

# Neural Network Quantization in Federated Learning at the Edge

Nicola Tonellotto<sup>b</sup>, Alberto Gotta<sup>a</sup>, Franco Maria Nardini<sup>a</sup>, Daniele Gadler<sup>d</sup>, Fabrizio Silvestri<sup>c</sup>

<sup>a</sup>*ISTI-CNR, Pisa, Italy*

<sup>b</sup>*University of Pisa, Pisa, Italy*

<sup>c</sup>*“Sapienza”, University of Rome*

<sup>d</sup>*ONE LOGIC GmbH, Frankfurt, Germany*

---

## Abstract

The massive amount of data collected in the Internet of Things (IoT) asks for effective, intelligent analytics. A recent trend supporting the use of Artificial Intelligence (AI) solutions in IoT domains is to move the computation closer to the data, i.e., from cloud-based services to edge devices. Federated learning (FL) is the primary approach adopted in this scenario to train AI-based solutions. In this work, we investigate the introduction of quantization techniques in FL to improve the efficiency of data exchange between edge servers and a cloud node. We focus on learning recurrent neural network models fed by edge data producers using the most widely adopted neural networks for time-series prediction. Experiments on public datasets show that the proposed quantization techniques in FL reduces up to 19× the volume of data exchanged between each edge server and a cloud node, with a minimal impact of around 5% on the test loss of the final model.

*Keywords:* Federated Learning, Artificial Neural Networks, Quantization, Internet of Things

---

## 1. Introduction

Internet of Things (IoT) is a disruptive technology that is pervasively extending the concept of data collection to everything around us through IoT devices, thus leading to a huge growth of the Internet traffic. Cisco Annual Internet Report forecasts that “*the share of Machine-To-Machine (M2M) connections will grow from 33 percent in 2018 to 50 percent by 2023. There will be 14.7 billion M2M connections by 2023*”. The huge amount of collected data enables the adoption of Artificial intelligence (AI) analytics currently provided by centralized cloud-based services. However, cloud-based solutions may raise major concerns among users about the privacy of online services. A recent approach to mitigate this issue is to decentralize the computation where data is, i.e., on personal IoT devices connected to the Internet. Edge computing [1] extends Cloud computing, according to the above decentralized approach: data processing and storage capabilities are not

---

*Email addresses:* nicola.tonellotto@unipi.it (Nicola Tonellotto), alberto.gotta@isti.cnr.it (Alberto Gotta), francomaria.nardini@isti.cnr.it (Franco Maria Nardini), daniele.gadler@yahoo.it (Daniele Gadler), fabrizio.silvestri@uniroma1.it (Fabrizio Silvestri)

13 exclusive characteristics of centralized data centers, but an additional layer, called *edge*,  
14 is placed in the middle between the Cloud and the IoT devices. This edge layer allows  
15 for storing data and executing applications on edge servers directly connected with IoT  
16 devices. Moreover, edge computing allows to preserve the confidentiality of private user  
17 data, which are outsourced to edge servers for storage and computational processes. A  
18 privacy perimeter is defined in the edge - hereinafter called *privacy domain* - wherein ac-  
19 cess control, authentication, encryption, and secure computation are performed. In this  
20 scenario, the recent introduction of *artificial-intelligence-as-a-service* [2] tackles signifi-  
21 cant innovations across all the industrial sectors in particular in the *artificial intelligence*  
22 *of things* (AIoT) [3, 4], where the data required to train AI solutions are kept local to  
23 the device without disclosing private or sensitive data.

24 In medical, industrial, and social IoT scenarios, AIoT solutions must analyze very  
25 large and dynamic time-series. Traditional time-series forecasting techniques, such as an  
26 autoregressive integrated moving average, rely on highly manually-tuned parameters [5]  
27 and are ill-equipped to learn long-range dependencies [4]. Recurrent Neural Networks  
28 (RNNs) have been successfully applied to tasks based on time-series data, such as stock  
29 forecasting, social behavior analysis, and natural language processing (NLP). Federated  
30 learning (FL) is a leading approach for training RNNs adopted in AIoT solutions [6].  
31 In FL, the computation involved in the training of an AI solution, e.g., a Neural Net-  
32 work (NN), is moved closer to where data are produced. FL can naturally apply to the  
33 IoT-Edge-Cloud scenario. This scenario may require preserving data privacy and data  
34 ownership. In an FL scenario applied to edge computing: i) every edge server receives a  
35 partially trained NN model from a cloud node, ii) every edge server performs additional  
36 training on data provided by the respective IoT devices to refine the previous model with-  
37 out disclosing any private data, iii) after the local training ends, the refined local model  
38 is sent back to the cloud node, and iv) the cloud node collects all the locally-trained  
39 models, generates a new global model, and broadcasts the global model back to the edge  
40 servers for a new round of local training. By doing so, the complex fusion of Machine  
41 Learning (ML) models on a cloud node is decoupled from the storage of training data on  
42 edge servers to preserve user data ownership and privacy. Yet, IoT devices are relieved of  
43 computationally intensive tasks, like AI services, allowing them to either publish data or  
44 execute a limited number of tasks. Learning state-of-the-art NN models in an FL scenario  
45 is challenging because of the models' size that range from hundreds of MBs to several  
46 GBs [7]. Hence, transferring such models from a cloud node to edge servers and vice-versa  
47 would lead to prolonged data exchange, high data transfer costs, and energy drain. In  
48 this work, we investigate the introduction of quantization techniques in the FL scenario  
49 to improve the efficiency of data exchange between edge servers and a cloud node. We  
50 focus on learning accurate NN models in an FL scenario with the model training steps  
51 performed on edge servers and the model aggregation performed on cloud nodes.

52 In detail, the novel and unpublished contributions of this article are the following:

- 53 • we propose the application of quantization techniques to the FL scenario by defining  
54 of the *Federated Learning with Quantization* (FLQ);
- 55 • we discuss how NN quantization techniques could apply to the FLQ scenario by  
56 introducing two new FL algorithms, namely FLQ and  $\Delta$ FLQ.
- 57 • we provide a comprehensive analysis of the performance of our FLQ and  $\Delta$ FLQ

58 algorithms with different quantization techniques in terms of i) the effectiveness of  
59 the learned NN model and ii) the data reduction attained on the public WikiText-2  
60 text dataset;

- 61 • we further assess the effectiveness and data efficiency performance of our FLQ  
62 and  $\Delta$ FLQ algorithms with the best quantization techniques on different public  
63 datasets and NN models, namely the MNIST image dataset with a Convolutional  
64 Neural Network (CNN) applied to the task of image classification, and the BAR  
65 CRAWL sensor dataset with an RNN applied to a regression task;
- 66 • we discuss the performance of the proposed  $\Delta$ FLQ algorithm with up to 10 edge  
67 servers and evaluate our algorithm’s performance in the presence of faults.

68 Our experimental results show that the application of quantization techniques to FL  
69 allows us to significantly reduce the total data exchanged between each edge server and  
70 a cloud node (up to 19 $\times$ ) with a minimal impact on the test loss of the final NN model  
71 (around 5%).

72 The remainder of the article is organized as follows. After an overview of the current  
73 state of the art in Section 2, we outline the NN quantization schemes and the proposed  
74 FLQ and  $\Delta$ FLQ algorithms in Section 3. We provide a thorough explanation of the  
75 experimental evaluation of our approach in Section 4, which includes experiments on  
76 multiple edge servers (up to 10) and a discussion of the robustness of our best approach  
77 in the presence of faults. Then, we discuss the comprehensive experimental evaluation of  
78 the quantization schemes and our FLQ and  $\Delta$ FLQ algorithms in Section 5. Section 5.4  
79 reports additional results on training networks for different datasets, containing namely  
80 image and sensor data. Finally, Section 7 draws the main conclusions and discusses future  
81 work.

## 82 2. Related Work

83 AI solutions and, in particular, NNs have been recently investigated in different IoT  
84 scenarios. Recently, Shanthamallu *et al.* survey AI methods and their applications to  
85 the IoT world [3], while Mohammadi *et al.* describe state-of-the-art methods in AI for  
86 IoT, big data, and streaming analytics [4]. Li *et al.* introduce a layered AI approach  
87 running on edge IoT devices [8], and Tang *et al.* describe methods to enable AI in IoT  
88 devices [9]. More recently, Lu *et al.* propose to apply federated learning in industrial  
89 IoT [10]. NNs are the current state-of-the-art model for such AIoT scenarios. Deep NNs  
90 are made up of multiple hidden layers of neurons connected through weighting matrices  
91 trained to accurately approximate a given objective function. Among them, convolutional  
92 neural networks (CNNs) are used to process image-based data, while RNNs are used to  
93 process variable-length sequences. RNNs are suitable for tasks involving time-series and  
94 time-segmented tasks. The Long-Short-Term-Memory (LSTM) is an RNN with long  
95 short-term memory blocks that consist of memory cell units. These memory cell units  
96 let the LSTM remember the state values for an arbitrarily long time sequence. LSTM  
97 networks have been successfully applied to applications with sequential data such as time  
98 series prediction [11], NLP [12], and social behavior analysis [13].

99 Deep NNs are over-parameterized to ease the training process. For example, AlexNet [14]  
100 has 60M parameters to be learned. NN pruning [15] consists of removing weights in a

101 neural network to reduce the storage requirements of the network parameters. NN prun-  
 102 ing is orthogonal to our proposed approach since every edge server can deploy its pruning  
 103 solution. NN quantization is another approach proposed to efficiently train NNs. This ap-  
 104 proach constrains the precision of floating-point 32-bit weights, activation values, and/or  
 105 gradients used in the training procedure to a fixed-point representation, using  $k < 32$   
 106 bits. Note that NN quantization only affects the massive computations performed to cal-  
 107 culate the updates of the matrix weights. In contrast, the matrix weights are maintained  
 108 at full precision during the training procedure. In this line of research, DoReFa-Net  
 109 quantizes the weights of a CNN to 1 bit, activations to 2 bits, and gradients to 6 bits to  
 110 preserve a high model accuracy [16]. Different quantization approaches for CNNs, i.e.,  
 111 one-bit [17] and multi-bit quantization [18] for limited-precision training, have also been  
 112 proposed. State-of-the-art methods for uniform quantization are XNOR-Net [19] and  
 113 binary weight nets [20], which propose a binary quantization mapping weights to  $-1$  and  
 114  $+1$  and replace operations on the weights with more efficient bit-wise operations. More  
 115 recently, ternary weight nets [21] have also been introduced, allowing weights to be zero.  
 116 Zhou *et al.* propose incremental NN quantization [22]: this approach iteratively converts  
 117 any pre-trained full-precision CNN into a low-precision version, whose weights are con-  
 118 strained to be either powers of two or zeros. Xu *et al.* propose to quantize the weights  
 119 of a full precision NN model to binary or ternary weights by leveraging an alternating  
 120 optimization approach applied at training time [23]: the accuracy loss of the resulting  
 121 model followed to be negligible both for binary and ternary quantization. Ardakani *et al.*  
 122 propose another binary and ternary quantization approach, where weights are sampled  
 123 from a Bernoulli distribution, and the obtained values are regularized [24]. While most  
 124 of the previous techniques applied to CNNs, very few of them applied to RNNs. As far  
 125 as RNN quantization is concerned, the two proposed techniques are alternating multi-bit  
 126 quantization [23] and Bernoulli sampling regularization [24]. NN quantization is orthog-  
 127 onal to our proposed solution since every edge server can deploy its quantization locally  
 128 to speed up the local model training. Nevertheless, we will exploit weight quantization  
 129 schemes to reduce the size of the NN (see Section 3).

130 FL and the corresponding decentralized training of NNs have been recently proposed  
 131 by McMahan *et al.* [25]: the main motivating example for FL arises when the training data  
 132 comes from users' interaction with mobile applications. Konečný *et al.* [26] propose an  
 133 FL collaborative approach, which enables smartphones to collaboratively learn a shared  
 134 prediction model while keeping all the training data on the device, thus decoupling the  
 135 ability to train an AI model from the need to store the data in the Cloud. The authors  
 136 describe two ways to reduce the uplink communication costs: i) using a smaller number  
 137 of parameters for the model and ii) compressing the parameters by using a combination  
 138 of quantization, random rotations, and sub-sampling before sending the model or the  
 139 model update to the cloud node. Experiments on both CNNs and RNNs show only  
 140 the accuracy of the proposed method, while the reduction in communication overhead is  
 141 only theoretically calculated in two orders of magnitude. More recently, Reisizadeh *et al.*  
 142 propose a communication-efficient FL method with periodic averaging and quantization  
 143 [27]: local CNN models are updated at each local device and only periodically averaged  
 144 at the cloud node; moreover, only a fraction of devices participate at each round of the  
 145 training; finally, the local devices quantize the parameters of their local models, before  
 146 uploading them to the cloud node.

147 Differently from [27], we introduce an intermediate computational layer, represented

148 by the edge servers, which performs the quantization of the NN model, which, being  
149 a computationally intensive task, cannot be deployed on constrained IoT devices. To  
150 this aim, a publish-subscribe paradigm is applied within each privacy domain to locally  
151 transfer data from publisher IoT devices to the relative subscribed edge server. Moreover,  
152 we propose to quantize both the NN model broadcasted by the cloud node and the NN  
153 models sent by the edge servers to significantly reduce the background traffic related to  
154 the NN model training.

### 155 3. Proposed Framework

156 In this section, we present our proposed framework to perform FL of NNs, where: (i)  
157 data are produced by IoT nodes and collected by the relative edge servers belonging to  
158 the private domain, and (ii) training is jointly performed both on edge servers and in  
159 the Cloud. Edge servers perform local training in their isolated privacy domains, while  
160 the cloud node performs central aggregation of the locally learned models. To reduce  
161 Edge-Cloud data exchange, we propose a novel FL algorithm with quantization (FLQ),  
162 aimed at quantizing the NN weights during the FL procedure, i.e., when the NN models  
163 are transmitted from the edge servers to the cloud node and vice-versa. In Section 3.1  
164 the communication paradigm underlying the proposed AI/FL algorithm is presented to  
165 evaluate the amount of data exchanged among nodes within the IoT-Edge-Cloud scenario.  
166 In Section 3.2, we describe our reference scenario and we outline how FL works in such  
167 a scenario; in Section 3.3 we describe the quantization approaches for NN that we will  
168 use to compress the NN models before transmission; finally, in Section 3.4 we describe  
169 our proposed FLQ and  $\Delta$ FLQ algorithms. For clarity, Table 1 summarizes all notations  
170 used herein.

#### 171 3.1. Communication Paradigm

172 IoT communication paradigms have been widely studied in the literature, and some  
173 studies have assessed the out-performance of information-centric networking, based on  
174 the publish/subscribe paradigm, w.r.t. the client/server paradigm [28]. More specifically,  
175 concerning IP-based data exchange solutions, the two most diffused IoT application proto-  
176 cols are the Constrained Application Protocol (CoAP) and the Message Queuing Telemetry  
177 Transport (MQTT) protocol, the latter being natively publish/subscribe. In MQTT,  
178 data producers (publishers) and data consumers (subscribers) are decoupled through a  
179 rendezvous node called *broker*. Data streams are organized into logical flows called *top-*  
180 *ics*. Each data packet is sent to the broker that maintains the list of active subscriptions  
181 and topics. Differently, CoAP adheres to the Representational State Transfer (REST)  
182 architectural style, providing support for resource-constrained environments. Resources  
183 are encapsulated by CoAP servers (data producers) and addressable by uniform resource  
184 identifiers. A CoAP client (data consumer) sends its request to retrieve a resource located  
185 on an IoT CoAP server. However, a publish/subscribe-like paradigm can be implemented  
186 also in CoAP, by exploiting the observer pattern in RFC 7641 and the proxy functionality  
187 in IETF RFC 7252. Therefore, both MQTT and CoAP can be deployed in the proposed  
188 framework to implement such a communication paradigm. In [29], the authors showed  
189 how, by implementing the publish/subscribe-like exchange advantage, CoAP can outper-  
190 form MQTT in terms of throughput, efficiency, and error resiliency. At present, however,  
191 MQTT is slightly more supported by a larger set of IoT embedded operating systems

Table 1: Table of symbols

Symbol	Definition
$n$	Number of edge servers
$C_1, \dots, C_n$	Edge servers
$S$	Cloud node
$A$	Generic dataset
$D, D_i$	A dataset
$M, M_i$	A NN model
$W, W_i, W'_i, \Delta W$	Weights matrices in $\mathbb{R}^{n \times m}$
$\hat{W}, \hat{W}_i, \hat{\Delta W}$	Quantized matrices in $\mathbb{R}^{n \times m}$
$\tau$	Number of local learning epochs
$T$	Number of federated learning rounds
$k$	number of quantization bits
$w_{\min}, w_{\max}$	Minimum and maximum values in $W$
$\mathbf{A}^{(i)}$	A matrix in $\mathbb{R}^{n \times m}$
$\mathbf{B}^{(i)}$	A matrix in $\{-1, +1\}^{n \times m}$
$\mathbf{1}_{n \times m}$	Element-wise product identity in $\mathbb{R}^{n \times m}$
$\epsilon_k$	$k$ -bit quantization error
$\alpha, \alpha_i, \alpha^*, \alpha_i^*$	Real values
$\hat{\mathbf{W}}^{(i)}$	Residual matrix

192 and micro-controllers, like Arduino, ESP6682, etc., and lightweight implementations of  
193 the MQTT broker are available for both constrained and unconstrained platforms [28].  
194 In our reference scenario, illustrated in Figure 1, a set  $\{C_1, \dots, C_n\}$  of *edge servers* are  
195 connected to a central *cloud node*  $S$ . Every edge server is connected to a certain number  
196 of *IoT devices* and defines a *privacy domain*, i.e., a domain wherein data generated by the  
197 relative devices can move without privacy and/or security threats. Data contained in a  
198 privacy domain cannot move outside of that. Within a privacy domain, an MQTT broker  
199 is responsible to dispatch data between producers, i.e. the IoT nodes, and consumers.  
200 The MQTT broker service is naturally located on the edge server. Thus, an MQTT topic  
201 is applied to messages to identify the data flows from producers to consumers. Yet, a  
202 database service is deployed within the security domain and subscribes to the relative  
203 MQTT topic. Therefore, each edge server  $C_i$  manages, privately, a local dataset  $D_i$ , fed  
204 by the devices publishing on the relative  $i^{th}$  topic. For example, a local dataset can be  
205 composed of video and/or audio samples collected at home. Finally, the parameters of a  
206 NN are stored on each edge server, which can access the dataset  $D_i$ : the training mecha-  
207 nism derives a model  $M_i$  from the  $D_i$  dataset, which is, then, delivered to the Cloud for  
208 the FL task.

### 209 3.2. Federated Learning Scheme

210 During FL, the cloud node  $S$  aims at training a NN model encoded as a set of param-  
211 eters  $W$ . The number of parameters depends on the structure of the NN being trained.

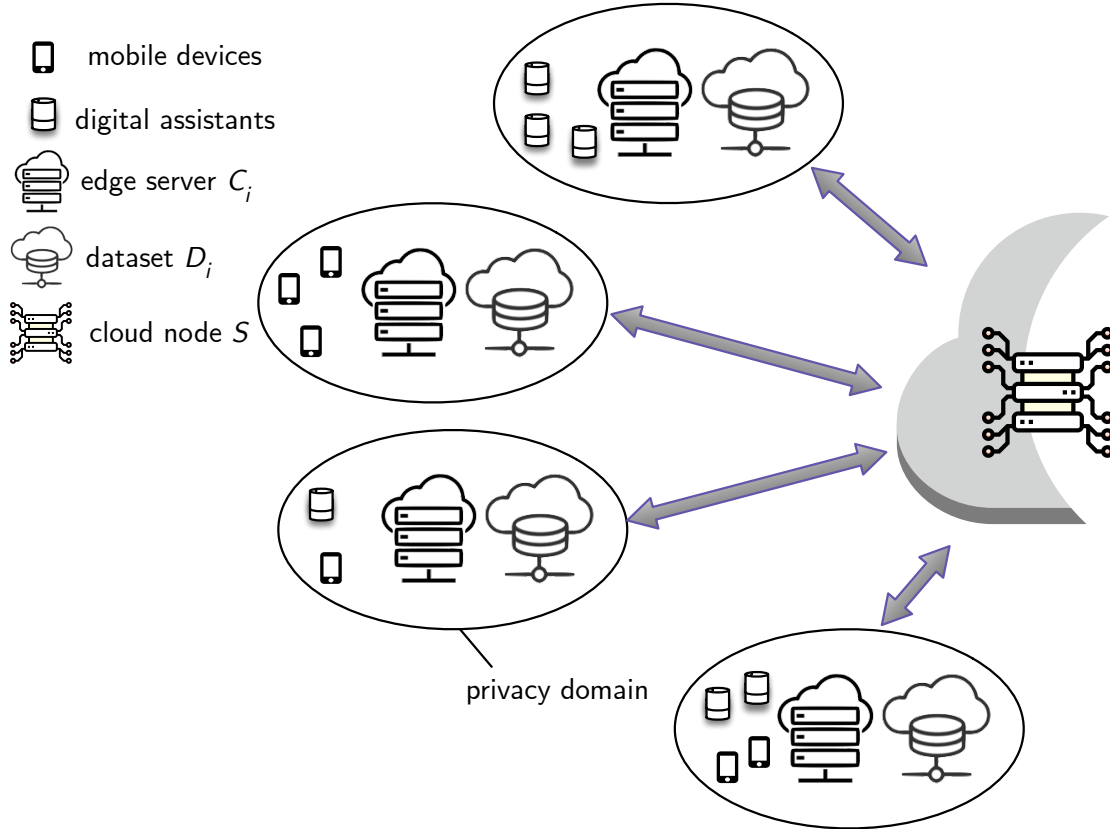


Figure 1: Federated learning reference scenario.

212 Algorithm 1 and Figure 2 illustrate the FL algorithm. The procedure  $\text{TRAIN}(W, D)$   
 213 performs the training of the NN model with parameters  $W$  exploiting a generic dataset  
 214  $D$  over a single training epoch, and returns the trained model. Namely, this procedure  
 215 computes, in each training round, the loss function produced by the NN model with cur-  
 216 rent parameters  $W$  when fed with the training data  $D$ . Then, the loss is back-propagated  
 217 through the network, and the updated network parameters  $W$  are returned [30]. After  
 218 the random initialization of the model  $W$  on the cloud node  $S$  (line 1 in Algorithm 1,  
 219 INIT box in Figure 2), the federated learning procedure advances in (federated learning)  
 220 rounds (line 2 in Algorithm 1). At the beginning of a round, the cloud node  $S$  distributes  
 221 the current global model  $W$  to all the edge servers  $C_1, \dots, C_n$  in a multicast manner  
 222 (line 4 in Algorithm 1, double-lined arrows in Figure 2). Then, each edge server  $C_i$  per-  
 223 forms an independent training of the global model  $W$  using its own local dataset  $D_i$  over  
 224  $\tau$  (local learning) epochs, generating a locally trained model  $W_i$  (lines 5-6 in Algorithm 1,  
 225 TRAIN box in Figure 2). The edge servers send back their locally trained models to the  
 226 cloud node (line 7 in Algorithm 1, double-lined arrows in Figure 2) and, finally, the cloud  
 227 node merges the  $n$  locally trained models into a new global model (FEDAVG box in Fig-  
 228 ure 2), to be distributed again during the next round. The merging of the local models  
 229 is an element-wise weighted average of the matrices  $W_1, \dots, W_n$  (line 8 in Algorithm 1).  
 230 This procedure is repeated until a maximum number  $T$  of rounds is reached, and the final  
 231 global model  $W$  is returned.

232 The FL algorithm just described allows each edge server to keep its training data,  
 233 collected through the respective server nodes, in its local privacy domain. In doing so,

---

**Algorithm 1:** The FL algorithm.

---

**Input** :  $n$  local datasets  $D_1, \dots, D_n$  at edge servers  $C_1, \dots, C_n$   
a number of rounds  $T$   
a number of epochs  $\tau$

**Output:** A matrix  $W$  of NN weights

FL( $D, D_1, \dots, D_n, T, \tau$ ):

```

1 Matrix  $W$  is randomly initialized
2 for  $T$  rounds do
3   for  $i \leftarrow 1$  to  $n$  do
4      $S$  sends  $W$  to  $C_i$  as  $W_i$ 
5     for  $\tau$  epochs do
6        $W_i \leftarrow \text{TRAIN}(W_i, D_i)$ 
7      $C_i$  sends  $W_i$  to  $S$ 
8      $W \leftarrow \sum_{i=1}^N \frac{|D_i|}{|D|} W_i$ 
9 return  $W$ 

```

---

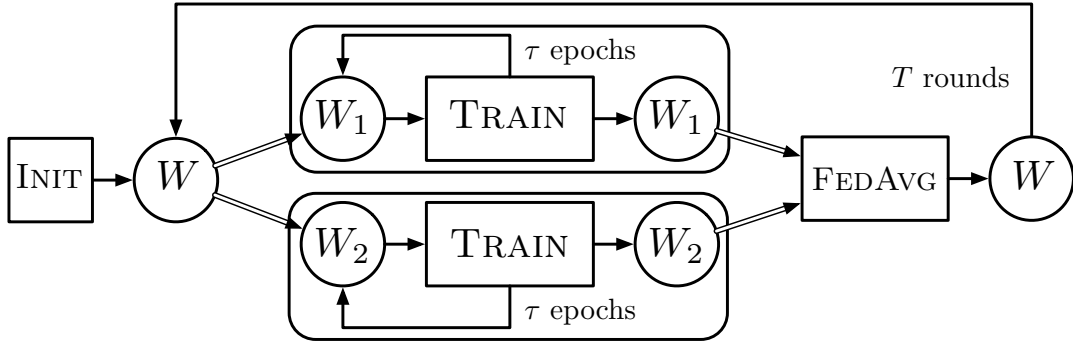


Figure 2: Main phases of the FL algorithm.

234 FL aims at producing a final NN model whose accuracy is as much as possible close to  
235 the accuracy of a NN model generated by a global training procedure performed on the  
236 aggregated dataset  $D_1 \cup \dots \cup D_n$ .

### 237 3.3. Model Quantization Schemes

238 Model quantization aims at computing a representation  $\hat{W}$  of the NN weights  $W$  with  
239 a smaller memory footprint. In particular,  $k$ -bit binary quantization maps the weights of  
240 a NN model to  $\{-1, +1\}^k$ .

241 For example, the binary quantization (BINQ), introduced in [31], assumes  $k = 1$ , and  
242 every weight is trivially quantized according to its sign. This quantization can be easily  
243 extended to  $k$  bits as shown in [32].

244 Besides this simple quantization heuristic, the NN quantization schemes can be grouped  
245 into *random* and *error minimization* quantizations.

246 **Random quantization.** Random quantization schemes select the bits representing  
247 every weight according to some probability. These approaches give theoretical guarantees  
248 on the expected value of the error introduced by the quantization scheme utilized.

249 The *probabilistic quantizer* (PROBQ) is a 1-bit random quantization scheme [33].  
 250 Let  $w_{\max}$  and  $w_{\min}$  be the maximum and minimum values among the elements of  $\mathbf{W}$ ,  
 251 respectively. The elements  $\hat{w}_{ij}$  of the quantized matrix  $\hat{\mathbf{W}}$  are quantized as follows:

$$\hat{w}_{ij} = \begin{cases} +w_{\max} & \text{with probability } p(w_{ij}), \\ -w_{\min} & \text{otherwise.} \end{cases} \quad (1)$$

252 where  $p(x)$  is the function:

$$p(x) = \frac{x - w_{\min}}{w_{\max} - w_{\min}}. \quad (2)$$

253 The *low-precision quantizer* (LOWQ) is a multi-level random quantization scheme [34].  
 254 Let  $s$  denote the number of quantization levels, i.e.,  $k = \lceil \log_2(s) \rceil$ , and let  $\omega_{ij} = s \frac{w_{ij}}{\|\mathbf{W}\|_F}$ ,  
 255 where  $\|\cdot\|_F$  is the Frobenius norm. The elements  $\hat{w}_{ij}$  of the quantized matrix  $\hat{\mathbf{W}}$  are

$$\hat{w}_{ij} = \text{sgn}(w_{ij})(\lfloor \omega_{ij} \rfloor + \delta_{ij})\|\mathbf{W}\|_F, \quad (3)$$

256 where the random variable  $\delta_{ij}$  has the following distribution:

$$\delta_{ij} = \begin{cases} 1 & \text{with probability } \omega_{ij} - \lfloor \omega_{ij} \rfloor, \\ 0 & \text{otherwise.} \end{cases} \quad (4)$$

257 The  $\text{sgn}(x)$  function is the most suited function to quantize positive and negative values  
 258 to a single bit.

259 **Error minimization quantization.** Random quantization schemes do not make any  
 260 assumption on the actual magnitude of the error introduced by the quantization schemes.  
 261 Now, given a weight matrix  $\mathbf{W} \in \mathbb{R}^{n \times m}$ , we look for an approximation matrix  $\hat{\mathbf{W}} \in \mathbb{R}^{n \times m}$   
 262 such that:

$$\mathbf{W} \simeq \hat{\mathbf{W}} = \sum_{i=1}^k \mathbf{A}^{(i)} \odot \mathbf{B}^{(i)}, \quad (5)$$

263 where  $\mathbf{A}^{(i)} \in \mathbb{R}^{n \times m}$  is a matrix with some special structure,  $\mathbf{B}^{(i)} \in \{-1, +1\}^{n \times m}$  is  
 264 a binary matrix,  $k$  is the number of bits used to represent every entry of the original  
 265 matrix  $\mathbf{W}$ , and  $\odot$  is the element-wise multiplication<sup>1</sup> (Hadamard product). We denote  
 266 the element-wise product identity as  $\mathbf{1}_{n \times m}$ .

267 To investigate the optimality of the approximation matrix  $\hat{\mathbf{W}}$ , let's define the  $k$ -bit  
 268 *quantization error*  $\varepsilon_k$  between the matrix  $\mathbf{W}$  and its quantized version  $\hat{\mathbf{W}}$  as:

$$\varepsilon_k(\mathbf{W}, \hat{\mathbf{W}}) = \|\mathbf{W} - \hat{\mathbf{W}}\|_F = \left\| \mathbf{W} - \sum_{i=1}^k \mathbf{A}^{(i)} \odot \mathbf{B}^{(i)} \right\|_F. \quad (6)$$

269 The optimal 1-bit quantization is defined as the solution  $\mathbf{A}^*$ ,  $\mathbf{B}^*$  of the following  
 270 optimization problem:

$$\begin{aligned} J(\mathbf{A}, \mathbf{B}) &= \varepsilon_1(\mathbf{W}, \hat{\mathbf{W}}) = \|\mathbf{W} - \mathbf{A} \odot \mathbf{B}\|_F, \\ \mathbf{A}^*, \mathbf{B}^* &= \arg \min_{\mathbf{A}, \mathbf{B}} J(\mathbf{A}, \mathbf{B}). \end{aligned} \quad (7)$$

---

<sup>1</sup>It produces a matrix with the same dimensions as the operands, where each output element  $i, j$  is the product of corresponding elements  $i, j$  of the two input matrices.

271 As reported in [32], the optimal solution to problem (7) for  $\mathbf{A} = \alpha \mathbf{1}_{n \times m}$ , where  $\alpha \in \mathbb{R}$ ,  
 272 is

$$\alpha^* = \frac{1}{nm} \sum_{i=1}^n \sum_{j=1}^m |w_{ij}| \quad \mathbf{B}^* = \text{sgn}(\mathbf{W}). \quad (8)$$

273 In general, directly minimizing the  $k$ -bit quantization error described in problem (6)  
 274 with  $k > 1$  is NP-hard [35]. All the proposed schemes minimizing the quantization error  
 275 address such a problem with algorithms based on some kind of heuristics. All these  
 276 schemes assume  $\mathbf{A}^{(i)} = \alpha_i \mathbf{1}_{n \times m}$ , where  $\alpha_i \in \mathbb{R}$ .

277 The *residual quantization* (RESQ) [36] is a  $k$ -bit error minimization quantization  
 278 leveraging the residual matrices  $\hat{\mathbf{W}}^{(i)}$ , defined as

$$\hat{\mathbf{W}}^{(i)} = \mathbf{W} - \sum_{j=1}^i \mathbf{A}^{(j)} \odot \mathbf{B}^{(j)} \quad \text{for } i = 1, \dots, k, \quad (9)$$

279 with  $\hat{\mathbf{W}}^{(0)} = \mathbf{W}$ .

280 Sequentially, for each  $i = 1, \dots, k$ , RESQ minimizes the residual errors  $\|\hat{\mathbf{W}}^{(i)} - \mathbf{W}\|_F$   
 281 one at a time. The optimal solutions to these  $k$  minimization problems are similar to the  
 282 optimal solution of the 1-bit quantization scheme, i.e.:

$$\alpha_i^* = \frac{1}{nm} \sum_{j=1}^n \sum_{k=1}^m |\hat{w}_{jk}^{(i-1)}|, \quad (10)$$

$$\mathbf{B}^{(i)*} = \text{sgn}(\hat{\mathbf{W}}^{(i-1)}).$$

283 In most cases, solutions to these  $k$  minimization problems will not be optimal for the  
 284 original minimization problem in (6). Hence, the  $\alpha_i$  values can be recomputed at every  
 285 step, once the first  $j$  optimal  $\mathbf{B}^{(i)*}$  matrices have been computed. This is carried out by  
 286 solving the following minimization problem:

$$J(\alpha_1, \dots, \alpha_j) = \left\| \mathbf{W} - \sum_{i=1}^j \alpha_i \mathbf{B}^{(i)*} \right\|_F, \quad (11)$$

$$\alpha_1^*, \dots, \alpha_j^* = \arg \min_{\alpha_1, \dots, \alpha_j} J(\alpha_1, \dots, \alpha_j).$$

287 Let us define the *vectorization* operator  $\text{vec}(\mathbf{X})$ , which returns a column vector, whose  
 288 elements are the stacking of the columns of the matrix  $\mathbf{X}$  on top of one another, and  
 289 let  $\mathbf{B}_j$  be the matrix, whose columns are  $\text{vec}(\mathbf{B}^{(1)*}), \dots, \text{vec}(\mathbf{B}^{(j)*})$ . The least squares  
 290 solution of problem (11) at step  $j$  is given by:

$$[\alpha_1, \dots, \alpha_j]^\top = (\mathbf{B}_j^\top \mathbf{B}_j)^{-1} \mathbf{B}_j^\top \text{vec}(\mathbf{W}). \quad (12)$$

291 During the re-computation of the  $\alpha_j$  values in Eq. (12), the computed  $\mathbf{B}^{(i)*}$  matrices  
 292 are no longer optimal for problem (7). Starting with the solution given by Eq. (10), the  
 293 *iterative quantization* (ITERQ) [23] iteratively re-computes these matrixes as follows:

294 1. compute the  $\alpha_1, \dots, \alpha_k$  values with Eq. 12 with all  $\mathbf{B}^{(1)*}, \dots, \mathbf{B}^{(k)*}$  matrices  
 295 known;

---

**Algorithm 2:** The proposed FLQ algorithm.

---

**Input** :  $n$  local datasets  $D_1, \dots, D_n$  at edge servers  $C_1, \dots, C_n$   
a number of rounds  $T$   
a number of epochs  $\tau$

**Output:** A matrix  $W$  of NN weights

FLQ( $D, D_1, \dots, D_n, T, \tau$ ):

```

1  Matrix  $W$  is randomly initialized
2  for  $T$  rounds do
3       $W \leftarrow \text{QUANTIZE}(W)$ 
4      for  $i \leftarrow 1$  to  $n$  do
5           $S$  sends  $W$  to  $C_i$  as  $W_i$ 
6          for  $\tau$  epochs do
7               $W_i \leftarrow \text{TRAIN}(W_i, D_i)$ 
8               $W_i \leftarrow \text{QUANTIZE}(W_i)$ 
9               $C_i$  sends  $W_i$  to  $S$ 
10          $W \leftarrow \sum_{i=1}^n \frac{|D_i|}{|D|} W_i$ 
11 return  $W$ 

```

---

- 296 2. build all possible  $2^k$  combinations of the  $\alpha_1, \dots, \alpha_k$  values with  $-1, +1$  and store  
297 them in a binary search tree data structure<sup>2</sup>;  
298 3. for each element of  $\mathbf{W}$ , select the closest combination and assign the corresponding  
299 values to the  $\mathbf{B}^{(1)*}, \dots, \mathbf{B}^{(k)*}$  matrices accordingly.

### 300 3.4. Federated Learning with Quantization

301 We propose to improve the efficiency of FL by reducing the amount of data being  
302 transferred from the edge servers to the cloud node. To do so, we propose to apply  
303 quantization schemes to the models being transferred during the execution of Algorithm 1.  
304 In particular, we propose to perform a first quantization after the model is updated by  
305 the edge servers, and a second quantization, before the global model is transferred from  
306 the cloud node to the edge servers.

307 Algorithm 2 and Figure 3 illustrate our proposed federated learning with quantization  
308 (FLQ) algorithm. FLQ aims at obtaining a NN model  $W$  while i) reducing the volume of  
309 data (in bytes) transferred among the cloud node and the edge servers while exchanging  
310 the model updates, and ii) preserving the privacy of the data stored on the edge servers.

311 In the FLQ algorithm, an FL round is now made up of 6 steps: (1) the cloud node  
312 applies quantization to its global model (line 3 in Algorithm 2, Q box on the left in  
313 Figure 3), (2) the cloud node sends its quantized global model (the red  $W$  on the left in  
314 Figure 3) to every edge server (line 5 in Algorithm 2, double-lined arrows in Figure 3),  
315 (3) each edge server trains its own local model, starting from the received quantized  
316 global model, on its data (line 7 in Algorithm 2, TRAIN box in Figure 3), (4) each edge  
317 server applies quantization to its locally trained model (line 8 in Algorithm 2, Q box  
318 on the right in Figure 3), (5) each edge server sends its quantized locally trained model

---

<sup>2</sup>For  $k = 2$ , assuming  $\alpha_1 > \alpha_2 > 0$ , the combinations are  $-\alpha_1 - \alpha_2$ ,  $-\alpha_1 + \alpha_2$ ,  $\alpha_1 - \alpha_2$  and  $\alpha_1 + \alpha_2$ .

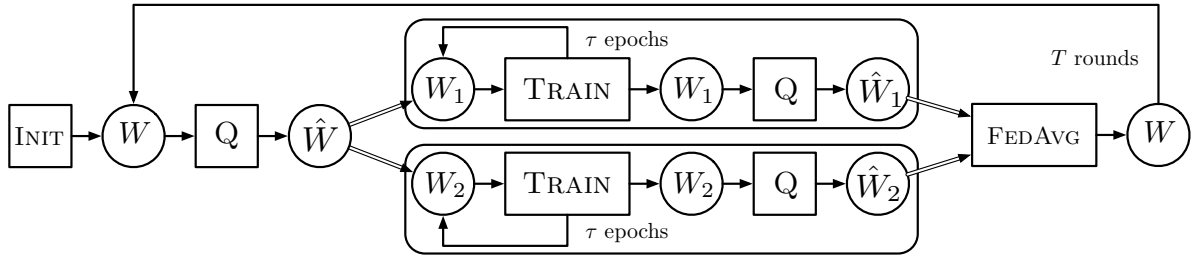


Figure 3: Main phases of FLQ algorithm.

319 ( $W$  on the right in Figure 3) back to the cloud node (line 9), and (6) the cloud node  
 320 applies weighted averaging to the quantized local models received (line 10 in Algorithm 2,  
 321 FEDAVG box in Figure 3). The procedure `QUANTIZE( $W$ )` performs the quantization of  
 322 the parameters  $W$  exploited by NN model. In principle, it may implement any of the  
 323 quantization approaches described in Section 3.3.

324 Since quantization can introduce large errors in the model weights in later rounds due  
 325 to the smaller and smaller changes applied to the weights during training, we also propose  
 326  $\Delta$ FLQ, a variant of the FLQ algorithm, illustrated in Algorithm 3. The  $\Delta$ FLQ algo-  
 327 rithm applies the FLQ algorithm not to the full models, but on their difference before  
 328 and after local training, at the edge servers, and before and after federated averaging, at  
 329 the cloud node. Initially, a random matrix is initialized on the cloud node and distributed  
 330 to all edge servers<sup>3</sup> (lines 1-3). Then, at every round, the  $\Delta W$  matrix, initially set to 0, is  
 331 quantized (line 6) and transferred to the edge servers (line 8). Each edge server sums this  
 332 matrix to its local model (line 9), and then it performs the local training on a *copy*  $W'_i$  of  
 333 their local models (lines 10-12). At the end, each edge server computes the  $\Delta W_i$  matrix,  
 334 i.e., the difference between the local model after and before the training, quantizes it and  
 335 transfers it back to the cloud node (lines 13-14). Finally, the cloud node computes the  
 336 weighted average of the received differences and stores it in the  $\Delta W$  matrix (line 15).

337 In the next section, we experimentally measure the accuracy of the final learned models  
 338 produced by FLQ and  $\Delta$ FLQ w.r.t. other training algorithms, such as local learning  
 339 and federated learning, and we evaluate the benefits of our FLQ and  $\Delta$ FLQ algorithms  
 340 at reducing the data transmitted among the cloud node and edge servers when different  
 341 quantization approaches are exploited.

#### 342 4. Experimental Setup

343 We experiment FLQ and  $\Delta$ FLQ on the *next word prediction* task that, given a  
 344 sentence context, i.e., a sequence of words, aims at learning a language model, i.e., a  
 345 probability distribution over the words conditioned on a given context, to predict the  
 346 *next* most likely word that appears after the context [23]. As such, this task can be seen  
 347 as a time series prediction problem where data points are words in a given vocabulary,  
 348 and given the first  $n$  words, we aim at predicting the  $n + 1$ -th word. It is a common use  
 349 case used in mobile applications like, for example, predictive keyboards or personal voice  
 350 assistants. In our scenario, we assume that every edge server has a private collection

<sup>3</sup>In practice the matrices can be initialized by communicating the random seed used to generate the random weights.

---

**Algorithm 3:** The proposed  $\Delta$ FLQ algorithm.

---

**Input** :  $n$  local datasets  $D_1, \dots, D_n$  at edge servers  $C_1, \dots, C_n$ ,  
a number of rounds  $T$   
a number of epochs  $\tau$

**Output:** A matrix  $W$  of NN weights

$\Delta$ FLQ( $D, D_1, \dots, D_n, T, \tau$ ):

```

1  Matrix  $W$  is randomly initialized
2  for  $i \leftarrow 1$  to  $n$  do
3     $S$  sends  $W$  to  $C_i$  as  $W_i$ 
4   $\Delta W \leftarrow 0$ 
5  for  $T$  rounds do
6     $\Delta W \leftarrow \text{QUANTIZE}(\Delta W)$ 
7    for  $i \leftarrow 1$  to  $n$  do
8       $S$  sends  $\Delta W$  to  $C_i$  as  $\Delta W_i$ 
9       $W_i \leftarrow W_i + \Delta W_i$ 
10      $W'_i \leftarrow W_i$ 
11     for  $\tau$  epochs do
12        $W'_i \leftarrow \text{TRAIN}(W'_i, D_i)$ 
13        $\Delta W_i \leftarrow \text{QUANTIZE}(W'_i - W_i)$ 
14        $C_i$  sends  $\Delta W_i$  to  $S$ 
15      $\Delta W \leftarrow \sum_{i=1}^n \frac{|D_i|}{|D|} \Delta W_i$ 
16      $W \leftarrow W + \Delta W$ 
17 return  $W$ 

```

---

351 of text written by the user, and we aim at learning a global language model without  
352 disclosing the private collections.

353 **Dataset.** We conduct the experimental evaluation using the WikiText-2 dataset [37].  
354 The WikiText-2 dataset is composed of 720 text articles: 600 articles in the training set  
355 and 60 in both the validation and the test sets, respectively. The training set consists of  
356 2, 088, 628 tokens, while the validation and test sets consist of 217, 646 and 245, 569 tokens,  
357 respectively. The vocabulary consists of 33, 278 words. In the following analysis, we use  
358 the training set to train our models while the validation and test sets are used for early  
359 stopping the training of the model, and to measure its final performance, respectively.

360 **Neural Architecture.** We perform next word prediction [12] by training a Long Short  
361 Term Memory Network (LSTM). In our experiments, following the setting and hyperpa-  
362 rameters laid out in [23], we train an LSTM model with an input embedding layer with  
363 200 neurons, two hidden LSTM layers, each one composed of 512 neurons, and an output  
364 linear layer with 33, 278 outputs, one per word in the vocabulary. We use a batch size of  
365 100 words. The training of the network is regularized with learning rate decay, i.e., if the  
366 validation loss does not decrease in an epoch, the learning rate, initially set to 20, is de-  
367 creased by a factor of 4, up to a minimum value of  $10^{-4}$ . For regularization purposes, we  
368 also employ a dropout strategy by setting the dropout rate to 0.5. Moreover, to address  
369 the gradient explosion problem when training the LSTM, we set the gradient clipping

370 to 0.25, the gradient norm clipping to 0.3, and the weight clipping to 1.0. The training  
371 of the LSTM is performed by minimizing the cross-entropy loss until the learning rate  
372 decreases below the minimum value of  $10^{-4}$  (early-stopping condition) or for a maximum  
373 of 80 epochs.

374 **Implementation Details.** The experimental framework is implemented in PyTorch  
375 1.4.0<sup>4</sup>. Experiments are performed using a Tesla T4 GPU on an AMD EPYC CPU  
376 clocked at 2.2 GHz and 24 GB of RAM. The machine works as both cloud node and edge  
377 server and we experiment with 2 edge servers. In the evaluation of federated scenarios,  
378 edge servers are emulated by training the LSTM models locally on the machine on equally-  
379 sized and disjoint partitions of the training set. All the quantization strategies tested in  
380 this paperwork quantize independently the weight matrices of the LSTM. No quantization  
381 is performed on the input embedding layer.

## 382 5. Experimental Evaluation

383 We now present a comprehensive analysis of the performance of the FLQ and  $\Delta$ FLQ  
384 algorithms by investigating two main research questions (RQs):

385 **RQ1.** What is the impact of the FLQ and  $\Delta$ FLQ algorithms on the accuracy of the  
386 learned models, measured in terms of validation and test losses?

387 **RQ2.** What is the reduction in terms of data transmitted between edge servers and cloud  
388 node by the FLQ and  $\Delta$ FLQ algorithms w.r.t. standard FL approaches, i.e., without  
389 quantization?

390 To investigate our RQs, we designed a comprehensive experimental setting structured  
391 in the following four scenarios.

392 **Local Learning (LL).** In this scenario, we perform the training of the LSTM locally on  
393 the cloud node on the full training dataset.

394 **Local Quantization (LQ).** In this scenario, we perform the training of the LSTM  
395 locally on the cloud node by locally applying the quantization schemes introduced in  
396 Section 3. As for LL, the training of the LSTM is performed on the full training dataset.

397 **Federated Learning (FL):** In this scenario, we experiment with the standard FL  
398 algorithm. We learn the LSTM on each edge server on its partition of the training data.  
399 At the end of each round, the edge servers send their LSTM model to the cloud node,  
400 which performs the federated averaging step. The model is then sent back to all the edge  
401 servers. The edge servers compute the loss on the training set, while the validation and  
402 test losses are computed on the cloud node after the federated averaging step.

403 **Federated Learning with Quantization (FLQ):** In this scenario, we experiment  
404 with the FLQ and  $\Delta$ FLQ algorithms. The edge servers perform model quantization  
405 before sending their model to the cloud node. Then, the cloud node performs a federated  
406 averaging of the received models and quantizes again the result before sending the model  
407 back to the edge servers.

408 Firstly, we assess the performance of the FL scenario w.r.t. the LL scenario. Sec-  
409 ondly, we apply quantization both in local scenario (LQ) and in the federated scenario  
410 (FLQ). In both scenarios, we first assess the performance of our FLQ algorithm, then we  
411 assess the performance of our  $\Delta$ FLQ algorithm. The rationale behind this experimental

---

<sup>4</sup>We plan to release the code upon acceptance of the paper.

412 methodology is to evaluate the impact of federated learning on the performance of the  
 413 LSTM, i.e., the effect of working on partitioned training data for each worker. Then,  
 414 we assess the impact of quantization techniques both in the LQ scenario, i.e., where no  
 415 data partitioning is introduced by federated learning, and in the FLQ scenario, i.e., with  
 416 both quantization and data partitioning. Finally, we compare the size of the models  
 417 transmitted in the federated scenarios to quantify the data volume reductions yielded by  
 418 the quantization schemes considered.

419 We perform our experiments by reporting the effectiveness of the LSTM in terms of:  
 420 i) test loss on the final model and ii) validation loss observed during the training process  
 421 to analyze the convergence speed of each method. The size of the models is computed by  
 422 counting the number of bits of both the quantized and un-quantized parameters of our  
 423 LSTM.

424 *A note on time units.* Local and federated scenarios employ different portions of the  
 425 training data, depending on the number of workers, and perform different local training  
 426 epochs and federated training rounds. The total wall-clock time required to train the  
 427 different models depends on four factors: (i) the number of workers processing their  
 428 portions of the dataset, (ii) the number of learning epochs used by each worker to locally  
 429 learn its model, (iii) the number of federated rounds used to collect the local models and  
 430 to compute their federated averaging, and (iv) the quantization process. In the following,  
 431 we report on the  $x$ -axis the training time measured in *time units*. We assume that the  
 432 time required to perform a single training epoch on a single worker on the whole dataset  
 433 corresponds to 1-time unit. On  $w$  workers, the whole dataset, divided among all workers,  
 434 requires  $1/w$  time units to be processed in a single learning epoch in parallel on  $w$  workers.  
 435 On  $w$  workers,  $\tau$  training epochs require  $\tau/w$  time units to perform the local training.  
 436 Finally, by measuring the time required by the quantization, and the model transmission  
 437 processing, we found out that their overhead is negligible w.r.t. the time required to  
 438 perform a single learning epoch. For this reason, we do not take it into account.

### 439 5.1. Local Learning and Federated Learning Scenarios

440 In this section, we propose an experimental analysis of the performance of the LSTM  
 441 network when trained both in LL and FL scenarios. Figure 4 reports the validation  
 442 loss of the LSTM during the training in the FL scenario according to Algorithm 1 by  
 443 varying  $\tau$  in  $\{1, 2, 4, 8, 16\}$ , i.e., the number of epochs performed by two edge servers  
 444 before transmitting the model to the cloud node for federated averaging.

445 We compare the performance achieved in the FL scenario with the one achieved in the  
 446 LL scenario (bold red line). When training in LL, the validation loss drops significantly  
 447 very soon, reaching the minimum value of 4.73 after 29 units of time. Moreover, the  
 448 plot shows a clear trend: the more epochs are performed locally to each worker in a  
 449 round of federated learning, i.e., the larger values of  $\tau$ , the slower is the convergence  
 450 of the validation loss to its minimum. The rationale of this behavior lies in the role of  
 451 federated averaging, i.e., (FEDAVG), which allows sharing the knowledge learned by the  
 452 two workers in isolation in a single LSTM model that is then used locally by edge servers to  
 453 continue the learning. Moreover, in the FL scenario, the neural network achieves a better  
 454 performance in terms of validation loss than the LL corresponding loss with small values  
 455 of  $\tau$ , i.e.,  $\{1, 2, 4, 8\}$ . This does not hold for  $\tau = 16$ , where the FL minimum validation  
 456 loss does not outperform the corresponding loss of the LL scenario. The rationale behind  
 457 the better performance shown in the FL scenario w.r.t. the LL scenario relies in the

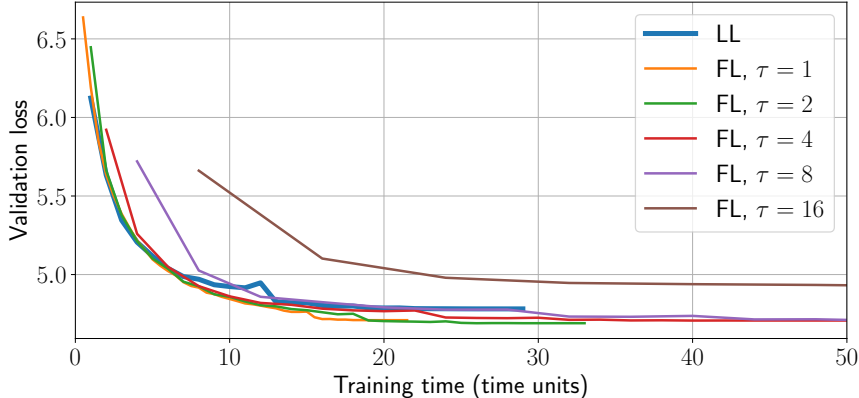


Figure 4: Validation loss of the LSTM network in the FL scenario by varying the number of epochs  $\tau$  in a round of federated learning. We also report the performance on the LL scenario (red bold line).

458 federated averaging step that introduces a regularization effect in the resulting merged  
 459 model, thus achieving a better generalization of the model learned.

Table 2: Rounds and time units to minimum validation loss (Rounds & Time), minimum validation loss (Val. loss) and test loss of the LSTM network model in the LL and FL scenarios, by varying  $\tau$ .

	$\tau$	Rounds	Time	Val. loss	Test loss
LL	–	27	27	4.78	4.73
FL	1	37	18.5	4.71	4.64
	2	25	25	4.69	4.64
	4	19	38	4.71	4.63
	8	27	108	4.70	4.66
	16	28	224	4.81	4.76

460 To better show the impact of  $\tau$  on the convergence speed, we deepen the analysis  
 461 by reporting, given  $\tau$ , the number of federated rounds needed by the LSTM to achieve  
 462 the minimum observed validation loss. Note that the time units required to perform a  
 463 federated learning round depend on  $\tau$ : a single federated round with 2 workers requires  
 464  $\tau/2$  time units. We report the results in Table 2. By increasing  $\tau$ , the results show a  
 465 clear slow down of the time units required by the NN to reach the minimum validation  
 466 error. In a LL scenario, the LSTM needs 27 epochs, i.e., 27 time units, to reach the  
 467 minimum validation loss. Instead, in a FL scenario, when  $\tau \in \{1, 2\}$ , the FL algorithm  
 468 is able to reach the minimum validation loss in just  $37 \times 1/2 = 18.5$  and  $25 \times 2/2 = 25$   
 469 time units, respectively, while it needs  $28 \times 16/2 = 224$  time units to reach the minimum  
 470 validation loss when  $\tau = 16$ . Table 2 also reports the final test loss achieved by the LSTM  
 471 in the LL and FL scenarios. The results consistently report the same trend identified  
 472 for the validation loss, i.e., the test loss achieved in the FL scenario outperforms the one  
 473 achieved in the FL one for  $\tau \in \{1, 2, 4, 8\}$  but not for  $\tau = 16$ .

474 To conclude, federated learning allows to gain effectiveness on the final LSTM model  
 475 with the big advantage of keeping data stored on edge servers in a decentralized manner.

476 *5.2. Local Quantization and Federated Learning with Quantization Scenarios*

477 We now focus our analysis on the impact of quantization when employed in the LQ  
 478 and FLQ scenarios. Table 3 reports the test loss achieved when training an LSTM model  
 479 with the 1-, 2-, and 3-bit quantization schemes discussed in Section 3 with  $\tau = 1$  in the  
 480 LQ and FLQ scenarios.

Table 3: Test loss of the LSTM network model in the LQ and FLQ scenarios for 1-, 2-, and 3-bit quantization schemes.

Quantization	LQ			FLQ		
	1 bit	2 bits	3 bits	1 bit	2 bits	3 bits
BINQ	33.19	-	-	39.32	-	-
PROBQ	13.28	-	-	9.71	-	-
LOWQ	-	6.91	6.88	-	6.99	6.95
RESQ	-	7.04	7.03	-	6.85	6.86
ITERQ	-	7.62	6.89	-	6.89	6.85

481 Comparing the test losses in the LQ scenario, the two 1-bit random quantization  
 482 schemes, i.e., BINQ and PROBQ, result in larger test losses than the other schemes,  
 483 i.e., LOWQ, RESQ, and ITERQ. This is due to the very limited precision obtained by  
 484 the former schemes in quantizing each model weight on a single bit, while the latter  
 485 schemes, using 2 bits, result in smaller quantization errors. For sake of space, we do  
 486 not include further results of BINQ and PROBQ schemes, as their performance is not  
 487 competitive with respect to LOWQ, RESQ, and ITERQ. Concerning the FLQ scenario, the  
 488 test losses of LOWQ, RESQ, and ITERQ are very close to the value in the corresponding  
 489 LQ scenario, i.e.,  $-1.2\%$ ,  $+2.8\%$ , and  $+0.3\%$  respectively.

490 Figure 5 reports the validation loss achieved by LOWQ, RESQ, and ITERQ when  
 491 training an LSTM model with 2-bit quantization and  $\tau = 1$  in the LQ and FLQ scenarios.  
 492 Similar results are obtained with higher values of  $\tau$  and with 3-bit quantization, with just  
 a slight increase in convergence speed, so we do not report them. In both scenarios, the

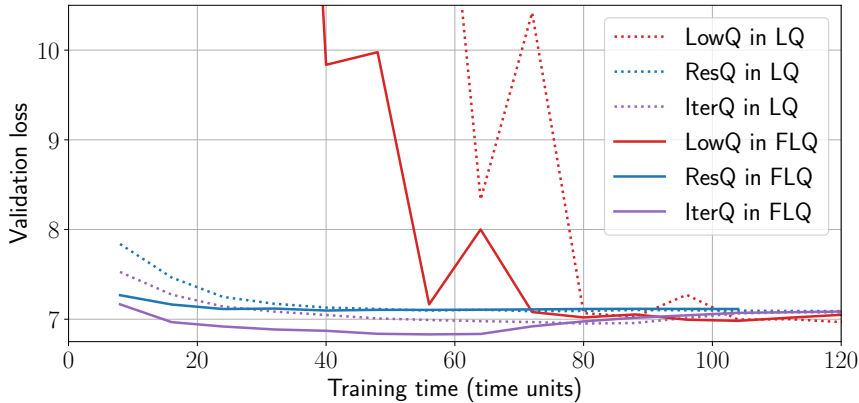


Figure 5: Validation loss of the LSTM network in the LQ and FLQ scenarios (dashed and solid lines, respectively) with 2-bit quantization schemes and  $\tau = 1$ .

494 RESQ and ITERQ quantization schemes converge to the minimum validation loss quickly,  
 495 i.e., in 4 time units or less. The LOWQ scheme starts with a higher loss value even if it  
 496 eventually converges towards the minimum value. However, the time needed by LOWQ  
 497 to converge is 33% more than the time needed by the previous two quantization schemes.  
 498 By comparing the LQ and FLQ scenarios, both the RESQ and ITERQ quantization  
 499 schemes in the FLQ scenario converge slightly slower than in the LQ scenario, while the  
 500 LOWQ quantization scheme converges faster in the FLQ scenario. We can conclude that  
 501 the considered quantization schemes can be successfully applied in federated learning  
 502 scenarios. Our experiments show that quantization in the FLQ scenario achieves the  
 503 same performance of the LQ scenario with no degradation in the effectiveness of the  
 504 LSTM model.

505 We now investigate the impact of the number of epochs  $\tau$  on the performance. In  
 506 particular, we assess if larger values of  $\tau$  allow improving the effectiveness of the LSTM  
 507 in the FLQ scenario, measured in terms of test loss. As  $\tau$  controls the number of epochs  
 508 performed on each edge server, the intuition is the following: the more epochs are per-  
 509 formed on the edge servers, the more the error introduced by quantization decreases, due  
 510 to the full precision used in the local learning process on the edge servers. In fact, no  
 511 quantization is performed during the training process on an edge server, as it is performed  
 512 only before sending the model or after receiving it from the cloud node.

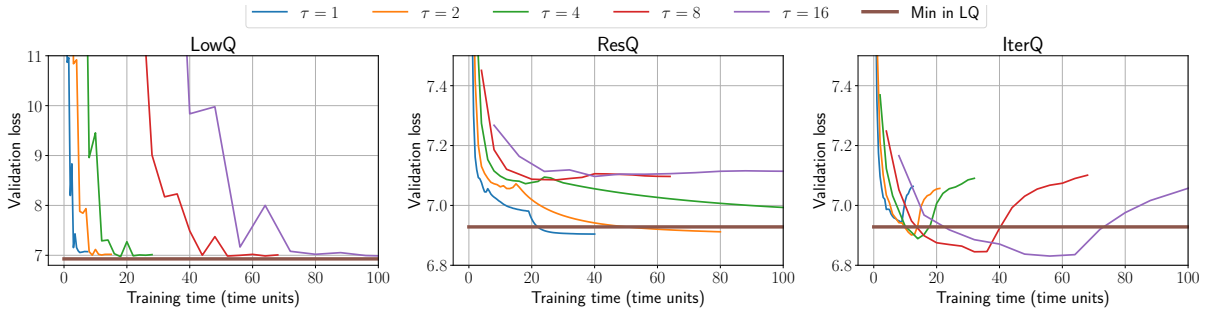


Figure 6: Validation loss of the LSTM network in an FLQ scenario using the FLQ Algorithm with 2-bit quantization schemes, by varying  $\tau$ .

Table 4: Minimum validation and test losses of the LSTM network model in the FLQ scenario using the FLQ Algorithm with 2-bit quantization schemes, by varying  $\tau$ .

$\tau$	Minimum validation loss			Test loss		
	LOWQ	RESQ	ITERQ	LOWQ	RESQ	ITERQ
1	7.06	6.90	6.95	6.99	6.85	6.89
2	7.00	6.91	6.90	6.94	6.86	6.86
4	6.97	6.97	6.89	6.92	6.93	6.85
8	6.99	7.09	6.85	6.93	7.04	6.80
16	6.98	7.10	6.83	6.93	7.05	6.77

513 Figure 6 reports the validation losses achieved using the FLQ Algorithm with 2-bit  
 514 schemes LOWQ, RESQ, and ITERQ, by varying  $\tau$  in  $\{1, 2, 4, 8, 16\}$ , and Table 4 reports  
 515 the corresponding minimum validation and test losses. As expected, the results show

516 that, for all quantization schemes, the convergence is slower as  $\tau$  increases. The results  
 517 for RESQ show that the validation and test losses are negatively affected by larger values  
 518 of  $\tau$ . Moreover, RESQ and LOWQ do not allow the LSTM to achieve the minimum  
 519 validation loss reported in the LQ scenario. In particular, only for  $\tau \in \{1, 2\}$ , RESQ is  
 520 able to reach and outperform the minimum validation loss observed in the LQ scenario,  
 521 even if, for  $\tau = 2$ , the convergence to the minimum is significantly slower than for  $\tau = 1$ .  
 522 Instead, the validation and test losses using ITERQ benefit from larger values of  $\tau$ . The  
 523 reported results show that, when increasing  $\tau$ , the effectiveness of the model, i.e., the  
 524 minimum validation loss, decreases significantly below the minimum achieved in the LQ  
 525 scenario. A side effect of the increase of  $\tau$  for ITERQ is the convergence speed, which  
 526 results to be slower. The convergence speed is important in the FLQ scenario because it  
 527 potentially affects the number of data transmissions between edge servers and the cloud  
 528 node, i.e., the slower the convergence, the larger the number of transmissions. When  
 529 increasing  $\tau$ , the slower convergence speed is reported for all the three quantization  
 530 schemes. However, when ITERQ is used, larger  $\tau$  values boost the final effectiveness.  
 531 Moreover, while the minimum validation loss for  $\tau = 1$  is reached after 8 units of time,  
 532 i.e.,  $2 \times 8 = 16$  model transmissions, with  $\tau = 16$  the minimum validation loss is reached  
 533 after 48 units of time, i.e.,  $2 \times 48/16 = 6$  model transmissions, thus more than halving  
 534 the number of transmissions required to achieve a more accurate model.

535 To conclude, we observe that a larger number of epochs performed locally on the edge  
 536 servers slows down the convergence of the whole training, even if this does not always  
 537 imply a larger number of transmissions. However, we found that large values of  $\tau$  impact  
 538 positively on the validation and test losses of the LSTM when employing the ITERQ  
 539 quantization scheme. In this case, the final model in the FLQ scenario outperforms the  
 540 final model in the LQ scenario for all  $\tau$  in  $\{2, 4, 8, 16\}$ .

### 541 5.3. Federated Learning with Delta Quantization

542 While in the FLQ scenario our FLQ algorithm produces a LSTM network with vali-  
 543 dation and test losses almost identical to the values of the LSTM network trained in the  
 544 LQ scenario, the performance of the models trained with FLQ (test loss 6.93) are higher  
 545 than the performance of the models trained with FL (test loss 4.73). We ascribe this  
 546 decrease in test loss to quantization, since it can introduce errors in the model weights.  
 547 In particular, during the training of the LSTM, the high variance in the small changes  
 548 applied to the weights at later rounds can have a negative impact on the whole training.  
 549 To mitigate this effect, in Section 3 we proposed the  $\Delta$ FLQ algorithm, a variant of the  
 550 FLQ algorithm. The  $\Delta$ FLQ algorithm quantizes the *changes* in the network weights ex-  
 551 change between the edge servers and the cloud node. We now investigate the performance  
 552 of  $\Delta$ FLQ when applied to the FLQ scenario. As before, we report the performance of  
 553 LOWQ, RESQ, and ITERQ when varying  $\tau$  in Figure 7 and Table 5.

554 The introduction of changes quantization to LOWQ negatively impacts the perfor-  
 555 mance of the final LSTM network. Indeed, the validation loss never converges to the  
 556 minimum achieved in LQ (green horizontal line). Moreover, for larger values of  $\tau$ , the  
 557 validation loss grows significantly, thus revealing the presence of heavy overfitting on the  
 558 local training data. On the other side, RESQ and ITERQ show a completely different  
 559 behaviour. For these two quantization schemes, the introduction of changes quantization  
 560 improves the performance of the FLQ algorithm. Indeed, RESQ and ITERQ outperform  
 561 the performance achieved using FLQ for all values of  $\tau$  considered. Moreover, when

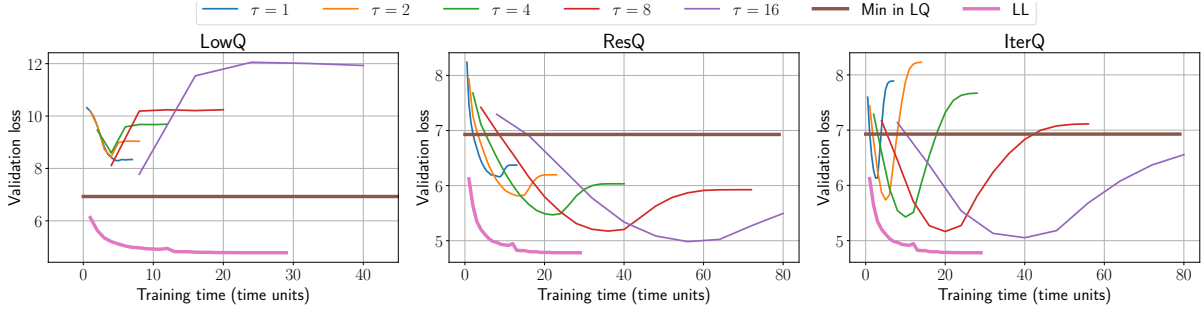


Figure 7: Validation loss of the LSTM network in the FLQ scenario using the  $\Delta$ FLQ Algorithm with 2-bit quantization schemes, by varying  $\tau$ .

Table 5: Minimum validation and test losses of the LSTM network model in the FLQ scenario using the  $\Delta$ FLQ Algorithm with 2-bit quantization schemes, by varying  $\tau$ .

$\tau$	Minimum validation loss			Test loss		
	LOWQ	RESQ	ITERQ	LOWQ	RESQ	ITERQ
1	8.29	6.16	6.13	8.27	6.10	6.06
2	8.45	5.82	5.74	8.42	5.76	5.66
4	8.59	5.47	5.43	8.54	5.41	5.36
8	8.17	5.17	5.17	8.07	5.11	5.10
16	7.78	4.98	5.05	7.73	4.91	4.99

562 increasing  $\tau$ , the validation loss gets closer to the value achieved in the LL scenario (pink  
563 line), even though convergence is slower.

564 Finally, we evaluate the impact of the number of edge servers on the best 2-bit quanti-  
565 zation scheme in the FLQ scenario using the  $\Delta$ FLQ Algorithm, i.e., ITERQ with  $\tau = 16$   
566 local learning epochs. We report the performance of ITERQ when varying the number of  
567 edge servers  $w$  in Figure 8. When more edge servers are deployed, the time units required  
568 to reach the minimum validation loss decrease with the number of servers, since more  
569 workers can carry out independently and in parallel the local learning over smaller por-  
570 tions of the whole dataset. At the same time, the smaller local dataset portions available  
571 at each edge server lead to a decrease in the performance of the whole learned model.  
572 In fact, the minimum validation loss of the final model increases as the number of edge  
573 servers increases.

574 Figure 9 illustrates the validation loss with 4 and 8 edge servers in the FLQ scenario  
575 using the FLQ and  $\Delta$ FLQ algorithms with 2-bit ITERQ quantization and  $\tau = 16$ , with  
576 a fraction  $f = 0\%$ ,  $25\%$ ,  $50\%$ , and  $75\%$  of faults. The faults happen after the second  
577 federated round, and the faulty edge servers are randomly selected at every federated  
578 round until the minimum validation loss is reached. For both 4 and 8 edge servers, the  
579  $\Delta$ FLQ algorithm reaches a smaller validation loss, confirming the better performance  
580 of the  $\Delta$ FLQ algorithm w.r.t. the FLQ algorithm reported in Tables 4 and 5 even in  
581 the presence of faults. As expected, as the number of faults increases, the minimum  
582 validation loss increases, since the dataset portions of the faulty edge servers no longer  
583 contribute to the accuracy of the global model learned via federated averaging. Moreover,  
584 the minimum validation loss increases with the number of edge servers also in the presence  
585 of faults, due to the smaller local dataset portions available to each server, confirming

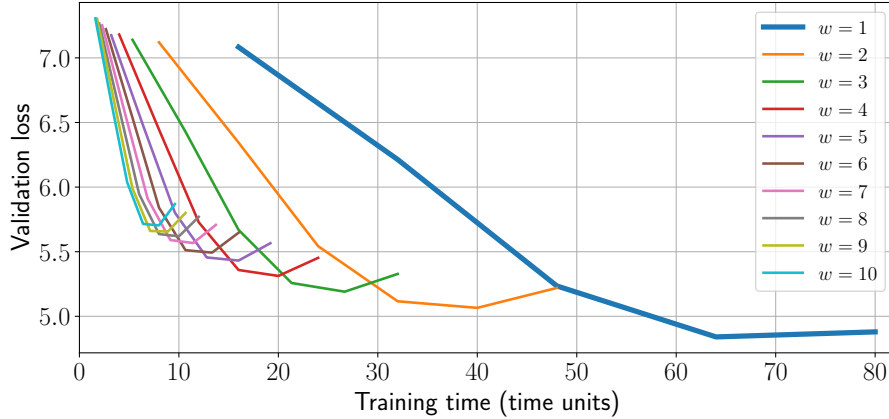


Figure 8: Validation loss of the LSTM network in the FLQ scenario using the  $\Delta$ FLQ algorithm with 2-bit ITERQ quantization and  $\tau = 16$  by varying the number of edge servers  $w$ .

586 the experimental results illustrated in Figure 8.

587 To conclude on RQ1, we experimentally showed that:

- 588 • the FLQ algorithm with LOWQ, RESQ, and ITERQ quantization schemes is able  
589 to train an LSTM network with performances similar to a LSTM network trained in  
590 the LQ scenario. In particular, the LSTM network trained with ITERQ in the FLQ  
591 scenario outperforms the LSTM network in the LQ scenario with any quantization  
592 scheme for all  $\tau$  in  $\{2, 4, 8, 16\}$ ;
- 593 • the  $\Delta$ FLQ algorithm with RESQ and ITERQ allows to train an LSTM network with  
594 performances similar to an LSTM network trained in the LL scenario. Moreover,  
595 we show that ITERQ allows for a faster convergence speed, thus allowing for a  
596 reduced number of model transmissions between edge servers and the cloud node,  
597 even in the presence of edge server faults.

#### 598 5.4. Analysis of Model Transmission Costs

599 We now evaluate the reduction in terms of data transmitted between the edge servers  
600 and the cloud node by the FLQ/ $\Delta$ FLQ algorithms w.r.t. the FL algorithm without  
601 quantization. Our LSTM model is composed of 27,249,264 parameters in total. As  
602 reported in Section 4, we do not apply quantization to the input embedding layer, whose  
603 weight matrix contains 6,655,600 elements. The other layers store 20,593,664 model  
604 weights, which are represented according to our quantization schemes.

605 By storing each parameter as a 32-bit floating point number, the LSTM network  
606 requires 109.16 MB  $\approx$  0.11 GB, while the 1-bit random quantization schemes BINQ and  
607 PROBQ store the quantized LSTM network in 29.36 MB, with a space reduction of  $3.71\times$ ,  
608 but their performance in term of test loss are poor, as we reported in Table 3. The 2-  
609 bit and 3-bit quantization schemes, i.e., LOWQ, RESQ, and ITERQ, store the quantized  
610 LSTM network in 31.94 MB (with a space reduction of  $3.42\times$ ) and 34.51 MB (with a  
611 space reduction of  $3.16\times$ ), respectively.

612 Table 6 reports the number of federated rounds required to reach the minimum val-  
613 idation loss with the FLQ and  $\Delta$ FLQ algorithms, for the 2-bit LOWQ, RESQ, and  
614 ITERQ quantization schemes. This number represents the optimal number of rounds to  
615 obtain the best LSTM network performance during training according to the validation

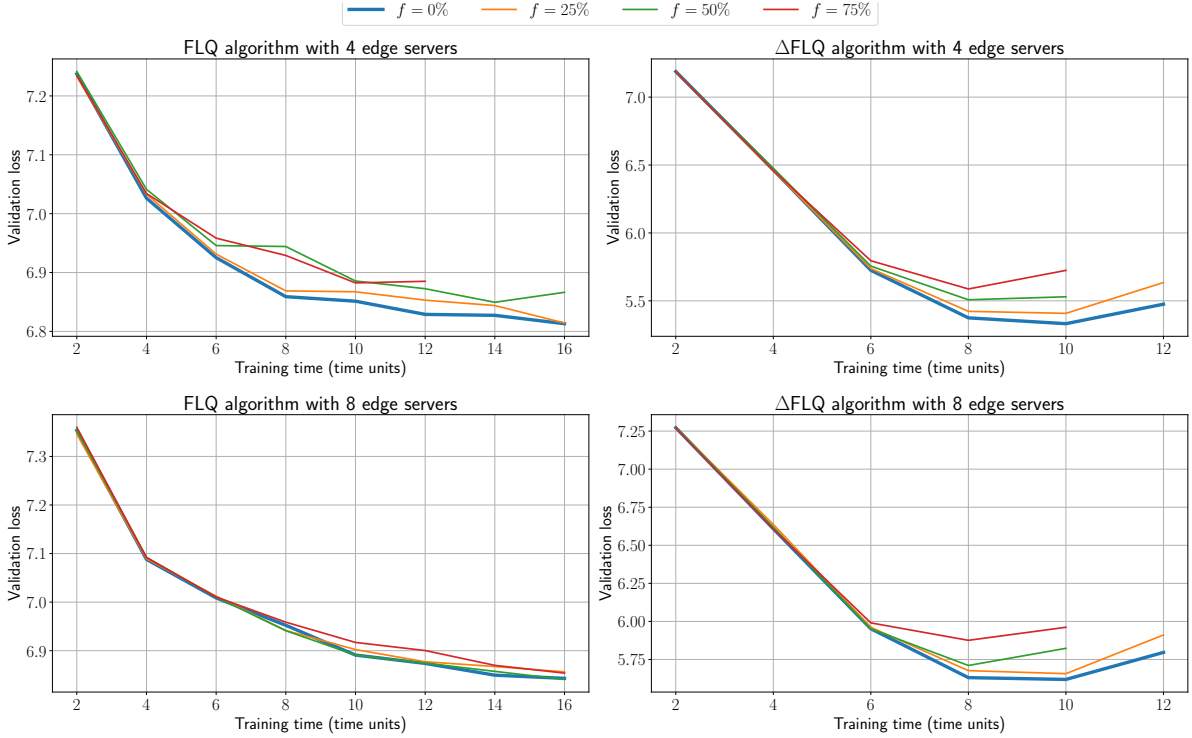


Figure 9: Validation loss of the LSTM network in the FLQ scenario using the FLQ (left) and  $\Delta$ FLQ (right) algorithms with 2-bit ITERQ quantization and  $\tau = 16$ , with 4 (top) and 8 (bottom) edge servers, with a fraction  $f$  of faults after the second federated round.

Table 6: Rounds to reach the minimum validation loss with the FLQ and  $\Delta$ FLQ algorithms, and total data exchanged between the cloud node and a single edge server (in GB) for the 2-bit quantization schemes.

$\tau$	Rounds to minimum loss			Total data exchanged		
	LOWQ	RESQ	ITERQ	LOWQ	RESQ	ITERQ
FLQ						
1	10	78	17	0.64	4.98	1.09
2	9	80	12	0.57	5.11	0.77
4	9	80	7	0.57	5.11	0.48
8	13	9	8	0.83	0.57	0.51
16	13	7	7	0.83	0.48	0.48
$\Delta$ FLQ						
1	9	17	5	0.57	1.09	0.32
2	4	14	5	0.26	0.89	0.32
4	2	11	5	0.13	0.70	0.32
8	1	9	5	0.06	0.57	0.32
16	1	7	5	0.06	0.48	0.32

616 loss, i.e., before the validation loss starts increasing and the local models start overfit-  
617 ting on the local data. The number of rounds  $T$  determines how many times the LSTM  
618 model is transferred from an edge server to the cloud node and vice-versa. Regarding

619 the FLQ algorithm, LOWQ is the best quantization scheme in terms of data exchanged  
 620 when  $\tau \in \{1, 2\}$ . For larger  $\tau$  values, the best scheme is ITERQ, requiring almost half of  
 621 the rounds required by LOWQ. As we have shown in Figure 5 and Table 4, in such cases  
 622 ITERQ produces a LSTM network with better performance than the LSTM produced by  
 623 LOWQ. With  $\tau = 16$ , FLQ with ITERQ quantization attains a validation loss of 6.83  
 624 in 7 rounds (56 time units) with a total of  $0.48 \text{ GB} \times 2 = 0.96 \text{ GB}$  transferred (taking  
 625 into account both workers), while FL attains a validation loss of 4.81 in 28 rounds (28  
 626 time units) with a total of  $0.11 \text{ GB} \times 2 \times 2 \times 28 = 12.23 \text{ GB}$  transferred, by taking into  
 627 account the full model size, the number of workers, the two-way transmissions and the  
 628 number of rounds.

629 Regarding  $\Delta$ FLQ, LOWQ still requires very few rounds to converge, but, as reported  
 630 in Figure 7 and Table 5, the performance of the resulting LSTM network is far from being  
 631 competitive. However, ITERQ is always producing a LSTM network with performance  
 632 close to the FL algorithm without quantization in very few rounds, namely 5, with just  
 633 0.32 GB of data transferred. With  $\tau = 16$ ,  $\Delta$ FLQ with ITERQ quantization gets a  
 634 validation loss of 5.05 in 5 rounds (20 time units) with a total of  $0.32 \text{ GB} \times 2 = 0.64 \text{ GB}$   
 635 transferred (taking into account both workers).

636 To conclude on RQ2, we experimentally showed that both the FLQ and  $\Delta$ FLQ  
 637 algorithms are able to largely reduce the data transferred between the cloud node and  
 638 the edge nodes in a federated learning scenario. In particular, when used with the ITERQ  
 639 quantization scheme and using 16 epochs for local training at each worker, the FLQ  
 640 algorithm trains the LSTM model with a 14% degradation in the validation loss while  
 641 reducing by a factor of  $13\times$  the total data transmitted over the network during federated  
 642 learning. The  $\Delta$ FLQ algorithm trains the LSTM model with just a 5% degradation in  
 643 the validation loss while reducing by a factor of  $19\times$  the total data transmitted over the  
 644 network during federated learning.

645

## 646 6. Additional Experiments

647 In this section, we conduct additional experiments to assess the performance of the  
 648 FLQ and  $\Delta$ FLQ algorithms on different datasets, namely image (Section 5.4) and sensor  
 649 data (Section 5.4). All experiments are conducted on the same experimental framework  
 650 detailed in Section 4. Given the results reported in Section 5, we limit our experiments  
 651 to 2-bit ITERQ quantization scheme and 2 edge servers, because IterQ yielded the best  
 652 performance among the quantization methods considered.

653

### 654 6.1. Experiments on the MNIST image dataset

655 We conduct experiments on the popular MNIST dataset. The dataset consists of a  
 656 training set of 60,000 and a test set of 10,000  $28 \times 28$  gray-scale images. We randomly  
 657 sample the 10% of the training set to be used as the validation set during training.  
 658 The remaining training set is uniformly split among the edge servers. We use a CNN  
 659 composed of two  $5 \times 5$  convolutional layers, followed by two linear layers. We choose  
 660 rectified linear units as activation functions, and we use  $2 \times 2$  max-pooling layers and two  
 661 dropout layers. Training is performed using a batch size of 10 and stochastic gradient  
 662 descent with a learning rate of 0.01 and momentum 0.5. The training of the CNN is  
 663 performed by minimizing the cross-entropy loss until the learning rate decreases below

664 the minimum value of  $10^{-4}$  (early-stopping condition) or for a maximum of 80 federated  
 665 learning rounds.

Table 7: Minimum validation, test losses, rounds to reach the minimum validation loss and the total data (both directions) exchanged between the cloud node and a single edge server (in MB) of the CNN network model in the FL and FLQ scenario using the FLQ and  $\Delta$ FLQ algorithms with 2-bit quantization scheme and 2 workers, by varying  $\tau$ .

$\tau$	Min. val. loss			Test loss			Rounds			Total data		
	FL	FLQ	$\Delta$ FLQ	FL	FLQ	$\Delta$ FLQ	FL	FLQ	$\Delta$ FLQ	FL	FLQ	$\Delta$ FLQ
1	0.0509	0.0596	0.0638	0.0351	0.0453	0.0460	15	16	11	23.10	2.00	1.37
2	0.0412	0.0461	0.0502	0.0357	0.0408	0.0400	14	13	11	21.56	1.62	1.37
4	0.0451	0.0482	0.0507	0.0353	0.0370	0.0357	10	10	7	15.40	1.25	0.87
8	0.0451	0.0441	0.0499	0.0320	0.0329	0.0340	8	8	7	12.32	1.00	0.87
16	0.0428	0.0429	0.0495	0.0308	0.0333	0.0381	5	7	5	7.70	0.87	0.62

666 Our CNN model is composed of 24,090 parameters in total. We do not apply quan-  
 667 tization to the 90 one-dimensional bias parameters. The 4D tensors of the convolutional  
 668 layers have been decomposed into matrices along the first two dimensions and each ma-  
 669 trix has been quantized independently. The uncompressed CNN model requires 770.88  
 670 KB  $\approx$  0.77 MB, while the 2-bit ITERQ quantization scheme stores the network in 62.40  
 671 KB  $\approx$  0.06 MB, with a space reduction of  $12.35\times$ . This higher space reduction w.r.t. the  
 672 LSTM model investigated in Section 5 takes into account the larger weight matrices in  
 673 the fully connected layers of the CNN model.

674 Table 7 reports the validation and test losses with the FLQ and  $\Delta$ FLQ algorithms.  
 675 Comparing the corresponding values with the validation and test losses in the FL scenario,  
 676 where no quantization is employed, the networks, trained by the two algorithms for  
 677 different  $\tau$  values, exhibit slightly worse performance. Moreover, Table 7 reports the  
 678 number of federated rounds required to reach the minimum validation loss with the FLQ  
 679 and  $\Delta$ FLQ algorithms. Recall that this number represents the optimal number of rounds  
 680 to obtain the best CNN network performance during training before the validation loss  
 681 starts increasing and the local models start overfitting on the local data. The number of  
 682 rounds determines how many times the CNN model is transferred from an edge server to  
 683 the cloud node and vice-versa. The table also reports the number of optimal rounds in  
 684 the FL scenario.

685 In the FL scenario, the best validation loss is 0.0428, obtained with  $\tau = 16$  and with  
 686 5 federated learning rounds, with a total of  $7.70 \text{ MB} \times 2 = 15.40 \text{ MB}$  transferred (taking  
 687 into account both edge servers). When the 2-bit ITERQ quantization scheme is adopted,  
 688 the best validation losses are 0.0429 (in 7 federated rounds) and 0.0495 (in 5 federated  
 689 rounds) for the FLQ and  $\Delta$ FLQ algorithms, respectively, obtained with  $\tau = 16$ . Hence,  
 690 FLQ transfers a total of 1.74 MB, while  $\Delta$ FLQ a total of 1.24 MB over the network,  
 691 taking into account both edge servers.

692 Finally, when used with the 2-bit ITERQ quantization scheme and using 16 epochs  
 693 for local training at each worker, the FLQ algorithm trains the CNN model with almost  
 694 no degradation in the validation loss while reducing by a factor of  $8.8\times$  the total data  
 695 transmitted over the network during federated learning. The  $\Delta$ FLQ algorithm trains the  
 696 CNN model with a 15% degradation in the validation loss while reducing by a factor of  
 697  $12.4\times$  the total data transmitted over the network during federated learning.

699 *6.2. Experiments on the Bar Crawl dataset*

700 We perform additional experiments on time series data by employing the Bar Crawl  
 701 dataset<sup>5</sup>. The dataset collects acquisitions from the accelerometer and the transdermal  
 702 alcohol content from a college bar crawl and it is used in literature to predict heavy drink-  
 703 ing episodes via mobile data [38]. For our experiments, we employ the raw data collected  
 704 from smartphones’ accelerometers at a sampling rate of 40Hz. We choose to make use of  
 705 this dataset because of its size (i.e., several millions of observations collected). In detail,  
 706 each data point of the accelerometer’s observations contains five fields, i.e., the times-  
 707 tamp, a participant ID, and a sample from each of the three axes of the accelerometer.  
 708 Data was collected from a mix of 11 iPhones and 2 Android phones. The whole dataset  
 709 is fully anonymized.

710 We preprocess the data by deriving, from each acquisition, the average acceleration.  
 711 We then perform a per-user min-max standardization to normalize the data in the range  
 712  $[0, 1]$ . We employ the derived time series to learn a model that predicts the “next”  
 713 acceleration value, i.e., given the sequence of previous  $N$  acquisitions we query the model  
 714 to predict the  $N + 1$  value. We model this task as a regression problem and we employ  
 715 an LSTM to solve it. The employed LSTM is made up of one layer with 2 hidden nodes  
 716 and one output node. We train the network by employing the Mean Squared Error loss  
 717 function and by using a window of  $N = 16$  previous acquisitions. We split the dataset in  
 718 train/validation/test set on a per-user basis, i.e., we employ 8 users as the training set,  
 719 2 users as the validation set and 3 users as test set. The preprocessed dataset divided in  
 720 train/validation/test is made available to allow for reproducibility of the results<sup>6</sup>. The  
 721 size of the model is 1,376 bytes for the non-compressed model, while the 2-bit ITERQ  
 722 quantization scheme stores the network in 1,096 bytes. Differently from the previous  
 723 experiments, we train the network for a maximum of 40 federated learning rounds with a  
 724 fixed learning rate. Early-stopping is performed when no improvement in the validation  
 725 loss is observed for 5 local training epochs. We also employ a batch size of 256 training  
 726 samples and a learning rate of 0.01.

Table 8: Minimum validation, test losses, rounds to reach the minimum validation loss and the total data (both directions) exchanged between the cloud node and a single edge server (in KB) of the LSTM network model in the FL and FLQ scenario using the FLQ and  $\Delta$ FLQ algorithms with 2-bit quantization scheme and 2 workers, by varying  $\tau$ .

$\tau$	Min. val. loss			Test loss			Rounds			Total data		
	FL	FLQ	$\Delta$ FLQ	FL	FLQ	$\Delta$ FLQ	FL	FLQ	$\Delta$ FLQ	FL	FLQ	$\Delta$ FLQ
1	0.0784	0.0956	0.1576	0.107	0.112	0.160	40	40	40	110.08	87.68	87.68
2	0.0711	0.0759	0.1550	0.101	0.103	0.151	40	40	22	110.08	87.68	48.22
4	0.0691	0.0728	0.1489	0.091	0.098	0.150	40	40	7	110.08	87.68	15.34
8	0.0690	0.0714	0.1191	0.089	0.092	0.121	24	38	4	66.04	83.29	8.76
16	0.0687	0.0701	0.0841	0.086	0.089	0.099	36	34	3	99.07	73.16	6.45

727 Table 8 reports the validation and test losses with the FL, FLQ, and  $\Delta$ FLQ algo-  
 728 rithms. Comparing the corresponding values with the validation and test losses, in the

<sup>5</sup><https://archive.ics.uci.edu/ml/datasets/Bar+Crawl:+Detecting+Heavy+Drinking>

<sup>6</sup>[http://hpc.isti.cnr.it/~nardini/datasets/data\\_timeseries.tar.gz](http://hpc.isti.cnr.it/~nardini/datasets/data_timeseries.tar.gz)

729 FL scenario, where no quantization is employed, the networks, trained by the two algo-  
730 rithms for different  $\tau$  values, exhibit slightly worse performance. Table 8 also reports the  
731 number of federated rounds required to reach the minimum validation loss with the FL,  
732 FLQ, and  $\Delta$ FLQ algorithms. We recall that this number represents the optimal num-  
733 ber of rounds to obtain the best LSTM network performance during training before the  
734 local models start overfitting on the local data and the validation loss starts increasing.  
735 The number of rounds determines how many times the LSTM model is transferred from  
736 an edge server to the cloud node and vice-versa. The table also reports the number of  
737 optimal rounds in the FL scenario.

738 In the FL scenario, the best validation loss is 0.0687, obtained with  $\tau = 16$  and with  
739 36 federated rounds, with a total of  $99.07 \text{ KB} \times 2 = 198.14 \text{ KB}$  transferred (taking into  
740 account both edge servers). The use of the 2-bit ITERQ quantization scheme allows us  
741 to achieve 0.0701 minimum validation loss (in 36 federated rounds) and 0.0841 minimum  
742 validation loss (in 3 federated rounds) for the FLQ and  $\Delta$ FLQ algorithms, respectively,  
743 obtained with  $\tau = 16$ . Hence, FLQ transfers a total of 146.32 KB, while  $\Delta$ FLQ transmits  
744 a total of 12.90 KB over the network, taking into account both edge servers.

745 To conclude, the 2-bit ITERQ quantization scheme allows to effectively train an LSTM  
746 model for the “next value” prediction task on time series data. In detail, when using  
747 16 epochs for local training at each edge server, the FLQ algorithm trains the LSTM  
748 model with no degradation in the validation loss while reducing by 25% the total data  
749 transmitted over the network during federated learning, and the  $\Delta$ FLQ algorithm trains  
750 the LSTM model with a 20% degradation in the validation loss while reducing by a factor  
751 of  $16\times$  the total data transmitted over the network during federated learning.

## 752 7. Conclusion

753 In this work, we presented a federated learning platform for the IoT built upon an edge  
754 computing framework and according to a publish/subscribe communication paradigm.  
755 We introduced the edge computing layer to let the neural network reside as close as pos-  
756 sible to the IoT data producers and within the privacy domain of IoT device owners.  
757 Data privacy has been assumed by design as a mandatory requirement, leading to data  
758 retention on the edge. We developed two novel federated learning quantization algo-  
759 rithms, namely FLQ and  $\Delta$ FLQ. The FLQ algorithm introduces quantization into the  
760 state-of-the-art FL algorithm. Instead,  $\Delta$ FLQ is designed from scratch to reduce the  
761 amount of traffic exchanged among edge servers and the Cloud. Differently from [23],  
762 our FLQ and  $\Delta$ FLQ algorithms quantize both the NN model broadcasted by the cloud  
763 node and the NN models sent by the edge servers to significantly reduce the background  
764 traffic related to the NN model training. The results achieved by  $\Delta$ FLQ outperform  
765 those of FLQ both in terms of validation loss and efficiency in data exchange and show  
766 good tolerance to server faults. We conducted experiments on public datasets by train-  
767 ing a Long Short Term Memory neural network on an edge node to solve a next word  
768 prediction task as a case study. Only the federation task is left to the Cloud, or, eventu-  
769 ally, to a centralized entity, not necessarily located in the core network. To validate the  
770 generality of the proposed algorithms, we also applied both  $\Delta$ FLQ and FLQ to different  
771 use cases and other deep learning models. In detail, we experimented our proposed al-  
772 gorithms with CNNs on image classification tasks using the popular MNIST dataset and  
773 with LSTMs on next value prediction using time series from the Bar Crawl dataset. In

774 both cases, the  $\Delta$ FLQ was able to reduce by  $8 - 16\times$  the total amount of traffic among  
775 edge servers and the Cloud. Last but not least, the proposed federated learning approach  
776 allows for reducing the time to reach the minimum validation loss w.r.t. a centralized  
777 Cloud approach, with a negligible exchange of data at each training round, thanks to  
778  $\Delta$ FLQ and FLQ algorithms. The comprehensive experimental evaluation compares our  
779 approach to state-of-the-art centralized learning, i.e., on the Cloud, namely local learning,  
780 and federated learning techniques employing both full precision and several quantization  
781 schemes. We remark that local learning violates the requirement of data privacy. Still, it  
782 stands here as the optimal performance metric of comparison, in terms of validation loss,  
783 when the whole dataset is centralized. The measurement campaign shows that the in-  
784 troduction of quantization techniques in federated learning allows to significantly reduce  
785 the data exchanged between each edge server and a cloud node, providing a measured  
786 improvement with a minimal impact on the validation loss of the final model. In detail,  
787 FLQ outperforms full-precision federated learning of an LSTM for next-word prediction  
788 by reducing by a factor of  $13\times$  the total data transmitted over the network, but leads  
789 to an increase in validation loss by 14%. On the other hand,  $\Delta$ FLQ can further reduce  
790 the total data transmitted up to  $19\times$ , with a validation loss increase of less than 5%.  
791 Subsection 5.4 also provides a numerical quantification of the effective reduction in data  
792 volume exchanged with the 2- and 3-bit quantization schemes. Such promising results  
793 pave the way for deepening research on federated learning quantization, exploring how  
794 different schemes and techniques could impact the efficiency of the federated learning  
795 process.  
796 Future investigations will also address more in-depth communication models taking into  
797 account 5G mobility scenarios and investigate the fault-tolerance of federated learning  
798 architectures in real emerging contexts.

## 799 Acknowledgements

800 This work is partially supported by the Italian Ministry of Education and Research  
801 (MIUR) in the framework of the CrossLab project (Departments of Excellence), by the  
802 BIGDATAGRAPES and the TEACHING projects funded by the EU Horizon 2020 re-  
803 search and innovation program under grant agreements No. 780751 and No. 871385,  
804 respectively, and by the OK-INSAID project funded by the Italian Ministry of Education  
805 and Research (MIUR) under grant agreement No. ARS01\_00917. This work is partially  
806 carried out in the framework of the project AUTENS (Sustainable Energy Autarky)  
807 funded by the University of Pisa (PRA 2020 program).

## 808 References

- 809 [1] Virraji Mothukuri, Reza M. Parizi, Seyedamin Pouriyeh, Yan Huang, Ali Dehghantanha, and Gau-  
810 tam Srivastava. A survey on security and privacy of federated learning. *Future Generation Computer*  
811 *Systems*, 115:619–640, 2021.
- 812 [2] Miguel De Prado, Jing Su, Rabia Saeed, Lorenzo Keller, Noelia Vallez, Andrew Anderson, David  
813 Gregg, Luca Benini, Tim Llewellynn, Nabil Ouerhani, Rozenn Dahyot, and Nuria Pazos. Bonseyes  
814 AI Pipeline – Bringing AI to You: End-to-End Integration of Data, Algorithms, and Deployment  
815 Tools. *ACM Transaction on Internet of Things*, 1(4), 2020.
- 816 [3] Farzad Samie, Lars Bauer, and Jörg Henkel. From cloud down to things: An overview of machine  
817 learning in internet of things. *IEEE Internet of Things Journal*, 6(3):4921–4934, 2019.

- 818 [4] Mehdi Mohammadi, Ala Al-Fuqaha, Sameh Sorour, and Mohsen Guizani. Deep learning for iot big  
819 data and streaming analytics: A survey. *IEEE Communications Surveys & Tutorials*, 20(4):2923–  
820 2960, 2018.
- 821 [5] Weishan Zhang, Wuwu Guo, Xin Liu, Yan Liu, Jiehan Zhou, Bo Li, Qinghua Lu, and Su Yang.  
822 LSTM-based analysis of industrial IoT equipment. *IEEE Access*, 6:23551–23560, 2018.
- 823 [6] Jed Mills, Jia Hu, and Geyong Min. Communication-efficient federated learning for wireless edge  
824 intelligence in iot. *IEEE Internet of Things Journal*, 7(7):5986–5994, 2019.
- 825 [7] Peter Kairouz and et al. Advances and open problems in federated learning. *ArXiv*, abs/1912.04977,  
826 2019.
- 827 [8] He Li, Kaoru Ota, and Mianxiong Dong. Learning IoT in edge: Deep learning for the Internet of  
828 Things with edge computing. *IEEE Network*, 32(1):96–101, 2018.
- 829 [9] Jie Tang, Dawei Sun, Shaoshan Liu, and Jean-Luc Gaudiot. Enabling deep learning on IoT devices.  
830 *Computer*, 50(10):92–96, 2017.
- 831 [10] Ahmed Imteaj and M. Hadi Amini. Distributed sensing using smart end-user devices: Pathway to  
832 federated learning for autonomous iot. In *Proc. CSCI*, pages 1156–1161, 2019.
- 833 [11] Felix A. Gers, Douglas Eck, and Jürgen Schmidhuber. Applying LSTM to time series predictable  
834 through time-window approaches. In *Neural Nets WIRN Vietri-01*, pages 193–200. Springer, 2002.
- 835 [12] Qian Chen, Xiaodan Zhu, Zhen-Hua Ling, Si Wei, Hui Jiang, and Diana Inkpen. Enhanced LSTM  
836 for Natural Language Inference. In *Proc. ACL*, pages 1657–1668, 2017.
- 837 [13] Alexandre Alahi, Kratarth Goel, Vignesh Ramanathan, Alexandre Robicquet, Li Fei-Fei, and Silvio  
838 Savarese. Social lstm: Human trajectory prediction in crowded spaces. In *Proceedings of the IEEE  
839 conference on computer vision and pattern recognition*, pages 961–971, 2016.
- 840 [14] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. ImageNet Classification with Deep Con-  
841 volutional Neural Networks. In *Proc. NIPS*, pages 1097–1105, USA, 2012.
- 842 [15] Russell Reed. Pruning algorithms-a survey. *IEEE transactions on Neural Networks*, 4(5):740–747,  
843 1993.
- 844 [16] Shuchang Zhou, Yuxin Wu, Zekun Ni, Xinyu Zhou, He Wen, and Yuheng Zou. DoReFa-Net:  
845 Training Low Bitwidth Convolutional Neural Networks with Low Bitwidth Gradients. *arXiv e-  
846 prints*, page arXiv:1606.06160, 2016.
- 847 [17] Itay Hubara, Matthieu Courbariaux, Daniel Soudry, Ran El-Yaniv, and Yoshua Bengio. Quantized  
848 neural networks: Training neural networks with low precision weights and activations. *The Journal  
849 of Machine Learning Research*, 18(1):6869–6898, 2017.
- 850 [18] Suyog Gupta, Ankur Agrawal, Kailash Gopalakrishnan, and Pritish Narayanan. Deep learning with  
851 limited numerical precision. In *Proc. ICML*, pages 1737–1746, 2015.
- 852 [19] Mohammad Rastegari, Vicente Ordonez, Joseph Redmon, and Ali Farhadi. Xnor-net: Imagenet  
853 classification using binary convolutional neural networks. In *Proc. ECCV*, pages 525–542. Springer,  
854 2016.
- 855 [20] Itay Hubara, Matthieu Courbariaux, Daniel Soudry, Ran El-Yaniv, and Yoshua Bengio. Binarized  
856 neural networks. In *Proc. NIPS*, pages 4107–4115, 2016.
- 857 [21] Chenzhuo Zhu, Song Han, Huizi Mao, and William J Dally. Trained ternary quantization. In *Proc.  
858 ICLR*, 2017.
- 859 [22] Aojun Zhou, Anbang Yao, Yiwen Guo, Lin Xu, and Yurong Chen. Incremental network quantiza-  
860 tion: Towards lossless CNNs with low-precision weights. In *Proc. ICLR*, 2017.
- 861 [23] Chen Xu, Jianqiang Yao, Zhouchen Lin, Wenwu Ou, Yuanbin Cao, Zhirong Wang, and Hongbin  
862 Zha. Alternating multi-bit quantization for recurrent neural networks. In *Proc. ICLR*, 2018.
- 863 [24] Arash Ardakani, Zhengyun Ji, Sean C Smithson, Brett H Meyer, and Warren J Gross. Learning  
864 recurrent binary/ternary weights. In *Proc. ICLR*, 2019.
- 865 [25] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas.  
866 Communication-efficient learning of deep networks from decentralized data. In *Proc. AISTATS*,  
867 2017.
- 868 [26] Jakub Konečný, H Brendan McMahan, Daniel Ramage, and Peter Richtárik. Federated optimiza-  
869 tion: Distributed machine learning for on-device intelligence. *arXiv preprint arXiv:1610.02527*,  
870 2016.
- 871 [27] Amirhossein Reisizadeh, Aryan Mokhtari, Hamed Hassani, Ali Jadbabaie, and Ramtin Pedarsani.  
872 Fedpaq: A communication-efficient federated learning method with periodic averaging and quanti-  
873 zation. In *Proc. PMLR*, pages 2021–2031, 2020.

- 874 [28] Jasenka Dizdarević, Francisco Carpio, Admela Jukan, and Xavi Masip-Bruin. A survey of communi-  
875 cation protocols for internet of things and related challenges of fog and cloud computing integration.  
876 *ACM Computing Surveys (CSUR)*, 51(6):1–29, 2019.
- 877 [29] Manlio Bacco, Marco Colucci, and Alberto Gotta. Application protocols enabling internet of remote  
878 things via random access satellite channels. In *Proc. ICC*, pages 1–6, 2017.
- 879 [30] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016.
- 880 [31] Matthieu Courbariaux, Yoshua Bengio, and Jean-Pierre David. Binaryconnect: Training deep  
881 neural networks with binary weights during propagations. In *Proc. NIPS*, pages 3123–3131, 2015.
- 882 [32] Mohammad Rastegari, Vicente Ordonez, Joseph Redmon, and Ali Farhadi. Xnor-net: Imagenet  
883 classification using binary convolutional neural networks. In *Proc. ECCV*, 2016.
- 884 [33] Jakub Konečný, H. Brendan McMahan, Felix X. Yu, Peter Richtarik, Ananda Theertha Suresh,  
885 and Dave Bacon. Federated learning: Strategies for improving communication efficiency. In *Proc.*  
886 *NIPS Workshop on Private Multi-Party Machine Learning*, 2016.
- 887 [34] Dan Alistarh, Demjan Grubic, Jerry Z. Li, Ryota Tomioka, and Milan Vojnovic. Qsgd:  
888 Communication-efficient sgd via gradient quantization and encoding. In *Proc. NIPS*, page  
889 1707–1718, 2017.
- 890 [35] Geoffrey Davis, Stephane G. Mallet, and Marco Avellaneda. Adaptive greedy approximations.  
891 *Constructive Approximation*, 13(1):57–98, 1997.
- 892 [36] Yiwen Guo, Anbang Yao, Hao Zhao, and Yurong Chen. Network sketching: Exploiting binary  
893 structure in deep cnns. In *Proc. CVPR*, pages 4040–4048, 2017.
- 894 [37] Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture  
895 models. *arXiv preprint arXiv:1609.07843*, 2016.
- 896 [38] Jackson A. Killian, Kevin M. Passino, Arnab Nandi, Danielle R. Madden, John D. Clapp, Nir-  
897 malie Wiratunga, Frans Coenen, and Sadiq Sani. Learning to detect heavy drinking episodes using  
898 smartphone accelerometer data. In *KHD@IJCAI*, pages 35–42, 2019.