_i = q(s(x_i, y_i))$  for  $i = 1, \dots, N$  and quantized NC score for the true test pair  $(x, y)$  as  $s_{N+1} = q(s(x, y))$ . We also introduce two sets of  $N+1$  NC scores. The first includes both calibration and test NC scores, i.e.,

We begin by studying the inner expectation over the ordering  $\pi(\mathcal{S}^*)$  after conditioning on the bag  $\{\mathcal{S}^*\}$  and the noise vector  $\tilde{z}$ . In the following, we write  $\mathbf{p}^* = \mathbf{p}(\{\mathcal{S}^*\})$  to simplify the notation. Recall that  $S_1, \dots, S_M$  are the  $M$  quantization levels. From exchangeability of the data, while the second replaces the test NC score  $s_{N+1}$  with the maximum NC score value  $S_M$ , we have the equality (see, e.g., [?]) i.e.,

$$\Pr[S_{N+1} = S_i | \tilde{z}, \{\mathcal{S}^*\}] = p_i^*. \quad (49)$$

For such a genie-aided set  $\mathcal{S}^* = \{s_1, \dots, s_{N+1}\}$  that has access to the test NC score, we use the bag notation  $\{\mathcal{S}^*\}$  to refer to a set of numerical values of the NC scores, which excludes the identity of the data point to which each NC score is assigned. Furthermore,

It follows that  $(\mathcal{S}^*)$  have the series of the quantization points assigned to each element in the bag  $\{\mathcal{S}^*\}$ .

We use the same notation for  $\mathcal{S}^{\max}$ . Based on these definitions, the set  $\mathcal{S}^*$  is unambiguously identified by the bag  $\{\mathcal{S}^*\}$  and by the assignment  $\pi(\mathcal{S}^*)$ .

Finally, given the bag  $\{\mathcal{S}^*\}$ , we introduce the  $\sum_{i=1}^{m_{\alpha_c}(\mathbf{p}^* + \tilde{\mathbf{z}})} 1$  vector as  $\mathbf{p}(\{\mathcal{S}^*\})$ , in which each  $m$ -th entry represents the fraction of NC scores in  $\{\mathcal{S}^*\}$  equal to quantization level  $S_m$ . With this definition, the vector  $\mathbf{p}^+$  in (??) can be equivalently defined as

$$\geq \min \left\{ 1, 1 - \alpha_c - \sum_{i=1}^{m_{\alpha_c}(\mathbf{p}^* + \tilde{\mathbf{z}})} z_i \right\}, \quad (54)$$

$$\mathbf{p}^+ = \mathbf{p}(\{\mathcal{S}^{\max}\}), \quad (50)$$

where the inequality follows from the definition of  $m_{\alpha_c}(\mathbf{p}(\{\mathcal{S}^{\max}\}) + \tilde{\mathbf{z}})$  in (??) and the  $\min\{\cdot\}$  operator is introduced to account for the case  $m_{\alpha_c}(\mathbf{p}^* + \tilde{\mathbf{z}}) = M$ .

$$\mathbf{v} = \mathbf{p}(\{\mathcal{S}^{\max}\}) + \tilde{\mathbf{z}}. \quad (51)$$

### Second step: Bounding $\mathbb{E}_{\tilde{\mathbf{z}}|\{\mathcal{S}^*\}}[\cdot]$

In (??) and (??), we have used the fact that histograms do not depend on the ordering of the defined set. Recall that we are interested in finding a lower bound on the probability

$$\begin{aligned} \mathbb{E}_{\tilde{\mathbf{z}}|\{\mathcal{S}^*\}} \mathbb{E}_{\pi(\mathcal{S}^*)|\tilde{\mathbf{z}}, \{\mathcal{S}^*\}} [\mathbb{1}(y \in \Gamma_{\alpha_c}^{\text{WFCP}}(x|\mathbf{v}))] &\geq \mathbb{E}_{\tilde{\mathbf{z}}|\{\mathcal{S}^*\}} \left[ \min \left\{ 1, 1 - \alpha_c - \sum_{i=1}^{m_{\alpha_c}(\mathbf{p}^* + \tilde{\mathbf{z}})} \tilde{z}_i \right\} \right] \\ \Pr(y \in \Gamma_{\alpha_c}^{\text{WFCP}}(x|\mathbf{v})) &= \mathbb{E}_{\mathcal{S}^*, \tilde{\mathbf{z}}} [\mathbb{1}(y \in \Gamma_{\alpha_c}^{\text{WFCP}}(x|\mathbf{v}))] \\ &= \mathbb{E}_{\mathcal{S}^*, \tilde{\mathbf{z}}} \left[ \mathbb{1}_{\frac{1}{2} \leq \frac{S_{m_{\alpha_c}(\mathbf{p}(\{\mathcal{S}^*\}) + \tilde{\mathbf{z}})}}{2} \sum_{i=1}^{m_{\alpha_c}(\mathbf{p}^* + \tilde{\mathbf{z}})} \tilde{z}_i} - \left| \alpha_c + \sum_{i=1}^{m_{\alpha_c}(\mathbf{p}^* + \tilde{\mathbf{z}})} \tilde{z}_i \right| \right], \\ &\geq \mathbb{E}_{\mathcal{S}^*, \tilde{\mathbf{z}}} [\mathbb{1}(s_{N+1} \leq S_{m_{\alpha_c}(\mathbf{p}(\{\mathcal{S}^*\}) + \tilde{\mathbf{z}})})] \end{aligned} \quad (55)$$

$$= \mathbb{E}_{\tilde{\mathbf{z}}|\{\mathcal{S}^*\}} \mathbb{E}_{\pi(\mathcal{S}^*)|\tilde{\mathbf{z}}, \{\mathcal{S}^*\}} [\mathbb{1}(s_{N+1} \leq S_{m_{\alpha_c}(\mathbf{p}(\{\mathcal{S}^*\}) + \tilde{\mathbf{z}})})], \quad (52)$$

in which we have used the identity  $\min\{x, y\} = y + \frac{x-y-|x-y|}{2}$  to obtain the last equality.

We now note that the index  $m_{\alpha_c}(\mathbf{p}^* + \tilde{\mathbf{z}})$  depends on  $\tilde{\mathbf{z}}$  and that, once conditioned on  $\{\mathcal{S}^*\}$  it depends only on the realization of the sequence  $\tilde{z}_1, \dots, \tilde{z}_{m_{\alpha_c}(\mathbf{p}^* + \tilde{\mathbf{z}})}$ . Therefore  $m_{\alpha_c}(\mathbf{p}^* + \tilde{\mathbf{z}})$  is a single value with the maximum value will never decrease the respective empirical quantile stopping time for the sequence  $\tilde{z}_1, \tilde{z}_2, \dots$ , and we are allowed to invoke first Wald's identity [?].

In the last equality, we have used the law of iterated expectations, which allows us to obtain

$$\mathbb{E}_{\tilde{\mathbf{z}}|\{\mathcal{S}^*\}} \left[ \sum_{i=1}^{m_{\alpha_c}(\mathbf{p}^* + \tilde{\mathbf{z}})} \tilde{z}_i \right] = \mathbb{E}_{\tilde{\mathbf{z}}|\{\mathcal{S}^*\}} [m_{\alpha_c}(\mathbf{p}^* + \tilde{\mathbf{z}})] \mathbb{E}[\tilde{z}_1] = 0, \quad (56)$$

First step: Bounding  $\mathbb{E}_{\mathcal{S}^*, \tilde{\mathbf{z}}|\{\mathcal{S}^*\}}[\cdot]$  where the last equality follows from the noise being zero mean. Furthermore, from Wald's second identity [?], we have

We begin by studying the inner expectation over the ordering  $\pi(\mathcal{S}^*)$  after conditioning on the bag  $\{\mathcal{S}^*\}$  and the noise vector  $\tilde{\mathbf{z}}$ . In the following, we write  $\mathbf{p}^* = \mathbf{p}(\{\mathcal{S}^*\})$  to simplify the notation. Recall that  $(S_1, \dots, S_M)$  are the  $M$  quantization levels. From exchangeability of the data, we have the equality (see, e.g., [?])

$$\Pr[s_{N+1} = S_i | \tilde{\mathbf{z}}, \{\mathcal{S}^*\}] = p_i^*. \quad (53)$$

It follows that applying Jensen's inequality  $\mathbb{E}[x]^2 \leq \mathbb{E}[x^2]$ , we can further bound (??),

$$\begin{aligned} \mathbb{E}_{\pi(\mathcal{S}^*)|\tilde{\mathbf{z}}|\mathcal{S}^*} [\mathbb{1}(s_{N+1} \leq S_{m_{\alpha_c}(\mathbf{p}^*+\tilde{\mathbf{z}})})] &= \Pr [s_{N+1} \leq S_{m_{\alpha_c}(\mathbf{p}^*+\tilde{\mathbf{z}})} | \tilde{\mathbf{z}}, \mathcal{S}^*] \\ &= 1 + \frac{1}{2} \mathbb{E}_{\tilde{\mathbf{z}}|\mathcal{S}^*} \left[ -\alpha_c - \sum_{i=1}^{m_{\alpha_c}(\mathbf{p}^*+\tilde{\mathbf{z}})} \tilde{z}_i - \left| \alpha_c + \sum_{i=1}^{m_{\alpha_c}(\mathbf{p}^*+\tilde{\mathbf{z}})} \tilde{z}_i \right| \right] \\ &\geq 1 - \frac{\alpha_c}{2} - \frac{1}{2} \sqrt{\mathbb{E}_{\tilde{\mathbf{z}}|\mathcal{S}^*} \left[ \alpha_c^2 + 2\alpha \min_{i=1}^{m_{\alpha_c}(\mathbf{p}^*+\tilde{\mathbf{z}})} \left| \tilde{z}_i \right| - \left( \alpha_c + \sum_{i=1}^{m_{\alpha_c}(\mathbf{p}^*+\tilde{\mathbf{z}})} \tilde{z}_i \right)^2 \right]}, \end{aligned} \quad (54)$$

where the inequality follows from the definition of  $m_{\alpha_c}(\mathbf{p}^*(\mathcal{S}^{\max}) + \tilde{\mathbf{z}})$  in (??) and (58).  $\min\{\cdot\}$  operator is introduced to account for the case  $m_{\alpha_c}(\mathbf{p}^* + \tilde{\mathbf{z}}) = M$ . Accordingly, we have

$$\text{Second step: Bounding } (\mathbb{E}_{\tilde{\mathbf{z}}|\mathcal{S}^*}[\mathbb{1}]) (y \in \Gamma_{\alpha_c}^{\text{WFCP}}(x|\mathbf{v})) \geq 1 - \frac{\alpha_c}{2} - \frac{1}{2} \sqrt{\alpha_c^2 + \sigma^2 M}. \quad (59)$$

We now marginalize over the noise vector  $\tilde{\mathbf{z}}$  given the bag  $\mathcal{S}^*$ . Given the bag  $\mathcal{S}^*$  we have

**Final step: Bounding  $(\mathbb{E}_{\mathcal{S}^*}[\cdot])$**

The final step follows directly from (??), in which the lower bound does not depend on the bag  $\mathcal{S}^*$ . This gives that

$$\mathbb{E}_{\tilde{\mathbf{z}}|\mathcal{S}^*} \mathbb{E}_{\pi(\mathcal{S}^*)|\tilde{\mathbf{z}}|\mathcal{S}^*} [\mathbb{1}(y \in \Gamma_{\alpha_c}^{\text{WFCP}}(x|\mathbf{v}))] \geq \mathbb{E}_{\tilde{\mathbf{z}}|\mathcal{S}^*} \left[ \min \left\{ 1, 1 - \alpha_c - \sum_{i=1}^{m_{\alpha_c}(\mathbf{p}^*+\tilde{\mathbf{z}})} \tilde{z}_i \right\} \right]$$

$$\Pr (y \in \Gamma_{\alpha_c}^{\text{WFCP}}(x|\mathbf{v})) \geq 1 - \frac{\alpha_c}{2} - \frac{1}{2} \sqrt{\alpha_c^2 + \sigma^2 M}, \quad (60)$$

Therefore, to satisfy the target coverage rate  $1 - \alpha$ , we set the corrected unreliability level as

$$\alpha_c = \alpha - \frac{\sigma^2 M}{4\alpha}. \quad (55)$$

in which we have used the identity  $\min\{x, y\} = y + \frac{x-y-|x-y|}{2}$  to obtain the last equality.

We now note that the index  $m_{\alpha_c}(\mathbf{p}^* + \tilde{\mathbf{z}})$  depends on  $\tilde{\mathbf{z}}$  and that, once conditioned on  $\mathcal{S}^*$  it depends only on the realization of the sequence  $\tilde{z}_1, \dots, \tilde{z}_{m_{\alpha_c}(\mathbf{p}^* + \tilde{\mathbf{z}})}$ . Therefore  $m_{\alpha_c}(\mathbf{p}^* + \tilde{\mathbf{z}})$  is a stopping time for the sequence  $\tilde{z}_1, \tilde{z}_2, \dots$ , and we are allowed to invoke first Wald's identity [?] to obtain

$$\mathbb{E}_{\tilde{\mathbf{z}}|\mathcal{S}^*} \left[ \sum_{i=1}^{m_{\alpha_c}(\mathbf{p}^*+\tilde{\mathbf{z}})} \tilde{z}_i \right] = \mathbb{E}_{\tilde{\mathbf{z}}|\mathcal{S}^*} [m_{\alpha_c}(\mathbf{p}^* + \tilde{\mathbf{z}})] \mathbb{E}[\tilde{z}_1] = 0, \quad (56)$$

where the last equality follows from the noise being zero mean. Furthermore, from Wald's second identity [?], we have

$$\mathbb{E}_{\tilde{\mathbf{z}}|\mathcal{S}^*} \left[ \left( \sum_{i=1}^{m_{\alpha_c}(\mathbf{p}^*+\tilde{\mathbf{z}})} \tilde{z}_i \right)^2 \right] \leq \sigma^2 \mathbb{E}_{\tilde{\mathbf{z}}|\mathcal{S}^*} [m_{\alpha_c}(\mathbf{p}^* + \tilde{\mathbf{z}})] \leq \sigma^2 M, \quad (57)$$