

1

## 2 **Supplementary Information for**

### 3 **Communication Efficient Federated Learning**

4 **Mingzhe Chen, Nir Shlezinger, H. Vincent Poor, Yonina C. Eldar, and Shuguang Cui**

5 **H. Vincent Poor**

6 **E-mail: [poor@princeton.edu](mailto:poor@princeton.edu)**

#### 7 **This PDF file includes:**

8     Supplementary text

9     Figs. S1 to S7

10    SI References

## 11 Supporting Information Text

### 12 1. Architecture of Wireless Federated Learning

13 Next, we introduce a practical scenario for the implementation of an federated learning (FL) algorithm, as shown in Fig. S1. In  
 14 this figure,  $U$  devices and one central controller (base station) cooperatively participate in the FL training process. During the  
 15 training process, each device first uses its collected data to train its local FL model. Then each device transmits its trained local  
 16 FL parameters to the base station (BS). The BS will aggregate the received local FL parameters to generate the global FL  
 17 model and broadcast it back to all devices. Since local FL parameters and the global FL model are transmitted over wireless  
 18 links, the FL training performance will be affected by wireless network performance. In particular, the number of resource  
 19 blocks (RBs) is limited and hence, the number of devices that can participate in FL is limited. Meanwhile, from Fig. S1,  
 20 we can see that the devices in other service areas may use the same RBs to transmit data thus affecting the FL parameter  
 21 transmission delay and FL convergence time.

### 22 2. Universal FL Model Parameter Compression

23 **A. Rationale.** Compression can be modeled as an encoding-decoding system. To faithfully represent the FL setup, we design  
 24 our quantization strategy in light of the following requirements and assumptions:

- 25 *A1* All devices share the same encoding function. This requirement significantly simplifies FL implementation.
- 26 *A2* No *a-priori* knowledge or distribution of the model updates is assumed.
- 27 *A3* The devices and the central controller (CC) share a source of common randomness. This is achieved by, e.g., letting the  
 28 CC share with each device a random seed along with the weights. Once a different seed is conveyed to each device, it can  
 29 be used to obtain a dedicated source of common randomness shared by the CC and each of the devices for the entire FL  
 30 procedure. These seeds can be conveyed along with the weight.

31 Requirement *A2* gives rise to the need for a *universal quantization* approach, namely, a scheme which operates reliably  
 32 regardless of the distribution of the FL parameter updates and without its prior knowledge.

33 **B. Model Compression Algorithm.** Here, we present the encoding and decoding functions. Following requirement *A1*, we utilize  
 34 universal vector quantization, i.e., a quantization scheme which maps each set of continuous-amplitude values into a discrete  
 35 representation in a manner which is ignorant of the underlying distribution. The source of common randomness assumed in *A3*  
 36 implies that the CC and the devices can generate the same realizations of a dither signal. We thus use a compression based on  
 37 dithered vector quantization, and particularly, on lattice quantization, detailed in the following.

38 Let  $L$  be a fixed positive integer, referred to henceforth as the lattice dimension, and let  $\mathbf{G}$  be a non-singular  $L \times L$  matrix,  
 39 which denotes the lattice generator matrix. For simplicity, we assume that  $M \triangleq \frac{m}{L}$  is an integer, where  $m$  is the number of  
 40 elements in the weight matrices of the FL parameters used to represent an FL model, although the scheme can also be applied  
 41 when this does not hold by replacing  $M$  with  $\lceil M \rceil$ . Next, we use  $\mathcal{L}$  to denote the lattice, which is the set of points in  $\mathbb{R}^L$  that  
 42 can be written as an integer linear combination of the columns of  $\mathbf{G}$ , i.e.,

$$43 \quad \mathcal{L} \triangleq \{\mathbf{x} = \mathbf{G}\mathbf{l} : \mathbf{l} \in \mathbb{Z}^L\}. \quad [1]$$

44 A lattice quantizer  $Q_{\mathcal{L}}(\cdot)$  maps each  $\mathbf{x} \in \mathbb{R}^L$  to its nearest lattice point, i.e.,  $Q_{\mathcal{L}}(\mathbf{x}) = \mathbf{l}_x$  where  $\mathbf{l}_x \in \mathcal{L}$  if  $\|\mathbf{x} - \mathbf{l}_x\| \leq \|\mathbf{x} - \mathbf{l}\|$   
 45 for every  $\mathbf{l} \in \mathcal{L}$ . Finally, let  $\mathcal{P}_0$  be the basic lattice cell, i.e., the set of points in  $\mathbb{R}^L$  which are closer to  $\mathbf{0}$  than to any other  
 46 lattice point:

$$47 \quad \mathcal{P}_0 \triangleq \{\mathbf{x} \in \mathbb{R}^L : \|\mathbf{x}\| < \|\mathbf{x} - \mathbf{p}\|, \forall \mathbf{p} \in \mathcal{L} \setminus \{\mathbf{0}\}\}. \quad [2]$$

48 For example, when  $\mathbf{G} = \Delta \cdot \mathbf{I}_L$  for some  $\Delta > 0$ , then  $\mathcal{L}$  is the square lattice, for which  $\mathcal{P}_0$  is the set of vectors  $\mathbf{x} \in \mathbb{R}^L$  whose  
 49  $\ell_{\infty}$  norm is not larger than  $\frac{\Delta}{2}$ . For this setting,  $Q_{\mathcal{L}}(\cdot)$  implements entry-wise scalar uniform quantization with spacing  $\Delta$ .

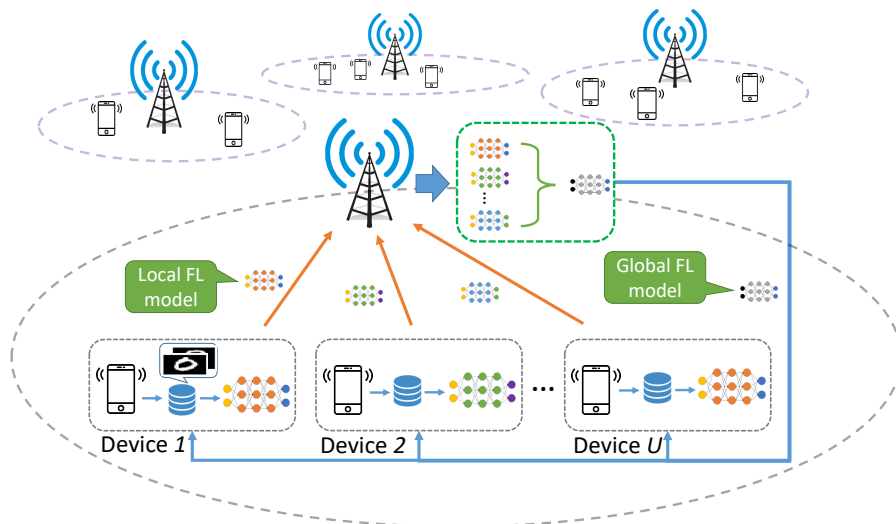
50 Define the model update of the  $k$ th device at iteration  $t + 1$  as  $\mathbf{h}_{t+1}^{(k)} \triangleq \sigma_{k,t+1}^{\tau} \mathbf{h}_t^{(k)} - \mathbf{b}_t$ . Using the above definitions in lattice  
 51 quantization, the encoding and decoding procedures, which are based on subtractive dithered lattice quantization, consist of  
 52 the following steps:

53 **Encoder:** The encoding function includes the following steps:

54 *E1 Normalize and partition:* The  $k$ th device scales  $\mathbf{h}_{t+1}^{(k)}$  by  $\zeta \|\mathbf{h}_{t+1}^{(k)}\|$  for some  $\zeta > 0$ , and divides the result into  $M$   
 55 distinct  $L \times 1$  vectors, denoted  $\{\bar{\mathbf{h}}_i^{(k)}\}_{i=1}^M$ . The scalar quantity  $\zeta \|\mathbf{h}_{t+1}^{(k)}\|$  is quantized separately from  $\{\bar{\mathbf{h}}_i^{(k)}\}_{i=1}^M$  using some  
 56 fine-resolution quantizer.

57 *E2 Dithering:* The encoder utilizes the source of common randomness, e.g., a shared seed, to generate the set of  $L \times 1$   
 58 dither vectors  $\{\mathbf{z}_i^{(k)}\}_{i=1}^M$ , which are randomized in an i.i.d. fashion, independently of  $\mathbf{h}_{t+1}^{(k)}$ , from a uniform distribution  
 59 over  $\mathcal{P}_0$ .

60 *E3 Quantization:* The vectors  $\{\bar{\mathbf{h}}_i^{(k)}\}_{i=1}^M$  are discretized by adding the dither vectors and applying lattice quantization, i.e.,  
 61 by computing  $\{Q_{\mathcal{L}}(\bar{\mathbf{h}}_i^{(k)} + \mathbf{z}_i^{(k)})\}$ .



**Fig. S1.** The architecture of an FL algorithm that is implemented over a wireless network that consists of multiple devices and one base station.

62 **E4 Entropy coding:** The discrete values  $\{Q_{\mathcal{L}}(\bar{\mathbf{h}}_i^{(k)} + \mathbf{z}_i^{(k)})\}$  are encoded into a digital codeword  $u_{t+1}^{(k)}$  in a lossless manner.

63 In order to utilize entropy coding in step E4, the discretized  $\{Q_{\mathcal{L}}(\bar{\mathbf{h}}_i^{(k)} + \mathbf{z}_i^{(k)})\}$  must take values on a *finite set*. This is  
 64 achieved by the normalization in Step E1, which guarantees that  $\{\bar{\mathbf{h}}_i^{(k)}\}_{i=1}^M$  all reside inside the  $L$ -dimensional ball with radius  
 65  $\zeta^{-1}$ , in which the number of lattice points is not larger than  $\frac{\pi^L/2}{\zeta^L \Gamma(1+L/2) \det(\mathcal{G})}$  (1, Ch. 2), where  $\Gamma(\cdot)$  is the Gamma function.  
 66 The overhead in accurately quantizing the single scalar quantity  $\zeta \|\mathbf{h}^{(k)}\|$  is typically negligible compared to the number of bits  
 67 required to convey the set of vectors  $\{\bar{\mathbf{h}}_i^{(k)}\}_{i=1}^M$ , hardly affecting the overall quantization rate.

68 **Decoder:** The decoding mapping implements the following:

69 **D1 Entropy decoding:** The CC first decodes each digital codeword  $u_{t+1}^{(k)}$  into the discrete value  $\{Q_{\mathcal{L}}(\bar{\mathbf{h}}_i^{(k)} + \mathbf{z}_i^{(k)})\}$ . Since  
 70 the encoding is carried out using a lossless source code, the discrete values are recovered without any errors.

71 **D2 Dither subtraction:** Using the source of common randomness, the CC generates the dither vectors  $\{\mathbf{z}_i^{(k)}\}$ , which can  
 72 be carried out rapidly and at low complexity using random number generators as the dither vectors obey a uniform  
 73 distribution. The CC then subtracts the corresponding vector from each lattice point, i.e., compute  $\{Q_{\mathcal{L}}(\bar{\mathbf{h}}_i^{(k)} + \mathbf{z}_i^{(k)}) - \mathbf{z}_i^{(k)}\}$ .  
 74 An illustration of the subtractive dithered lattice quantization procedure is illustrated in Fig. S2.

75 **D3 Collecting and scaling:** The values  $\{Q_{\mathcal{L}}(\bar{\mathbf{h}}_i^{(k)} + \mathbf{z}_i^{(k)}) - \mathbf{z}_i^{(k)}\}$  are collected into an  $m \times 1$  vector  $\hat{\mathbf{h}}_{t+1}^{(k)}$  using the inverse  
 76 operation of the partitioning and normalization in Step E1.

77 **D4 Model recovery:** The recovered matrices are combined into an updated model. Namely,

$$78 \quad \mathbf{b}_{t+1} = \frac{1}{\sum_k N_k} \sum_k N_k \hat{\mathbf{h}}_{t+1}^{(k)} + \mathbf{b}_t. \quad [3]$$

79 The usage of subtractive dithered lattice quantization in Steps E2-E3 and D2 allow obtaining a digital representation which  
 80 is relatively close to the true quantity, as illustrated in Fig. S2, without relying on prior knowledge of its distribution. The joint  
 81 decoding aspect of the proposed scheme is introduced in the final model recovery Step D4. The remaining encoding-decoding  
 82 procedure, i.e., Steps E1-D3 is carried out independently for each device.

### 83 3. Optimization of Resource Block Allocation

84 Next, we explain how to solve the RB allocation optimization problem. We first define the data rate of device  $i$  transmitting its  
 85 compressed local FL model to the CC as  $c_i^U(\chi_{i,t}) = \sum_{r=1}^R \chi_{i,t}^r B \log_2 \left(1 + \frac{P h_i}{I_r + B N_0}\right)$  where  $\chi_{i,t}^r \in \{0, 1\}$  is the RB allocation index  
 86 with  $\chi_{i,t}^r = 1$  implying that RB  $r$  is allocated to device  $i$  at iteration  $t$ , otherwise, we have  $\chi_{i,t}^r = 0$ , and  $\boldsymbol{\chi}_{i,t} = [\chi_{i,t}^1, \dots, \chi_{i,t}^R]$ .  
 87 Then, the transmission delay at each FL iteration is given by  $\max_{i \in \mathcal{U}_{\mathbf{p}_t}} \frac{Z}{c_i^U(\chi_{i,t})}$ , where  $\mathcal{U}_{\mathbf{p}_t}$  is the subset of devices that transmit  
 88 their compressed local FL model parameters to the CC at iteration  $t$  and  $Z$  is the data size of FL parameters. Given these  
 89 definitions, we rewrite the optimization problem as follows:

$$90 \quad \min_{\boldsymbol{\chi}_t} \max_{i \in \mathcal{U}_{\mathbf{p}_t}} \frac{Z}{c_i^U(\boldsymbol{\chi}_{i,t})} \quad [4]$$

$$91 \quad \text{s. t. } \chi_{i,t}^r \in \{0, 1\}, \quad \forall i \in \mathcal{U}_{\mathbf{p}_t}, r \in \mathcal{R}, \quad [4a]$$

$$92 \quad \sum_{i \in \mathcal{U}_{\mathbf{p}_t}} \chi_{i,t}^r = 1, \forall r \in \mathcal{R}, \quad [4b]$$

$$93 \quad \sum_{r=1}^R \chi_{i,t}^r = 1, \forall i \in \mathcal{U}_{\mathbf{p}_t}, \quad [4c]$$

94 where  $\mathcal{R}$  is the set of RBs that can be allocated to the devices and  $\boldsymbol{\chi}_t = [\boldsymbol{\chi}_{1,t}, \dots, \boldsymbol{\chi}_{|\mathcal{U}_{\mathbf{p}_t}|,t}]$  with  $|\mathcal{U}_{\mathbf{p}_t}|$  being the number of  
 95 devices in  $\mathcal{U}_{\mathbf{p}_t}$ . We assume that a variable  $q$  exists such that  $\frac{Z}{c_i^U(\boldsymbol{\chi}_{i,t})} \leq q, \forall i \in \mathcal{U}_{\mathbf{p}_t}$ . Hence, we can rewrite the problem in  
 96 Eq. (4) as

$$97 \quad \min_{\boldsymbol{\chi}_t} q \quad [5]$$

$$98 \quad \text{s. t. } \chi_{i,t}^r \in \{0, 1\}, \quad \forall i \in \mathcal{U}_{\mathbf{p}_t}, r \in \mathcal{R}, \quad [5a]$$

$$99 \quad \sum_{i \in \mathcal{U}_{\mathbf{p}_t}} \chi_{i,t}^r = 1, \forall r \in \mathcal{R}, \quad [5b]$$

$$100 \quad \sum_{r=1}^R \chi_{i,t}^r = 1, \forall i \in \mathcal{U}_{\mathbf{p}_t}, \quad [5c]$$

$$101 \quad q \geq \frac{Z}{c_i^U(\boldsymbol{\chi}_{i,t})}. \quad [5d]$$

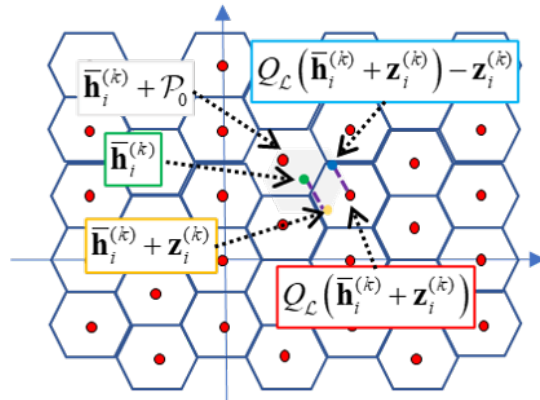


Fig. S2. Subtractive dithered lattice quantization illustration.

95 Here,  $q \geq \frac{Z}{c_i^U(\mathbf{x}_{i,t})}$  is nonconvex. Therefore, we first need to transform it to a convex equation. We assume that the data rate  
 96 of device  $i$  using RB  $r$  to transmit local FL model parameters is  $c_{ir,t}^U = B \log_2 \left(1 + \frac{P h_i}{I_r + B N_0}\right)$  where  $I_r$  is the interference over RB  
 97  $r$ . Then, we have  $c_i^U(\mathbf{x}_{i,t}) = \sum_{r=1}^R \chi_{i,t}^r c_{ir,t}^U$ . In consequence, the problem in Eq. (5) can be expressed by

$$98 \quad \min_{\mathbf{x}_t} q \quad [6]$$

$$\text{s. t. } \chi_{i,t}^r \in \{0, 1\}, \quad \forall i \in \mathcal{U}_{\mathbf{p}_t}, r \in \mathcal{R}, \quad [6a]$$

$$\sum_{i \in \mathcal{U}_{\mathbf{p}_t}} \chi_{i,t}^r = 1, \forall r \in \mathcal{R}, \quad [6b]$$

$$\sum_{r=1}^R \chi_{i,t}^r = 1, \forall i \in \mathcal{U}_{\mathbf{p}_t}, \quad [6c]$$

$$\sum_{r=1}^R \chi_{i,t}^r c_{ir,t}^U \geq \frac{Z}{q}. \quad [6d]$$

99 The problem in Eq. (6) is an integer linear programming problem, which can be solved by using Matlab toolbox. In this paper,  
 100 we use Matlab intlinprog function to solve the problem in Eq. (6).

#### 101 4. Convergence Analysis

102 Here, we can analyze the convergence, proving Theorem 1. Our proof follows a similar outline to that used in (2, 3), with the  
 103 introduction of additional arguments for handling the quantization constraints. The unique characteristics of the quantization  
 104 error which arise from the dithered strategy allow us to rigorously incorporate its contribution into the overall flow of the proof.  
 105 As a preliminary step, we first recall the assumptions in light of which our analysis is carried out:

106 *AS1* The expected squared  $\ell_2$  norm of the random vector  $\nabla f_k^i(\mathbf{o})$ , representing the stochastic gradient evaluated at  $\mathbf{b}$ , is  
 107 bounded by some  $\xi^2 > 0$  for all  $\mathbf{o} \in \mathbb{R}^m$ .

*AS2* The local objective functions  $\{f_k(\cdot)\}$  are all  $\rho_s$ -smooth, namely, for all  $\mathbf{v}_1, \mathbf{v}_2 \in \mathbb{R}^m$  it holds that

$$f_k(\mathbf{v}_1) - f_k(\mathbf{v}_2) \leq (\mathbf{v}_1 - \mathbf{v}_2)^T \nabla f_k(\mathbf{v}_2) + \frac{1}{2} \rho_s \|\mathbf{v}_1 - \mathbf{v}_2\|^2.$$

*AS3* The local objective functions  $\{f_k(\cdot)\}$  are all  $\rho_c$ -strongly convex, namely, for all  $\mathbf{v}_1, \mathbf{v}_2 \in \mathbb{R}^m$  it holds that

$$f_k(\mathbf{v}_1) - f_k(\mathbf{v}_2) \geq (\mathbf{v}_1 - \mathbf{v}_2)^T \nabla f_k(\mathbf{v}_2) + \frac{1}{2} \rho_c \|\mathbf{v}_1 - \mathbf{v}_2\|^2.$$

108 *AS4* The probabilistic device selection method selects the set  $\mathcal{U}_{\mathbf{p}_t}$  of devices, such that  $|\mathcal{U}_{\mathbf{p}_t}| = R$  and  $\mathcal{U}_{\mathbf{p}_t}$  is uniformly  
 109 distributed over all  $R$ -sized subsets of  $\mathcal{U}$ .

**A. Recursive Bound on Weights Error.** From (4), it follows that the effect of subtractive dithered quantization can be modeled  
 as additive noise, independent of the quantized value, whose distribution depends only on the properties of the lattice. In  
 particular, it holds that the distortion induced in quantizing the model update  $\mathbf{h}_t^{(k)}$ , denoted  $\boldsymbol{\epsilon}_t^{(k)}$ , is an  $m \times 1$  zero-mean  
 additive noise vector *independent of*  $\mathbf{h}_{t+1}^{(k)}$ . Consequently, by defining the sequence  $\mathbf{e}_t^{(k)}$  such that  $\mathbf{e}_t^{(k)} = \boldsymbol{\epsilon}_t^{(k)}$  if  $t$  is an integer  
 multiple of  $\tau$  and  $\mathbf{e}_t^{(k)} = \mathbf{0}$  otherwise, the instantaneous weights at the  $k$ th device, defined as  $\tilde{\mathbf{o}}_t^{(k)} \triangleq \mathbf{o}_{k, \lfloor t/\tau \rfloor \tau}^{(k)}$ , can be written as

$$\tilde{\mathbf{o}}_{t+1}^{(k)} = \begin{cases} \tilde{\mathbf{o}}_t^{(k)} - \tilde{\lambda}_t \nabla f_k^{i_t^{(k)}}(\tilde{\mathbf{o}}_t^{(k)}) + \mathbf{e}_{t+1}^{(k)} & t+1 \notin \mathcal{T}_\tau, \\ \sum_{k'=1}^U \alpha_{k'} \left( \tilde{\mathbf{o}}_t^{(k')} - \tilde{\lambda}_t \nabla f_k^{i_t^{(k')}}(\tilde{\mathbf{o}}_t^{(k')}) + \mathbf{e}_{t+1}^{(k')} \right) & t+1 \in \mathcal{T}_\tau, \end{cases} \quad [7]$$

110 where  $\tilde{\lambda}_t \triangleq \lambda_{\lfloor t/\tau \rfloor \tau}^{t - \lfloor t/\tau \rfloor \tau}$ ,  $\mathcal{T}_\tau$  is the set of integer multiples of  $\tau$ , and  $f_k^i$  is the objective evaluated at the  $i$ th sample of the  $k$ th  
 111 device.

112 The equivalent model update representation in Eq. (7) allows us to model the effect of subtractive dithered quantization  
 113 on the overall FL procedure as additional noise corrupting the computation of the stochastic gradients. Building upon this  
 114 representation, we now follow the strategy proposed in (2) and adapted to heterogeneous data in (3). This is achieved by  
 115 defining a virtual sequence  $\{\mathbf{v}_t\}$  from  $\{\tilde{\mathbf{o}}_t^{(k)}\}$  which can be shown to behave almost like mini-batch SGD with batch size  $\tau$ ,  
 116 while being within a bounded distance of the FL model weights  $\{\tilde{\mathbf{o}}_t^{(k)}\}$ , by properly setting the step size  $\tilde{\lambda}_t$ . In particular, we  
 117 define the virtual sequence  $\{\mathbf{v}_t\}$  via

$$\bar{\mathbf{v}}_t \triangleq \begin{cases} \sum_{k=1}^U \alpha_k \tilde{\mathbf{o}}_t^{(k)} & t \notin \mathcal{T}_\tau, \\ \sum_{k \in \mathcal{U}_{\mathbf{p}_t/\tau}} \alpha_k \tilde{\mathbf{o}}_t^{(k)} & t \in \mathcal{T}_\tau, \end{cases} \quad [8]$$

which coincides with  $\tilde{\mathbf{o}}_t^{(k)}$  when  $t$  is an integer multiple of  $\tau$ . Also, let  $\mathbf{v}_t \triangleq \sum_{k=1}^U \alpha_k \tilde{\mathbf{o}}_t^{(k)}$  be the virtual sequence representing the averaged model over all devices (both participating and non-participating) at each time instance. Further define the averaged noisy stochastic gradients and the averaged full gradients as

$$\tilde{\mathbf{g}}_t \triangleq \sum_{k=1}^U \alpha_k \left( \nabla f_k^{i_t^{(k)}}(\tilde{\mathbf{o}}_t^{(k)}) - \frac{1}{\lambda_t} \mathbf{e}_{t+1}^{(k)} \right), \quad [9a]$$

$$\mathbf{g}_t \triangleq \sum_{k=1}^U \alpha_k \nabla f_k(\tilde{\mathbf{o}}_t^{(k)}), \quad [9b]$$

119 respectively. Note that since the quantization error is zero-mean and each mini-batch consists of a single sample, whose  
 120 indexes  $\{i_t^{(k)}\}$  are independent and uniformly distributed, it holds that  $\mathbb{E}\{\tilde{\mathbf{g}}_t\} = \mathbf{g}_t$ . Additionally, the virtual sequence satisfies  
 121  $\mathbf{v}_{t+1} = \mathbf{v}_t - \tilde{\lambda}_t \tilde{\mathbf{g}}_t$ . We use the following lemmas, proved in (3, Appendix B.4).

122 **Lemma 1** Under assumption AS4,  $\bar{\mathbf{v}}_t$  is an unbiased estimation of  $\mathbf{v}_t$ , i.e.  $\mathbb{E}_{\mathcal{U}_{\mathbf{p}_t/\tau}}\{\bar{\mathbf{v}}_t\} = \mathbf{v}_t$ .

123 **Lemma 2** The expected difference between  $\mathbf{v}_t$  and  $\bar{\mathbf{v}}_t$  is bounded by

$$\mathbb{E}_{\mathcal{U}_{\mathbf{p}_t/\tau}}\{\|\bar{\mathbf{v}}_t - \mathbf{v}_t\|^2\} \leq \frac{4(U-R)}{(U-1)R} \eta_t^2 \tau^2 \xi^2. \quad [10]$$

We next use these lemmas to bound the distance between the FL model parameters and the optimal one, as

$$\begin{aligned} \|\bar{\mathbf{v}}_{t+1} - \mathbf{b}^*\|^2 &= \|\bar{\mathbf{v}}_{t+1} - \mathbf{v}_{t+1} + \mathbf{v}_{t+1} - \mathbf{b}^*\|^2 \\ &= \|\bar{\mathbf{v}}_{t+1} - \mathbf{v}_{t+1}\|^2 + \|\mathbf{v}_{t+1} - \mathbf{b}^*\|^2 + \underbrace{2\langle \bar{\mathbf{v}}_{t+1} - \mathbf{v}_{t+1}, \mathbf{v}_{t+1} - \mathbf{b}^* \rangle}_A. \end{aligned} \quad [11]$$

125 The term  $\mathbb{E}_{\mathcal{U}_{\mathbf{p}_t/\tau}}\{A\} = 0$  since  $\bar{\mathbf{v}}_t$  is unbiased by Lemma 1. Further, using Lemma 2, it follows from Eq. (11) that

$$\mathbb{E}\{\|\bar{\mathbf{v}}_{t+1} - \mathbf{b}^*\|^2\} \leq \mathbb{E}\{\|\mathbf{v}_{t+1} - \mathbf{b}^*\|^2\} + \frac{4(U-R)}{(U-1)R} \eta_t^2 \tau^2 \xi^2. \quad [12]$$

The resulting term which has to be bounded in Eq. (12) is thus equivalent to that used in (3), and as a result, by assumptions AS2-AS3, it follows from (3, Lemma 1) that if  $\tilde{\lambda}_t \leq \frac{1}{4\rho_s}$  then

$$\mathbb{E}\{\|\mathbf{v}_{t+1} - \mathbf{b}^*\|^2\} \leq (1 - \tilde{\lambda}_t \rho_c) \mathbb{E}\{\|\mathbf{v}_t - \mathbf{b}^*\|^2\} + 6\rho_s \tilde{\lambda}_t^2 \psi + \tilde{\lambda}_t^2 \mathbb{E}\{\|\tilde{\mathbf{g}}_t - \mathbf{g}_t\|^2\} + 2\mathbb{E}\left\{\sum_{k=1}^U \alpha_k \left\| \mathbf{v}_t - \tilde{\mathbf{o}}_t^{(k)} \right\|^2\right\}. \quad [13]$$

127 Eq. (13) bounds the expected distance between the virtual sequence  $\{\mathbf{v}_t\}$  and the optimal weights  $\mathbf{b}^*$  in a recursive manner.  
 128 We further bound the summands in Eq. (13), using the following lemmas:

129 **Lemma 3** If the step size  $\tilde{\lambda}_t$  is non-increasing and satisfies  $\tilde{\lambda}_t \leq 2\tilde{\lambda}_{t+\tau}$  for each  $t \geq 0$ , then, when assumption AS1 is satisfied,  
 130 it holds that

$$\tilde{\lambda}_t^2 \mathbb{E}\{\|\tilde{\mathbf{g}}_t - \mathbf{g}_t\|^2\} \leq (1 + 4M\zeta^2 \bar{\sigma}_L^2 \tau^2) \tilde{\lambda}_t^2 \xi^2 \sum_{k=1}^U \alpha_k^2. \quad [14]$$

132 **Lemma 4** If the step size  $\tilde{\lambda}_t$  is non-increasing and satisfies  $\tilde{\lambda}_t \leq 2\tilde{\lambda}_{t+\tau}$  for each  $t \geq 0$ , then, when assumption AS1 is satisfied,  
 133 it holds that

$$\mathbb{E}\left\{\sum_{k=1}^U \alpha_k \left\| \mathbf{v}_t - \tilde{\mathbf{o}}_t^{(k)} \right\|^2\right\} \leq 4(\tau - 1)^2 \tilde{\lambda}_t^2 \xi^2. \quad [15]$$

135 Next, we define  $\delta_t \triangleq \mathbb{E}\{\|\bar{\mathbf{v}}_t - \mathbf{b}^*\|^2\}$ . When  $t \in \mathcal{T}_\tau$ , the term  $\delta_t$  represents the  $\ell_2$  norm of the error in the weights of the  
 136 global model. Using Lemmas 3-4, while substituting Eq. (15) and Eq. (14) into Eq. (13) and Eq. (12), we obtain the following  
 137 recursive relationship on the weights error:

$$\delta_{t+1} \leq (1 - \tilde{\lambda}_t \rho_c) \delta_t + \tilde{\lambda}_t^2 \varpi, \quad [16]$$

where

$$\varpi \triangleq (1 + 4M\zeta^2 \bar{\sigma}_L^2 \tau^2) \xi^2 \sum_{k=1}^U \alpha_k^2 + 6\rho_s \psi + 8(\tau - 1)^2 \xi^2 + \frac{4(U-R)}{(U-1)R} \tau^2 \xi^2.$$

139 The relationship in Eq. (16) is used in the sequel to prove the FL convergence bound stated in Theorem 1.

140 **B. FL Convergence Bound.** Here, we prove Theorem 1 based on the recursive relationship in Eq. (16). This is achieved by  
 141 properly setting the step-size and the FL systems parameters in Eq. (16) to bound  $\delta_t = \mathbb{E} \{ \|\bar{\mathbf{v}}_t - \mathbf{b}^*\|^2 \}$ , and combining the  
 142 resulting bound with the strong convexity of the objective AS3 to prove the theorem.

143 In particular, we set the step size  $\tilde{\lambda}_t$  to take the form  $\tilde{\lambda}_t = \frac{\beta}{t+\gamma}$  for some  $\beta > 0$  and  $\gamma \geq \max(4\rho_s\beta, \tau)$ , for which  $\tilde{\lambda}_t \leq \frac{1}{4\rho_s}$   
 144 and  $\tilde{\lambda}_t \leq 2\tilde{\lambda}_{t+\tau}$ , implying that Eq. (13) and Eq. (15) hold.

Under such settings, we show that there exists a finite  $\nu$  such that  $\delta_t \leq \frac{\nu}{t+\gamma}$  for all integer  $l \geq 0$ . We prove this by induction,  
 noting that setting  $\nu \geq \gamma\delta_0$  guarantees that it holds for  $t = 0$ . Consequently, we next show that if  $\delta_t \leq \frac{\nu}{t+\gamma}$ , then  $\delta_{t+1} \leq \frac{\nu}{t+1+\gamma}$ .  
 It follows from Eq. (16) that

$$\begin{aligned} \delta_{t+1} &\leq \left(1 - \frac{\beta}{t+\gamma}\rho_c\right) \frac{\nu}{t+\gamma} + \left(\frac{\beta}{t+\gamma}\right)^2 b \\ &= \frac{1}{t+\tau} \left( \left(1 - \frac{\beta}{t+\gamma}\rho_c\right) \nu + \frac{\beta^2}{t+\gamma} b \right). \end{aligned} \quad [17]$$

Consequently,  $\delta_{t+1} \leq \frac{\nu}{t+1+\gamma}$  holds when

$$\frac{1}{t+\tau} \left( \left(1 - \frac{\beta}{t+\gamma}\rho_c\right) \nu + \frac{\beta^2}{t+\gamma} \varpi \right) \leq \frac{\nu}{t+1+\gamma},$$

145 or, equivalently,

$$\left(1 - \frac{\beta}{t+\gamma}\rho_c\right) \nu + \frac{\beta^2}{t+\gamma} \varpi \leq \frac{t+\gamma}{t+1+\gamma} \nu. \quad [18]$$

By setting  $\nu \geq \frac{1+\beta^2\varpi}{\beta\rho_c}$ , the left hand side of Eq. (18) satisfies

$$\begin{aligned} \left(1 - \frac{\beta}{t+\gamma}\rho_c\right) \nu + \frac{\beta^2}{t+\gamma} \varpi &= \frac{t-1+\gamma}{t+\gamma} \nu + \left(\frac{1-\beta\rho_c}{t+\gamma} \nu + \frac{\beta^2}{t+\gamma} \varpi\right) \\ &= \frac{t-1+\gamma}{t+\gamma} \nu + \frac{1}{t+\gamma} \left( (1-\beta\rho_c) \nu + \beta^2 \varpi \right) \\ &\stackrel{(a)}{\leq} \frac{t-1+\gamma}{t+\gamma} \nu, \end{aligned} \quad [19]$$

where (a) holds since  $\nu \geq \frac{1+\beta^2\varpi}{\beta\rho_c}$ . As the right hand side of Eq. (19) is not larger than that of Eq. (18), it follows that Eq. (18)  
 holds for the current setting, which in turn proves that  $\delta_{t+1} \leq \frac{\nu}{t+1+\gamma}$ . Finally, the smoothness of the objective AS2 implies that

$$\begin{aligned} \mathbb{E}\{f(\mathbf{b}_t)\} - f(\mathbf{b}^*) &= \mathbb{E}\{f(\bar{\mathbf{v}}_{t\tau})\} - f(\mathbf{b}^*) \\ &\leq \frac{\rho_s}{2} \delta_{t\tau} \leq \frac{\rho_s \nu}{2(t\tau + \gamma)}, \end{aligned} \quad [20]$$

147 which, in light of the above setting, holds for  $\nu \geq \max\left(\frac{1+\beta^2\varpi}{\beta\rho_c}, \gamma\delta_0\right)$ ,  $\gamma \geq \max(\tau, 4\beta\rho_s)$ , and  $\beta > 0$ . In particular, setting  
 148  $\beta = \frac{\tau}{\rho_c}$  results in  $\gamma \geq \tau \max(1, 4\rho_s/\rho_c)$  and  $\nu \geq \max\left(\frac{\rho_c + \tau^2\varpi}{\tau\rho_c}, \gamma\delta_0\right)$ , which, when substituted into Eq. (20), proves the theorem.  
 149 ■

150 **B.1. Deferred Proofs.** Here we detail the proofs of the intermediate lemmas used for obtaining the recursion Eq. (16).

**Proof of Lemma 3** To prove Eq. (14), we note that since the quantization noise and the stochastic gradients are mutually  
 independent, it follows from the definition of the gradient vectors in Eq. (9) that

$$\begin{aligned} \tilde{\lambda}_t^2 \mathbb{E} \{ \|\tilde{\mathbf{g}}_t - \mathbf{g}_t\|^2 \} &= \sum_{k=1}^U \alpha_k^2 \mathbb{E} \left\{ \left\| \mathbf{e}_{t+1}^{(k)} \right\|^2 \right\} + \tilde{\lambda}_t^2 \sum_{k=1}^U \alpha_k^2 \mathbb{E} \left\{ \left\| \nabla f_k^{i_t^{(k)}}(\tilde{\mathbf{o}}_t^{(k)}) - \nabla f_k(\tilde{\mathbf{o}}_t^{(k)}) \right\|^2 \right\} \\ &\stackrel{(a)}{\leq} \sum_{k=1}^U \alpha_k^2 \mathbb{E} \left\{ \left\| \mathbf{e}_{t+1}^{(k)} \right\|^2 \right\} + \tilde{\lambda}_t^2 \xi^2 \sum_{k=1}^U \alpha_k^2, \end{aligned} \quad [21]$$

where (a) holds since the uniform distribution of the random index  $i_k$  implies that the expected value of the stochastic gradient is  
 the full gradient, i.e.,  $\mathbb{E}\{\nabla f_k^{i_t^{(k)}}(\mathbf{b})\} = \nabla f_k(\mathbf{b})$ , and consequently,  $\mathbb{E}\{\|\nabla f_k^{i_t^{(k)}}(\tilde{\mathbf{o}}_t^{(k)}) - \nabla f_k(\tilde{\mathbf{o}}_t^{(k)})\|^2\} \leq \mathbb{E}\{\|\nabla f_k^{i_t^{(k)}}(\tilde{\mathbf{o}}_t^{(k)})\|^2\} \leq \xi^2$



by assumption [AS1](#). Furthermore, the definition of  $\mathbf{e}_{t+1}^{(k)}$  implies that  $\mathbb{E}\{\|\mathbf{e}_{t+1}^{(k)}\|^2\} = 0$  for  $t+1 \notin \mathcal{T}$ , while for  $t+1 \in \mathcal{T}$  it holds that  $\mathbb{E}\{\|\mathbf{e}_{t+1}^{(k)}\|^2\} = \mathbb{E}\{\|\boldsymbol{\epsilon}_{t+1}^{(k)}\|^2\} = M\sigma_{\mathcal{L}}^2$ . Now, the quantization error satisfies

$$\begin{aligned} \mathbb{E}\{\|\mathbf{e}_{t+1}^{(k)}\|^2\} &\leq M\zeta^2\bar{\sigma}_{\mathcal{L}}^2\mathbb{E}\left\{\left\|\sum_{t'=t+1-\tau}^{t+1}\tilde{\lambda}_{t'}\nabla f_k^{i_{t'}^{(k)}}(\tilde{\boldsymbol{o}}_{t'}^{(k)})\right\|^2\right\} \\ &\stackrel{(a)}{\leq} M\zeta^2\bar{\sigma}_{\mathcal{L}}^2\tau\sum_{t'=t+1-\tau}^{t+1}\tilde{\lambda}_{t'}^2\mathbb{E}\left\{\left\|\nabla f_k^{i_{t'}^{(k)}}(\tilde{\boldsymbol{o}}_{t'}^{(k)})\right\|^2\right\} \\ &\stackrel{(b)}{\leq} M\zeta^2\bar{\sigma}_{\mathcal{L}}^2\tau^2\tilde{\lambda}_{t+1-\tau}^2\xi^2\stackrel{(c)}{\leq} 4M\zeta^2\bar{\sigma}_{\mathcal{L}}^2\tau^2\tilde{\lambda}_t^2\xi^2, \end{aligned} \quad [22]$$

151 where in (a) we used the inequality  $\|\sum_{t'=t+1-\tau}^{t+1}\mathbf{r}_{t'}\|^2 \leq \tau\sum_{t'=t+1-\tau}^{t+1}\|\mathbf{r}_{t'}\|^2$ , which holds for any multivariate sequence  
152  $\{\mathbf{r}_t\}$ ; (b) is obtained from assumption [AS1](#); and (c) follows since  $\tilde{\lambda}_{t+1-\tau} \leq 2\tilde{\lambda}_{t+1} \leq 2\tilde{\lambda}_t$ . Substituting Eq. (22) into Eq. (21)  
153 proves the lemma. ■

**Proof of Lemma 4** Note that for  $t_0 = \lfloor t/\tau \rfloor \tau$ , which is an integer multiple of  $\tau$ , it holds that  $\mathbf{v}_{t_0} = \tilde{\boldsymbol{o}}_{t_0}^{(k)}$ . Since Eq. (15) trivially holds for  $t = t_0$ , we henceforth focus on the case where  $t > t_0$ . We now write

$$\begin{aligned} \mathbb{E}\left\{\sum_{k=1}^U\alpha_k\left\|\tilde{\boldsymbol{o}}_t^{(k)}-\mathbf{v}_t\right\|^2\right\} &= \mathbb{E}\left\{\sum_{k=1}^U\alpha_k\left\|\tilde{\boldsymbol{o}}_t^{(k)}-\tilde{\boldsymbol{o}}_{t_0}^{(k)}-(\mathbf{v}_t-\mathbf{v}_{t_0})\right\|^2\right\} \\ &\stackrel{(a)}{\leq} \mathbb{E}\left\{\sum_{k=1}^U\alpha_k\left\|\tilde{\boldsymbol{o}}_t^{(k)}-\tilde{\boldsymbol{o}}_{t_0}^{(k)}\right\|^2\right\} \\ &= \sum_{k=1}^U\alpha_k\mathbb{E}\left\{\left\|\tilde{\boldsymbol{o}}_t^{(k)}-\tilde{\boldsymbol{o}}_{t_0}^{(k)}\right\|^2\right\}, \end{aligned} \quad [23]$$

where in (a) we used the fact that for every set  $\{\mathbf{r}^{(k)}\}$ , one can define a random vector  $\mathbf{r}$  such that  $\Pr(\mathbf{r} = \mathbf{r}^{(k)}) = \alpha_k$ , and thus

$$\begin{aligned} \sum_{k=1}^U\alpha_k\left\|\mathbf{r}^{(k)}-\sum_{l=1}^U\alpha_l\mathbf{r}^{(l)}\right\|^2 &= \mathbb{E}\{\|\mathbf{r}-\mathbb{E}\{\mathbf{r}\}\|^2\} \\ &\leq \mathbb{E}\{\|\mathbf{r}\|^2\} = \sum_{k=1}^U\alpha_k\|\mathbf{r}^{(k)}\|^2. \end{aligned}$$

Next, we recall that  $\mathbf{e}_{t'} = \mathbf{0}$  for each  $t' = t_0 + 1, \dots, t$ . Consequently, similarly to the derivation in Eq. (22),

$$\begin{aligned} \mathbb{E}\left\{\left\|\tilde{\boldsymbol{o}}_t^{(k)}-\tilde{\boldsymbol{o}}_{t_0}^{(k)}\right\|^2\right\} &= \mathbb{E}\left\{\left\|\sum_{t'=t_0}^{t-1}\tilde{\lambda}_{t'}\nabla f_k^{i_{t'}^{(k)}}(\tilde{\boldsymbol{o}}_{t'}^{(k)})\right\|^2\right\} \\ &\stackrel{(a)}{\leq} (\tau-1)\sum_{t'=t_0}^{t-1}\tilde{\lambda}_{t'}^2\mathbb{E}\left\{\left\|\nabla f_k^{i_{t'}^{(k)}}(\tilde{\boldsymbol{o}}_{t'}^{(k)})\right\|^2\right\} \\ &\stackrel{(b)}{\leq} (\tau-1)^2\tilde{\lambda}_{t_0}^2\xi^2\stackrel{(c)}{\leq} 4(\tau-1)^2\tilde{\lambda}_t^2\xi^2, \end{aligned} \quad [24]$$

154 where in (a) we used the inequality  $\|\sum_{t'=t_0}^{t-1}\mathbf{r}_{t'}\|^2 \leq (t-1-t_0)\sum_{t'=t_0}^{t-1}\|\mathbf{r}_{t'}\|^2 \leq (\tau-1)\sum_{t'=t_0}^{t-1}\|\mathbf{r}_{t'}\|^2$ , which holds for any  
155 multivariate sequence  $\{\mathbf{r}_t\}$ ; (b) is obtained from assumption [AS1](#); and (c) follows since  $\tilde{\lambda}_{t_0} \leq \tilde{\lambda}_{t-\tau} \leq 2\tilde{\lambda}_t$ . Substituting Eq. (24)  
156 into Eq. (23) proves the lemma. ■

## 157 5. Handwritten Digit Identification

158 Fig. [S3](#) shows the architecture of the FL model used for handwritten digit identification. From Fig. [S3](#), we can see that the FL  
159 model is a shallow fully-connected neural network that consists of 50 neurons. Since the size of one digit image is  $28 \times 28$ , the  
160 input size is 784. The purpose of training FL is to identify 10 class handwritten digits and hence, the output size is 10. The  
161 activation functions in the hidden and output layers are *trasiq* and softmax functions.

## 162 6. Object Recognition in Images

163 Fig. [S4](#) shows the architecture of the FL model used for object recognition in images. The FL model is a convolutional neural  
164 network (CNN) that consists of 16 layers. In this figure, Conv2D, MaxPooling, AvgPooling, and FullyConnected are short for  
165 2D convolutional layer, maximum pooling layer, average pooling layer, and fully connected layer. The stride of each pooling or  
166 convolutional layer is 1 unless specified otherwise.

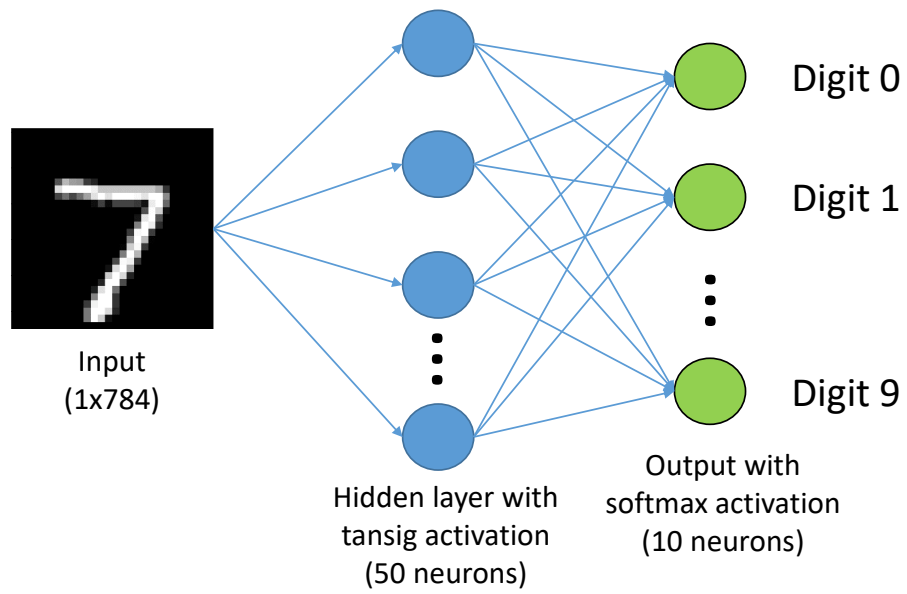


Fig. S3. The architecture of the FL model used for handwritten digit identification.

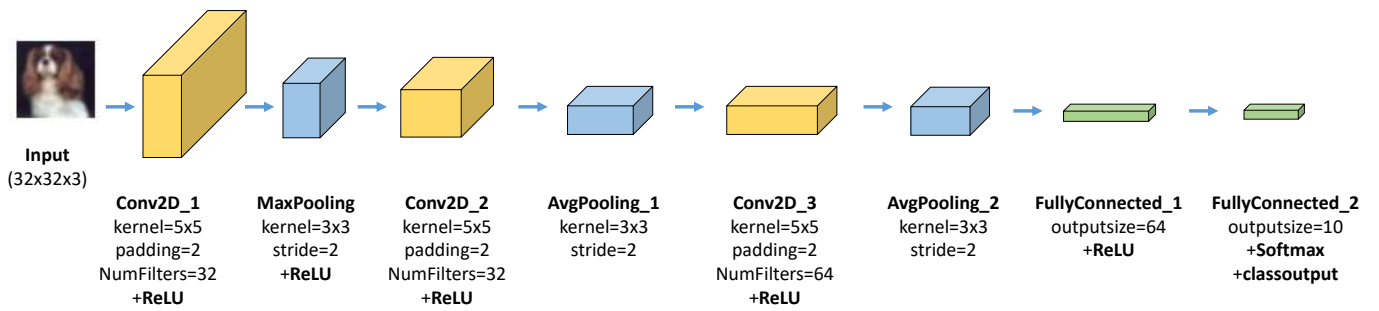


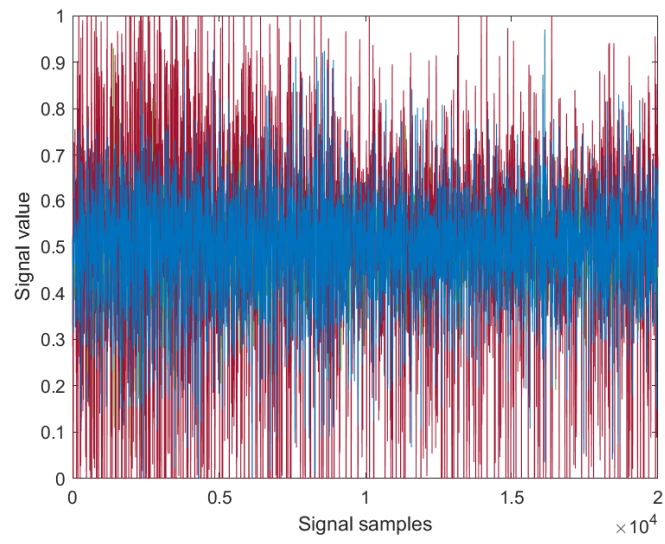
Fig. S4. The architecture of the FL model used for object recognition in images.

## 167 7. Finger Movement Detection

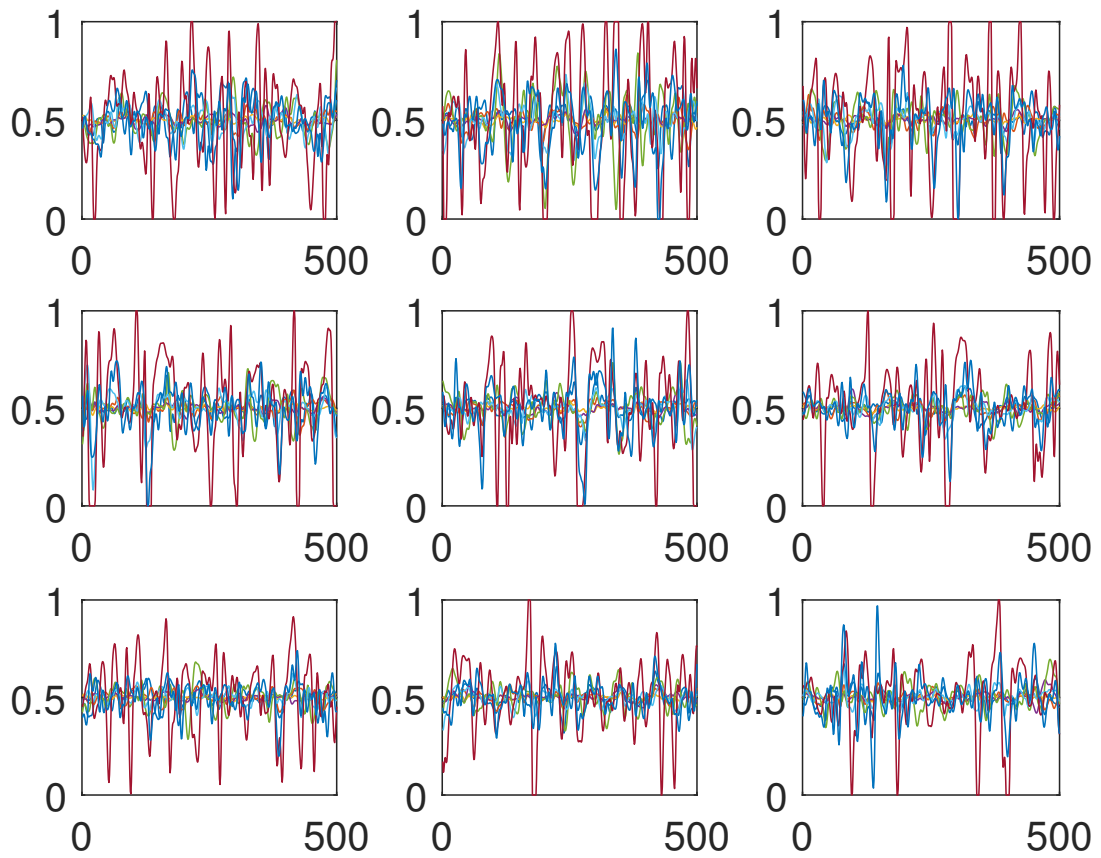
168 Fig. S5 shows the raw data collected from a hand close finger movement within five seconds. From Fig. S5, we can see that the  
169 raw data consists of 20000 signal samples. Hence, it is impossible to directly use the raw data as an input vector to train the  
170 FL model. A windowing method is used to split the raw data. We assume that the window size is 500 and the interval between  
171 two windows is 50. Therefore, 20000 signal samples can be divided into 391 windows. Fig. S6 shows the hand close finger  
172 movement data processed by the windowing method. From Fig. S6, we can see that the signal samples at each window are  
173 different which increases the complexity of training FL models. Fig. S7 shows the architecture of the FL model used for finger  
174 movement detection. The FL model is a CNN that consists of 19 layers. We consider the signal data in one window as an input  
175 vector of the CNN and hence, the size of a CNN input is  $500 \times 8$ . Since a CNN is used to detect 15 class finger movement, the  
176 size of a CNN output is  $15 \times 1$ . In Fig. S7, BatchNormalization is short for batch normalization layer. Meanwhile, the stride is  
177 a vector  $[2, 1]$  with 2 being the vertical stride and 1 being the horizontal stride.

## 178 References

- 179 1. JH Conway, NJA Sloane, *Sphere Packings, Lattices and Groups*. (Springer Science & Business Media) Vol. 290, (2013).
- 180 2. SU Stich, Local SGD converges fast and communicates little. *arXiv preprint arXiv:1805.09767* (2018).
- 181 3. X Li, K Huang, W Yang, S Wang, Z Zhang, On the convergence of fedavg on non-iid data. *arXiv preprint arXiv:1907.02189*  
182 (2019).
- 183 4. R Zamir, M Feder, On lattice quantization noise. *IEEE Transactions on Inf. Theory* **42**, 1152–1159 (1996).



**Fig. S5.** Raw data collected from a close finger movement.



**Fig. S6.** The close finger movement data processed by the windowing method.

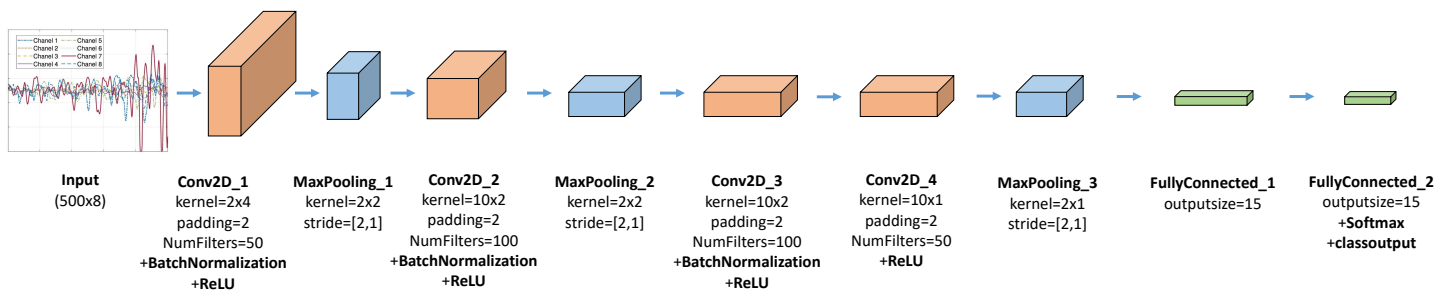


Fig. S7. The architecture of the FL model used for finger movement detection.