

Communication efficient model-aware federated learning for visual crowd counting and density estimation in smart cities

A. Armacki, N. Milosevic, D. Bajovic, S. Kar, D. Jakovetic, A. Bakhtiarnia, L. Esterle, A. Muscat, and T. Festi

Abstract—Federated learning (FL) is an attractive paradigm where a number of users can improve their local models via sharing trained models or model increments with a central server, while the users’ data is kept private. However, when model sizes are huge, FL incurs a significant communication overhead. In such scenarios, strategies that perform user sampling, so that only informative users communicate their models to the server, are desired. In this paper, we make several contributions on user sampling in FL. On the theoretical side, we consider a general framework that exhibits user heterogeneity across several dimensions: activation probabilities, gradient noise variance, number of updates per epoch, and communication channel quality. In this setting, we derive convergence rate of the FedAvg method. The rate explicitly characterizes the effects of heterogeneity and enables us to derive optimal user sampling probabilities in an offline setting, when the sampling probabilities are pre-computed. We then show how these derived probabilities naturally connect with existing optimized sampling strategies in an adaptive-online setting. On the practical side, we study visual crowd counting (VCC) as a representative deep learning application with huge-sized models. We provide an implementation of the FL system over real-world data across three pilot sites – Valetta, Trento and Novi Sad. The evaluation results demonstrate significant model accuracy benefits through employing FL over the multiple cities, and significant communication savings via non-uniform user sampling strategies.

Index Terms—Federated learning, Optimized client selection, Visual crowd counting, Communication efficient protocol.

I. INTRODUCTION

Federated learning (FL) refers to the concept where a number of users collaboratively learn a global machine learning model with the assistance of a global coordinator (server), e.g., [1]. During the FL process, users do not share their raw data

but only share locally trained models or model increments with the server, hence making FL attractive from data protection and privacy points of view in various applications [1], e.g., visual crowd counting (VCC) in smart cities [2].

In this paper, the communication-efficiency strategy of main interest is optimized user sampling, i.e., selecting a subset of appropriately defined most informative users at each training epoch. This is in contrast with the standard uniform sampling strategy in FL that selects S out of N users uniformly at random; see, e.g., [3]–[5]. In more detail, we consider a general non-uniform sampling framework for FL in the presence of various sources of system heterogeneity. Therein, each user i is active, i.e., transmits its model increment, with a user-dependent probability q_i . The users perform the federated averaging (FedAvg) training method [6], where each user at each epoch makes a user-dependent number K_i of local stochastic gradient (SGD) steps. The user-dependent K_i ’s model heterogeneous users’ capabilities in terms of computational and storage power, processor speed, etc, e.g., [7]. In addition, the noise variances introduced with users’ local SGD steps are user-dependent, hence modeling different data quality or different batch sizes across users. Finally, we allow that the user-server links may be unreliable, so that transmission of an active user is received at the server with a user-dependent probability k_i . The k_i ’s here account for communication channel imperfections, such as finite user transmit power, packet dropouts, etc.

There have been several recent works that consider non-uniform user sampling in FL and propose optimized user sampling strategies, e.g., [8]–[10]. A user sampling strategy can be considered in an offline setting, where user activation probabilities q_i ’s are determined beforehand. On the other hand, in an online setting, user activation probabilities are calculated adaptively at each training epoch based on the FL algorithm progress.

In this paper, we make several theoretical and practical contributions towards better understanding of how users should be sampled for communication-efficient FL. First, on the theoretical side, we consider a general heterogeneous offline user sampling framework, as described above. For this framework, we derive convergence rate of FedAvg assuming either strongly convex or convex user losses. We then explicitly quantify the achieved rate with respect to various sources of

A. Armacki and S. Kar are with Carnegie Mellon University, Pittsburgh, PA, USA. N. Milosevic and D. Jakovetic are with Faculty of Sciences, University of Novi Sad, Serbia. D. Bajovic is with Faculty of Technical Sciences, University of Novi Sad, Serbia. A. Bakhtiarnia and L. Esterle are with Aarhus University, Denmark. A. Muscat is with Greenroads, Malta. T. Festi is with Comune di Trento, Italy. Authors’ emails: aarmacki@andrew.cmu.edu, nmilosevic@dm.uns.ac.rs, dbajovic@uns.ac.rs, soumyak@andrew.cmu.edu, dusan.jakovetic@dm.uns.ac.rs, [arianbakh,lukas.esterle]@ece.au.dk, adrian@greenroadsmalta.com, thomas.festi@comune.trento.it. The work of NM, DB, AB, LE, AM and TF was supported in part by the European Union’s Horizon 2020 Research and Innovation program under grant agreement No 957337. The paper reflects only the view of the authors and the Commission is not responsible for any use that may be made of the information it contains. The work of NM, DB and DJ was also supported in part by the Serbian Ministry of Education, Science and Technological Development.

system heterogeneity. Next, optimal user activation probabilities q_i 's are derived that maximize the convergence rate subject to a given per-epoch communication budget and user-server link capacities, i.e., probabilities k_i 's. These results generalize existing bounds [8] to the case of heterogeneous local user variances, and, more importantly, to heterogeneous numbers of local SGD updates K_i 's. This reveals several useful insights regarding practical user sampling strategies in FL. First, a user sampling probability should be inversely proportional to its number of local updates K_i ; intuitively, performing more local work leads to better variance reduction and more informative updates, allowing users to communicate with the FL server less frequently. On the other hand, interestingly, the user activation probability should be proportional to its local variance, i.e., higher variance users should be polled more frequently. This is explained intuitively by the fact that users with more inherent local noise should be active more frequently in order to filter out the noise, in view of the underlying law of large numbers phenomenon. For the online (adaptive) client sampling setting, we describe how the results we derive are in accordance with the optimal adaptive user selection strategy in [8].

Second, on the implementation and application side, we develop a FL system for VCC in smart cities, with a real-world system deployment over three cities, Valetta, Trento, and Novi Sad. VCC is an important task with numerous applications in smart cities, such as people density estimation in public squares, pedestrian counting in traffic areas, and monitoring large-scale events by unmanned aerial vehicles (UAVs). The deployed FL system for VCC supports heterogeneous client sampling. Specifically, we employ the optimal adaptive (online) client sampling strategy from [8]. Extensive numerical experiments demonstrate that, for the training sites (cities) that lack sufficient data beforehand, FL leads to significant accuracy improvements of VCC. Moreover, the results demonstrate significant communication savings when non-uniform, optimized user sampling is performed.

II. PROBLEM SETUP

A. Federated learning model

We consider a FL system with N users that wish to collaboratively minimize the following function:

$$f(w) = \sum_{i=1}^N p_i f_i(w). \quad (1)$$

Here, quantity $0 < p_i < 1$, $\sum_{i=1}^N p_i = 1$, and $f_i : \mathbb{R}^d \mapsto \mathbb{R}$ is the local user i 's cost function, e.g., the empirical or population loss of user i with respect to its local data (local data distribution).

In order to solve (1) in the FL setting, users and the server perform the FedAvg algorithm with user-dependent number of local SGD updates and a heterogeneous probabilistic user activation scheme. We denote by $t = 0, 1, \dots$ the outer iteration (epoch) counter, and by $j = 0, 1, \dots$ the inner iteration index that counts the number of local SGD updates at each user. To be more precise, we denote by $w_{t,j}^i \in \mathbb{R}^d$ the user i 's estimate

of the solution to (1) at outer iteration (epoch) t and inner SGD iteration j , $j = 1, \dots, K_i$. Similarly, denote by $w_t \in \mathbb{R}^d$ the server's solution estimate at epoch t , $t = 0, 1, \dots$. Further, denote by $g_{t,j}^i$ user i 's noisy estimate of gradient $\nabla f_i(w_{t,j}^i)$ at point $w_{t,j}^i$, and let $g_t^i = \frac{1}{K_i} \sum_{j=0}^{K_i-1} g_{t,j}^i$. Next, let X_t^i be a Bernoulli random variable that encodes the information whether the server receives the locally updated model w_{t,K_i}^i from user i at epoch t . We let $X_t^i \sim \text{Bernoulli}(q_i)$, with parameter q_i . In other words, q_i is the probability that the server successfully receives the model w_{t,K_i}^i at the end of the inner iteration process at epoch t from user i . Quantity X_t^i in our setting models two effects: user i 's random activation (user decides to be active or not); and message reception failure if the model was transmitted but was not received by the server.¹ The algorithm works as follows. For all $i = 1, \dots, N$, at each epoch t , each user $i = 1, \dots, N$ performs the following update:

$$w_{t,j+1}^i = w_{t,j}^i - \alpha g_{t,j}^i, \quad (2)$$

for $j = 0, \dots, K - 1$. Here, $\alpha > 0$ is the step-size, and $w_{t,0}^i = w_t$. Then, the server aggregates model increments $g_t^i = \frac{w_{t,K_i}^i - w_{t,0}^i}{\alpha K_i}$ from all users i for which $X_t^i = 1$, and computes

$$w_{t+1} = w_t - \alpha \sum_{i=1}^N p_i \frac{X_{t+1}^i}{q_i} g_t^i. \quad (3)$$

Note that we implicitly assume that the server knows quantities q_i 's, $i = 1, \dots, N$. The division of each user's contributing term by q_i in (3) enables elimination (on average) of the non-uniform sampling bias, similarly to how dividing the local aggregate updates by $\frac{1}{K_i}$ ensures objective consistency, e.g., [7]. In other words, there holds: $\mathbb{E}[w_{t+1} | \mathcal{F}_t] = w_t - \alpha \sum_{i=1}^N p_i g_t^i$. Here, $\mathbb{E}[\cdot | \mathcal{F}_t]$ denotes conditional expectation, \mathcal{F}_t is the history of the algorithm, including all the algorithmic steps history prior to the generation of the instances of the variables X_t^i , $i = 1, \dots, N$.

B. Visual crowd counting

VCC corresponds to counting the total number of people that are present in a given scene. The input data X to the model is a high-resolution RGB colour image of a scene containing people, and the output Y that this model provides is a density map, which is a single-channel image that specifies the density of the crowd at each pixel of the input image. The values in this density map can be summed to get the total count in the form of a single number. The annotations used for training this model are in the form of head annotations, where the location of the center of each person's head is specified. The model utilizes the SASNet neural network architecture [11]. In other words, we employ a model-aware approach, where each FL user pre-assumes the specific SASNet architecture that is

¹More precisely, we can let $X_t^i = Z_t^i Y_t^i$, where Z_t^i and Y_t^i are both Bernoulli random variables, Z_t^i indicating whether agent i was selected to send its update during communication round t , while Y_t^i indicating if the communication was successful or not.

in prior works proven to be effective in VCC in a single-user setting [11]. Then, the weights of SASNet are learned via FL, harnessing data from all FL users. The FL users here correspond to $N = 3$ sites in 3 different cities (Vareta, Trento, Novi Sad). In view of (1), the unknown model $w \in \mathbb{R}^d$ thus represents the unknown weights of the SASNet neural network that maps inputs X to outputs Y , while function f_i represents the training loss with respect to the input-output data acquired at FL user i . We refer to [11] for further SASNet model and loss function details.

III. CONVERGENCE RATE AND OPTIMAL ACTIVATION

A. Convergence rate

We next present convergence rate results for the FL algorithm (2)–(3) under the heterogeneous setting described in Section II. For a cleaner presentation, we let $p_i = 1/N$, for all i (see (1)), while the results generalize to arbitrary p_i 's. We make the following assumptions.

Assumption 1. Each f_i is L -smooth and μ -strongly convex, i.e., the following inequalities are satisfied for all $x, y \in \mathbb{R}^d$, for some constants $0 < \mu \leq L$: $\|\nabla f_i(x) - \nabla f_i(y)\| \leq L\|x - y\|$, and $\langle \nabla f_i(x) - \nabla f_i(y), x - y \rangle \geq \mu\|x - y\|^2$, where $\|\cdot\|$ denotes the Euclidean norm.

Assumption 2. Each $g_{t,j}^i$ is an unbiased estimate of the true gradient, and the variance is bounded, i.e.:

$$\begin{aligned} \mathbb{E}[g_{t,j}^i | w_{t,j-1}^i] &= \nabla f_i(w_{t,j-1}^i), \\ \mathbb{E}[\|g_{t,j}^i - \nabla f_i(w_{t,j-1}^i)\|^2 | w_{t,j-1}^i] &\leq \sigma_i^2, \end{aligned}$$

for all $j = 1, \dots, K_i$, for some constant $\sigma_i^2 > 0$, $i = 1, \dots, N$.

Assumption 3. The indicator variables X_t^i are independent across time and across users.

For future reference, we introduce the following quantities: the maximal number of local updates $\bar{K} = \max_{i \in [N]} K_i$ (here, $[N] = \{1, 2, \dots, N\}$); the local communication probability parameter $\bar{q}_i = \frac{1-q_i}{q_i}$; the minimal communication probability parameter $q_{\min} = \min_{i \in [N]} q_i$. We also denote by w^* a solution to (1).

We have the following result. The proof is provided in Appendix A.

Theorem 1. Let assumptions 1-3 hold. For the step-size satisfying $\alpha \leq \min\left\{\frac{q_{\min}}{64L}, \frac{1}{10L\bar{K}}\right\}$, we have the following convergence guarantees:

1) For f general convex ($\mu = 0$), $R \geq 1$, there holds:

$$\mathbb{E}(f(\bar{w}_R) - f^*) = \mathcal{O}\left(\frac{A_1}{\sqrt{RN}} + \frac{A_2}{R^{2/3}} + \frac{A_3}{R}\right),$$

where $\bar{w}_R = \frac{1}{R} \sum_{r=0}^{R-1} w_r$, with the problem related constants being

$$A_1 = \sqrt{\frac{\Delta_0}{N} \sum_{i=1}^N \left(\frac{\sigma_i^2}{q_i K_i} + \bar{q}_i \|\nabla f_i(w^*)\|^2 \right)},$$

$$\begin{aligned} A_2 &= \sqrt[3]{\frac{L\Delta_0^2}{N} \sum_{i=1}^N \frac{(K_i - 1)}{q_i} (\sigma_i^2 + \|\nabla f_i(w^*)\|^2)}, \\ A_3 &= \Delta_0 \max \left\{ \frac{64L}{q_{\min}}, 10\bar{K}L \right\}. \end{aligned}$$

2) For f μ -strongly convex, $R \geq \max\left\{\frac{64L}{q_{\min}\mu}, \frac{10L\bar{K}}{\mu}\right\}$, there exists a sequence of weights $\{v_r\}_{r \geq 0}$, such that there holds:

$$\mathbb{E}(f(\bar{w}_R) - f^*) = \tilde{\mathcal{O}}\left(\mu\Delta_0 \exp(-\beta R) + \frac{B_1}{NR} + \frac{B_2}{R^2}\right),$$

where $\bar{w}_R = \frac{1}{V_R} \sum_{r=0}^{R-1} v_r w_r$, $V_R = \sum_{r=0}^{R-1} v_r$, with the problem related constants being

$$\begin{aligned} \beta &= \min \left\{ \frac{q_{\min}\mu}{128L}, \frac{\mu}{20L\bar{K}} \right\}, \\ B_1 &= \frac{1}{\mu N} \sum_{i=1}^N \left(\frac{\sigma_i^2}{q_i K_i} + \bar{q}_i \|\nabla f_i(w^*)\|^2 \right), \\ B_2 &= \frac{L}{\mu^2 N} \sum_{i=1}^N \frac{(K_i - 1)}{q_i} (\sigma_i^2 + \|\nabla f_i(w^*)\|^2). \end{aligned}$$

Here, $\Delta_0 = \|w_0 - w^*\|^2$ is the distance of the initial model from the true minimizer.

Several comments on Theorem 1 are now in order. First, the Theorem provides convergence rate guarantees for FedAvg for a more general heterogeneous setting than provided in existing studies such as [4], [6], [12] that assume equal number of local updates K_i 's and uniform sampling probabilities q_i 's. Next, the rates achieved for both smooth convex, and smooth strongly convex cases are order-optimal, i.e., when considering the slowest-decaying terms with respect to the number of epochs R , they scale order-optimally, as $1/\sqrt{R}$ and $1/R$, respectively. Furthermore, the established rates exhibit a linear speedup arising from FL in the number of users N . In more detail, with smooth convex functions, the leading error term scales as $1/\sqrt{NR}$. This rate is N -times faster than a central processing SGD method that would process at each iteration one of the functions f_i 's sequentially. A similar conclusion also holds for the strongly convex case. Hence, we generalize the linear speedup effect established in [6] to heterogeneous settings. In addition, through constants $A_1 - A_3$ and B_1, B_2 , we explicitly quantify the heterogeneity effects on performance, in terms of the σ_i 's, K_i 's, q_i 's and $\|\nabla f_i(w^*)\|$.

B. Optimal offline sampling probabilities q_i

Note that the convergence rate established in Theorem 1 depends on the parameters q_i . In particular, the dominant terms associated with q_i are of the form $\sum_{i=1}^N \frac{c_i}{q_i}$, where $c_i = \|\nabla f_i(w^*)\|^2 + \frac{\sigma_i^2}{K_i}$. Therefore, we would like to minimize

²Note that the communication probability associated with $\|\nabla f_i(w^*)\|^2$ is $\bar{q}_i = \frac{1-q_i}{q_i} = \frac{1}{q_i} - 1$. In the context of optimizing the communication probabilities, only the value $\frac{1}{q_i}$ is relevant, hence the expression $\sum_{i=1}^N \frac{c_i}{q_i}$ is correct.

quantity $\sum_{i=1}^N \frac{c_i}{q_i}$ with respect to probabilities q_i 's. As it is expected, without constraints on the q_i 's, the latter is minimized for $q_i = 1$, $i = 1, \dots, N$. However, we impose two types of constraints on the q_i 's. First, we have that $q_i \leq k_i$, where $k_i \in (0, 1]$ characterizes the capacity (quality) of the channel from user i to the server. Second, we let $\sum_{i=1}^N q_i = S$, for some $S < N$. This means that, on average, the server receives model increments from S users. This constraint mimics the uniform sampling scenario where a limited, small number of exactly S users are sampled at each epoch, in order to avoid the communication bottleneck at the server. This leads to the following optimization problem formulation:

$$\begin{aligned} \min_{q_1, \dots, q_N} \quad & \sum_{i=1}^N \frac{c_i}{q_i} \\ \text{s.t.} \quad & 0 \leq q_i \leq k_i, i \in [N], \sum_{i=1}^N q_i \leq S \end{aligned} \quad (4)$$

It can be shown (see Appendix B) that the solution is given by $q_i^* = k_i$, $i \in \mathcal{M}$, and $q_i^* = (S - k(\mathcal{M})) \frac{\sqrt{c_i}}{\sum_{j \in \mathcal{M}} \sqrt{c_j}}$, $i \notin \mathcal{M}$, where \mathcal{M} is the set of $|\mathcal{M}| = m$ indices corresponding to the largest values of c_i , $k(\mathcal{M}) = \sum_{i \in \mathcal{M}} k_i$, and m is either the largest positive integer satisfying $S - k(\mathcal{M}) \geq \frac{\sum_{i \in \mathcal{M}} \sqrt{c_i}}{\sqrt{\tilde{c}_m}}$, or $m = 0$. Here, \tilde{c}_m represents the m -th largest value of c_i 's. We now comment on the q_i^* 's. We can see that, the larger the c_i , the larger q_i^* should be. Moreover, q_i^* is either greedily equal to k_i , or it is proportional to the value of $\sqrt{c_i}$. Consider quantity $c_i = \|\nabla f_i(w^*)\|^2 + \frac{\sigma_i^2}{K_i}$. We can see that, therefore, the server should sample less frequently the users with large K_i 's. This is intuitive, as the users that make more local work reduce the local variance more and hence can transmit to the server less frequently. On the other hand, a user with a larger local SGD variance should communicate more frequently. This may appear counter-intuitive, but it is explained as follows: users with larger variance need to be reflected at the server side with more updates in order to filter out the local noise they incur. Finally, users with larger $\|\nabla f_i(w^*)\|$ should communicate more frequently with the server. Intuitively, a user with a very small $\|\nabla f_i(w^*)\|$ may safely skip its transmission when w_t is close to w^* , because its contribution would not significantly change w_t .

C. Online adaptive communication protocol

Note that quantities c_i 's depend on parameters that may be difficult to evaluate, such as $\|\nabla f_i(w^*)\|$ and σ_i^2 . These quantities can be estimated beforehand, so that a sub-optimal offline sampling probability scheme is devised by replacing the exact quantities with their estimates. For example, σ_i^2 relates to the local mini-batch size used by user i , while $\|\nabla f_i(w^*)\|$ may be replaced with $\|\nabla f_i(w')\|$, where w' is a solution estimate available at the server. In alternative, an online adaptive user selection policy may be utilized. Reference [8] devises an online user sampling strategy assuming reliable communication links. When adapted to our framework and notation, this strategy works as follows. At epoch t , user i is selected with

probability q_i^t , where q_i^t is the solution of (4) with quantity c_i replaced with $\|g_t^i\|^2$, and k_i set to one for all i . (We recall that $g_t^i = \frac{1}{K_i} \sum_{j=0}^{K_i-1} g_{t,j}^i$.) We next give a rationale that relates the optimal offline sampling derived here with the online strategy in [8]. When w_t is close to w^* , we can approximate g_t^i as follows: $g_t^i \approx \frac{1}{K_i} \sum_{j=0}^{K_i-1} (\nabla f_i(w^*) + n_{i,j}^t) = \nabla f_i(w^*) + \frac{1}{K_i} \sum_{j=0}^{K_i-1} n_{i,j}^t$, where $n_{i,j}^t$ is the user i SGD noise at the relevant inner and outer iteration. Taking the squared norm and expectation, while using independence of SGD noises across inner iterations, we can see that quantity $\mathbb{E}[\|g_t^i\|^2]$ can be approximated with c_i . In other words, when close to the solution, our offline sampling strategy is approximately the same as the adaptive strategy in [8]. The rest of the paper is devoted to a real-world implementation of the adaptive strategy in [8] on VCC tasks.

IV. EVALUATION RESULTS

We now present implementation and evaluation studies of heterogeneous user sampling on real-world deployments for VCC. The main purpose of the section is two fold. First, we show that employing FL in VCC can improve model performance (accuracy) over the VCC models working in isolation. Second, we show that non-uniform user sampling can significantly improve FL communication efficiency.

We consider three training schemes; 1) each user (pilot site) trains the VCC model in isolation on its local data (this scheme is termed *local training*); 2) FL algorithm (2)–(3) with full participation, i.e., all users transmit to the server at all epochs (this scheme is termed here *FedAvg*); and 3) FL algorithm (2)–(3) with the optimal adaptive sampling strategy [8] (this scheme is termed here *NUS – Non-Uniform Sampling*). We implement each of the three training schemes using the *Flower* Federated Learning framework³ for the Python programming language. With all schemes, all users receive the same initial model from the server which has its weights initialized to random values.

Each of the three pilot sites, namely, Valetta, Trento and Novi Sad, possess a local VCC input-output training dataset. For MT 170 videos and for GRN 660 videos were collected using static cameras. The corresponding data frames have been anonymized and subsequently annotated using the CVAT tool⁴. For the Novi Sad site, a dataset was collected within the staged recording that was carried out at the Petrovaradin fortress. No anonymisation was necessary as all human participants signed consent for the recording; we also note that in this type of experimental setup, identification of individuals is generally difficult due to the higher and overhead camera position. In total, around 800 annotated frames are provided for training. The model and loss functions adopted are as described in Section II-B. For the NUS policy, we set $S = 1$ (expected number of user transmissions per epoch). For both FedAvg and NUS, the total number of epochs is set to 100.

³<https://flower.dev/>

⁴<https://cvat.org/>

Table I shows the mean absolute error (MAE) at the three pilot sites at the end of training, for each of the three schemes. We can see that NUS improves performance over both local training and FedAvg for Novi Sad and Trento. On the other hand, the Valetta site incurs performance degradation. While monitoring the user selection process, we noticed that NUS rarely selects the Valetta user. Also, the distribution of the data samples in Trento and Novi Sad are quite similar, while for Valetta the corresponding distribution is different. For this reason, the Valetta site becomes “starved” of FL updates. In practice, the degradation due to FL can be easily fixed, by comparing the validation accuracy for both the joint FL model and the locally trained model and selecting the better performing model.

Figure 2 compares communication costs of FedAvg and NUS. We see the NUS cost reduction in model parameter transfers of around 17 gigabytes (17183MB). Therefore, NUS significantly reduces communication cost while at the same time improving MAE for the majority of sites. See also Figure 1 for evaluation of various losses versus execution time across the three sites for FedAvg and NUS.

TABLE I
FINAL VALIDATION MAE FOR LOCAL TRAINING, FEDAVG, AND NUS.

	NUS	FedAvg	Local training
Novi Sad	2.0890	3.1086	2.3622
Trento	27.6216	30.0928	40.3407
Valetta	9.1261	0.6900	0.6269

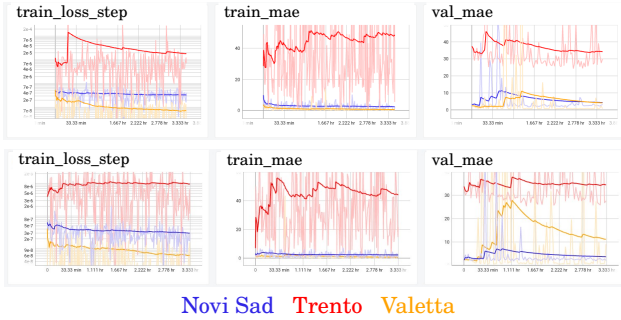


Fig. 1. Loss versus execution time. Top row: FedAvg; bottom: NUS. We monitor train_loss_step (training loss), train_mae (Training set MAE value) and val_mae (Validation set MAE value)

V. CONCLUSION

We considered a general framework for FL with various degrees of users heterogeneity, namely in terms of local gradient variances, activation probabilities, local number of updates, and (unreliable) link qualities. For this setting, we derived convergence rates for the FedAvg method, as well as optimal user activation probabilities that are pre-computed offline. We then connected these probabilities with previously derived optimal probabilities for an online setting [8]. A real-world data FL system for visual crowd counting across three

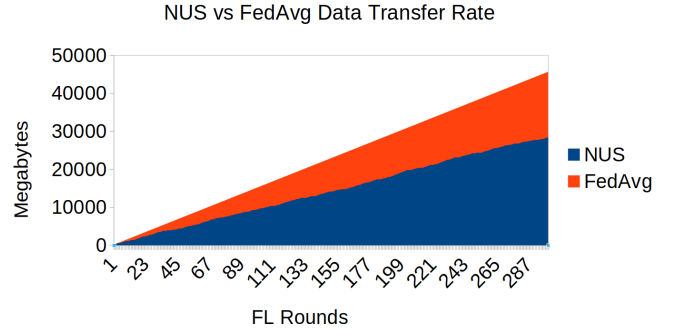


Fig. 2. Data transfers cost for FedAvg and NUS.

sites (Valetta, Trento, Novi Sad) was then deployed and evaluated, featuring optimized user activation probabilities. The evaluation shows significant accuracy benefits of federation at poorly performing local sites, and significant communication savings due to the optimized, non-uniform sampling employed.

REFERENCES

- [1] P. K. et al., “Advances and open problems in federated learning,” *Foundations and Trends in Machine Learning*, vol. 14, no. 1–2, 2021.
- [2] A. Armacki, D. Bajovic, D. Jakovetic, and S. Kar, “Personalized federated learning via convex clustering,” in *2022 IEEE International Smart Cities Conference (ISC2)*, Pafos, Cyprus, 2022.
- [3] X. Li, K. Huang, W. Yang, S. Wang, and Z. Zhang, “On the convergence of fedavg on non-iid data,” in *ICLR*, 2020.
- [4] S. P. Karimireddy, S. Kale, M. Mohri, S. Reddi, S. Stich, and A. T. Suresh, “SCAFFOLD: Stochastic controlled averaging for federated learning,” in *Proceedings of the 37th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, vol. 119. PMLR, 13–18 Jul 2020, pp. 5132–5143.
- [5] H. Yang, M. Fang, and J. Liu, “Achieving linear speedup with partial worker participation in non-iid federated learning,” *ICLR*, 2021.
- [6] Z. Qu, K. Lin, Z. Li, and J. Zhou, “Federated learning’s blessing: Fedavg has linear speedup,” in *ICLR 2021- DPML Workshop*, 2021.
- [7] J. Wang, Q. Liu, H. Liang, G. Joshi, and H. V. Poor, “Tackling the objective inconsistency problem in heterogeneous federated optimization,” in *Advances in Neural Information Processing Systems*, vol. 33. Curran Associates, Inc., 2020, pp. 7611–7623.
- [8] W. Chen, S. Horváth, and P. Richtárik, “Optimal client sampling for federated learning,” *Transactions on Machine Learning Research*, 2022.
- [9] Y. Fraboni, R. Vidal, L. Kamení, and M. Lorenzi, “Clustered sampling: Low-variance and improved representativity for clients selection in federated learning,” in *38th International Conference on Machine Learning*, vol. 139. PMLR, 18–24 Jul 2021, pp. 3407–3416.
- [10] Y. Jee Cho, J. Wang, and G. Joshi, “Towards understanding biased client selection in federated learning,” in *25th International Conference on Artificial Intelligence and Statistics*, vol. 151. PMLR, Mar 2022, pp. 10 351–10 375.
- [11] Q. Song, C. Wang, Y. Wang, Y. Tai, C. Wang, J. Li, J. Wu, and J. Ma, “To choose or to fuse? scale selection for crowd counting,” *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, no. 3, pp. 2576–2583, May 2021.
- [12] X. Gu, K. Huang, J. Zhang, and L. Huang, “Fast federated learning in the presence of arbitrary device unavailability,” in *Advances in Neural Information Processing Systems*, vol. 34. Curran Associates, Inc., 2021, pp. 12 052–12 064.
- [13] Y. Nesterov, *Lectures on Convex Optimization*, 2nd ed. Springer Publishing Company, Incorporated, 2018.

APPENDIX A CONVERGENCE ANALYSIS

In order to solve (1), the baseline algorithm works as follows: at the start of an epoch, the data center sends the

current parameter estimate w_t to the agents. Each agent, performs K_i local updates as follows

$$w_{t,j+1}^i = w_{t,j}^i - \alpha \nabla f_i(w_{t,j}^i, \xi_{t,j}^i), \quad j = 0, \dots, K_i - 1, \quad (5)$$

where $w_{t,0}^i = w_t$ and $\xi_{t,j}^i$ is the stochastic noise at agent i , epoch t , and local round j . Note that each agent is allowed to perform a different number of local steps, for example due to different number of available samples, or local computational budget. The master then averages the updates, i.e. performs

$$w_{t+1} = \sum_{i=1}^N \frac{p_i}{K_i} w_{t,K}^i. \quad (6)$$

Unrolling (5) in (6), we get

$$\begin{aligned} w_{t+1} &= \sum_{i=1}^N \frac{p_i}{K_i} \left(w_t - \alpha \sum_{j=0}^{K_i-1} \nabla f_i(w_{t,j}^i, \xi_{t,j}^i) \right) \\ &= w_t - \alpha \sum_{i=1}^N \frac{p_i}{K_i} \sum_{j=0}^{K_i-1} \nabla f_i(w_{t,j}^i, \xi_{t,j}^i). \end{aligned} \quad (7)$$

For the ease of notation, we define

$$\begin{aligned} g_{t,j}^i &= \nabla f_i(w_{t,j}^i, \xi_{t,j}^i), \\ g_t^i &= \frac{1}{K_i} \sum_{j=0}^{K_i-1} g_{t,j}^i, \end{aligned}$$

so that the master update (7) can be rewritten as

$$w_{t+1} = w_t - \alpha \sum_{i=1}^N p_i g_t^i. \quad (8)$$

We assume arbitrary communication likelihood on the side of the agents. To model this phenomena, we introduce Bernoulli random variables $X_t^i \sim \text{Bernoulli}(q_i)$, with parameter q_i . In this context, q_i represents the probability of agent i communicating their update to the master. Using the update rule (8), we then get

$$w_{t+1} = w_t - \alpha \sum_{i=1}^N X_{t+1}^i p_i g_t^i.$$

Taking the expectation with respect to the communication likelihood, \mathcal{C}_t , we get

$$\mathbb{E}_{\mathcal{C}_t}[w_{t+1}] = w_t - \alpha \sum_{i=1}^N q_i p_i g_t^i,$$

which is a biased version of the original update (8). Hence, we use the following debiased update rule

$$w_{t+1} = w_t - \alpha \sum_{i=1}^N p_i \frac{X_{t+1}^i}{q_i} g_t^i. \quad (9)$$

Note that our model subsumes the different agent sampling schemes. For example, in the case of uniform sampling, setting $q_i = \frac{S}{N}$, with $p_i = \frac{1}{N}$, using (9), we get

$$w_{t+1} = w_t - \alpha \sum_{i=1}^N \frac{1}{N} \frac{N}{S} X_{t+1}^i g_t^i = w_t - \frac{\alpha}{S} \sum_{l=1}^S g_t^{i_l} = \frac{1}{S} \sum_{l=1}^S w_t^{i_l}.$$

In this section, we analyze the convergence of the proposed algorithm. For the sake of simplicity, in what follows, we will assume that the cost function is of the form

$$f(w) = \frac{1}{N} \sum_{i=1}^N f_i(w). \quad (10)$$

Note that the original cost function can be restated in this form, by defining

$$f(w) = \frac{1}{N} \sum_{i=1}^N \tilde{f}_i(w),$$

where $\tilde{f}_i(w) := N p_i f_i(w)$.

The following quantities appear throughout the proofs and in our convergence results: the maximal number of local updates, $\bar{K} := \max_{i \in [N]} K_i$; the local communication probability parameters, $\bar{q}_i := \frac{1-q_i}{q_i}$, $q_{\min} := \min_{i \in [N]} q_i$. Next, we introduce some technical results used in our proofs.

Lemma 1. For any $a, b \in \mathbb{R}^d$, and any $\tau > 0$, we have

$$\begin{aligned} \|a + b\|^2 &\leq (1 + \tau) \|a\|^2 + \left(1 + \frac{1}{\tau}\right) \|b\|^2, \\ \left\| \sum_{i=1}^N a_i \right\|^2 &\leq N \sum_{i=1}^N \|a_i\|^2. \end{aligned}$$

The following lemma is a simple consequence of Theorem 2.1.5 from [13].

Lemma 2. Let $f_i : \mathbb{R}^d \mapsto \mathbb{R}$, $i \in [N]$ be convex and L -smooth. Define $f(x) := \frac{1}{N} \sum_{i=1}^N f_i(x)$, and let x^* be a minimizer of f . Then, for any $x \in \mathbb{R}^d$, we have:

$$\frac{1}{N} \sum_{i=1}^N \|\nabla f_i(x) - \nabla f_i(x^*)\|^2 \leq 2L(f(x) - f^*)$$

The following are restatements of Lemmas 1, 2, 4 and 5 from [4], respectively.

Lemma 3. For any non-negative sequence $\{d_r\}_{r \geq 0}$ and any parameters $\eta_{\max} > 0$, $c \geq 0$, $R \geq 0$, there exists a constant step-size $\eta \leq \eta_{\max}$, and weights $w_r = 1$, such that

$$\begin{aligned} \Psi_R &:= \frac{1}{R+1} \sum_{r=1}^{R+1} \left(\frac{d_{r-1}}{\eta} - \frac{d_r}{\eta} + c_1 \eta + c_2 \eta^2 \right) \\ &\leq \frac{d_0}{\eta_{\max}(R+1)} + \frac{2\sqrt{c_1 d_0}}{\sqrt{R+1}} + 2 \left(\frac{d_0}{R+1} \right)^{\frac{2}{3}} c_2^{\frac{1}{3}}. \end{aligned}$$

Lemma 4. For any non-negative sequence $\{d_r\}_{r \geq 0}$ and any parameters $\mu > 0$, $\eta_{\max} \in (0, \frac{1}{\mu}]$, $c \geq 0$, $R \geq \frac{1}{2\eta_{\max}\mu}$, there exists a constant step-size $\eta \leq \eta_{\max}$, and weights $w_r = (1 - \eta\mu)^{1-r}$, such that, for $W_R = \sum_{r=1}^{R+1} w_r$

$$\begin{aligned} \Psi_R &:= \frac{1}{W_R} \sum_{r=1}^{R+1} \left(\frac{w_r}{\eta} (1 - \mu\eta) d_{r-1} - \frac{w_r}{\eta} d_r + c\eta w_r \right) \\ &= \tilde{\mathcal{O}} \left(\mu d_0 \exp(-\mu\eta_{\max} R) + \frac{c}{\mu R} \right), \end{aligned}$$

where $\tilde{O}(\cdot)$ is used to hide poly-logarithmic terms and global, problem independent constants.

Lemma 5. Let f be μ -strongly convex, L -smooth. Then, for any $x, y, z \in \mathbb{R}^d$, we have:

$$\langle \nabla f(x), z - y \rangle \geq f(z) - f(y) + \frac{\mu}{4} \|z - y\|^2 - L \|z - x\|^2.$$

Lemma 6. Let X_1, \dots, X_n be n random variables in \mathbb{R}^n , not necessarily independent. Suppose their means are $\mathbb{E}[X_i] = \xi_i$, and the variance is bounded $\mathbb{E}\|X_i - \xi_i\|^2 \leq \sigma_i^2$. Then, the following holds:

$$\mathbb{E}\left\|\sum_{i=1}^n X_i\right\|^2 \leq \left\|\sum_{i=1}^n \xi_i\right\|^2 + n \sum_{i=1}^n \sigma_i^2.$$

Now suppose that $\mathbb{E}[X_i | X_1, \dots, X_{i-1}] = \xi_i$, i.e., the variables $\{X_i - \xi_i\}$ form a martingale difference sequence. Then, the following tighter bound holds:

$$\mathbb{E}\left\|\sum_{i=1}^n X_i\right\|^2 \leq 2\left\|\sum_{i=1}^n \xi_i\right\|^2 + 2\sum_{i=1}^n \sigma_i^2.$$

The following is a simple consequence of L -smoothness.

Lemma 7. Let f be L -smooth. Then, for any $x, y \in \mathbb{R}^d$, the following holds

$$\|\nabla f(x)\|^2 \leq 2L^2 \|x - y\|^2 + 2\|\nabla f(y)\|^2.$$

We will prove Theorem 1 by bounding two terms: the communication variance caused by agents' inability to communicate, and the optimality gap in the ideal scenario when everyone can send their updates. To this end, let v_{t+1} be the update in the ideal communication case, i.e.

$$v_{t+1} = w_t - \frac{\alpha}{N} \sum_{i=1}^N g_t^i. \quad (11)$$

Note that, by definition

$$\mathbb{E}_{C_t}[w_{t+1}] = v_{t+1}. \quad (12)$$

Next, we have

$$\begin{aligned} \mathbb{E}\|w_{t+1} - w^*\|^2 &= \mathbb{E}\|w_{t+1} - v_{t+1}\|^2 + \mathbb{E}\|v_{t+1} - w^*\|^2 \\ &\quad + 2\mathbb{E}\langle w_{t+1} - v_{t+1}, v_{t+1} - w^* \rangle. \end{aligned} \quad (13)$$

Using (12), we can conclude that the middle term on the right-hand side of (13) is zero, hence we can proceed to bound the remaining two terms.

A. Bounding the communication variance

In this section, we proceed to bound the first term, i.e., $\mathbb{E}\|w_{t+1} - v_{t+1}\|^2$. We will prove the following lemma.

Lemma 8. Let assumptions 1-3 hold. We then have the following bound on the communication variance

$$\begin{aligned} \mathbb{E}\|v_{t+1} - w_{t+1}\|^2 &\leq \frac{8\alpha^2 L^2}{N^2} \sum_{i,j} \frac{\bar{q}_i}{K_i} \mathbb{E}\|w_{t,j}^i - w_t\|^2 \\ &\quad + \frac{2\alpha^2}{N^2} \sum_{i=1}^N \frac{\bar{q}_i \sigma_i^2}{K_i} + \frac{4\alpha^2}{N^2} \sum_{i=1}^N \bar{q}_i \|\nabla f_i(w^*)\|^2 \\ &\quad + \frac{8\alpha^2}{N^2} \sum_{i=1}^N \bar{q}_i \mathbb{E}\|\nabla f_i(w_t) - \nabla f_i(w^*)\|^2. \end{aligned}$$

Proof. Using (9) and (11), we get

$$\begin{aligned} \mathbb{E}\|v_{t+1} - w_{t+1}\|^2 &= \alpha^2 \mathbb{E}\left\|\frac{1}{N} \sum_{i=1}^N \left(1 - \frac{X_{t+1}^i}{q_i}\right) g_t^i\right\|^2 \\ &= \alpha^2 \sum_{i=1}^N \mathbb{E}\left\|\frac{1}{N} \left(1 - \frac{X_{t+1}^i}{q_i}\right) g_t^i\right\|^2 \\ &= \frac{\alpha^2}{N^2} \sum_{i=1}^N \mathbb{E}\left(1 - \frac{X_{t+1}^i}{q_i}\right)^2 \mathbb{E}\|g_t^i\|^2 \\ &= \frac{\alpha^2}{N^2} \sum_{i=1}^N \frac{1 - q_i}{q_i} \mathbb{E}\|g_t^i\|^2, \end{aligned}$$

where the second equality follows from Assumption 3 and the fact that $\mathbb{E}_{C_t}\left(1 - \frac{X_t^i}{q_i}\right) = 0$. We proceed to bound $\mathbb{E}\|g_t^i\|^2$, as follows

$$\begin{aligned} \mathbb{E}\|g_t^i\|^2 &= \mathbb{E}\left\|\frac{1}{K_i} \sum_{j=0}^{K_i-1} g_{t,j}^i\right\|^2 \leq \frac{2\sigma_i^2}{K_i} + 2\mathbb{E}\left\|\frac{1}{K_i} \sum_{j=0}^{K_i-1} \nabla f_i(w_{t,j}^i)\right\|^2 \\ &\leq 4\mathbb{E}\left\|\frac{1}{K_i} \sum_{j=0}^{K_i-1} \nabla f_i(w_{t,j}^i) - \nabla f_i(w^*)\right\|^2 + \frac{2\sigma_i^2}{K_i} + 4\|\nabla f_i(w^*)\|^2 \\ &\leq \frac{4}{K_i} \sum_{j=0}^{K_i-1} \mathbb{E}\|\nabla f_i(w_{t,j}^i) - \nabla f_i(w^*)\|^2 + \frac{2\sigma_i^2}{K_i} + 4\|\nabla f_i(w^*)\|^2 \\ &\leq 4\|\nabla f_i(w^*)\|^2 + \frac{8L^2}{K_i} \sum_{j=0}^{K_i-1} \mathbb{E}\|w_{t,j}^i - w_t\|^2 + \frac{2\sigma_i^2}{K_i} \\ &\quad + 8\mathbb{E}\|\nabla f_i(w_t) - \nabla f_i(w^*)\|^2, \end{aligned} \quad (14)$$

where we used the second part of Lemma 6 in the first inequality. Plugging (14) back in the original expression completes the proof. \square

B. Bounding the optimality gap

In this section, we proceed to bound the first term, i.e., $\mathbb{E}\|v_{t+1} - w^*\|^2$. We will prove the following lemma.

Lemma 9. Let assumptions 1-3 hold. If the step-size satisfies $\alpha \leq \frac{1}{8L}$, we then have the following bound on the optimality gap

$$\begin{aligned} \mathbb{E}\|v_{t+1} - w^*\|^2 &\leq \left(1 - \frac{\mu\alpha}{2}\right) \mathbb{E}\|w_t - w^*\|^2 + \frac{2\alpha^2}{N^2} \sum_{i=1}^N \frac{\sigma_i^2}{K_i} \\ &\quad + \frac{4\alpha L}{N} \sum_{i,j} \frac{1}{K_i} \mathbb{E}\|w_{t,j}^i - w_t\|^2 - \alpha \mathbb{E}(f(w_t) - f(w^*)). \end{aligned}$$

Proof. Using (11), we get

$$\begin{aligned} \mathbb{E}\|v_{t+1} - w^*\|^2 &= \mathbb{E}\left\|w_t - w^* - \frac{\alpha}{N} \sum_{i=1}^N g_t^i\right\|^2 \\ &= \mathbb{E}\|w_t - w^*\|^2 + \alpha^2 \mathbb{E}\left\|\frac{1}{N} \sum_{i,j} \frac{1}{K_i} g_{t,j}^i\right\|^2 \\ &\quad - 2\alpha \mathbb{E}\left\langle w_t - w^*, \frac{1}{N} \sum_{i=1}^N g_t^i \right\rangle \\ &\leq \mathbb{E}\|w_t - w^*\|^2 - \frac{2\alpha}{N} \sum_{i,j} \frac{1}{K_i} \mathbb{E}\langle w_t - w^*, \nabla f_i(w_{t,j}^i) \rangle \\ &\quad + \frac{2\alpha^2}{N^2} \sum_{i=1}^N \frac{\sigma_i^2}{K_i} + 2\alpha^2 \mathbb{E}\left\|\frac{1}{N} \sum_{i,j} \frac{1}{K_i} \nabla f_i(w_{t,j}^i)\right\|^2, \end{aligned} \tag{15}$$

where we applied the second part of Lemma 6 in the last inequality. Next, we bound the term

$$\begin{aligned} \mathbb{E}\left\|\frac{1}{N} \sum_{i,j} \frac{1}{K_i} \nabla f_i(w_{t,j}^i)\right\|^2 &\leq 2\mathbb{E}\left\|\frac{1}{N} \sum_i \nabla f_i(w_t)\right\|^2 \\ &\quad + 2\mathbb{E}\left\|\frac{1}{N} \sum_{i,j} \frac{1}{K_i} (\nabla f_i(w_{t,j}^i) - \nabla f_i(w_t))\right\|^2 \\ &\leq \frac{2}{N} \sum_{i,j} \frac{1}{K_i} \mathbb{E}\|\nabla f_i(w_{t,j}^i) - \nabla f_i(w_t)\|^2 + 2\mathbb{E}\|\nabla f(w_t)\|^2 \\ &\leq \frac{2L^2}{N} \sum_{i,j} \frac{1}{K_i} \mathbb{E}\|w_{t,j}^i - w_t\|^2 + 4L\mathbb{E}(f(w_t) - f^*). \end{aligned}$$

Applying Lemma 5, and plugging in the expression above, we have

$$\begin{aligned} \mathbb{E}\|v_{t+1} - w^*\|^2 &\leq \left(1 - \frac{\mu\alpha}{2}\right) \mathbb{E}\|w_t - w^*\|^2 + \frac{2\alpha^2}{N^2} \sum_{i=1}^N \frac{\sigma_i^2}{K_i} \\ &\quad + \frac{2\alpha L(1 + 2\alpha L)}{N} \sum_{i,j} \frac{1}{K_i} \mathbb{E}\|w_{t,j}^i - w_t\|^2 \\ &\quad - 2\alpha(1 - 4\alpha L) \mathbb{E}(f(w_t) - f(w^*)). \end{aligned}$$

The claim then follows by choosing $\alpha \leq \frac{1}{8L}$. \square

C. Bounding the local drift

In this section, we bound the local drift, i.e., $\mathbb{E}\|w_{t,j}^i - w_t\|^2$. We use Lemma B.4 from [12], which is a slightly modified version of Lemma 8 from [4]. For completeness, we provide the full proof.

Lemma 10. Let assumptions 1-3 hold. If the step-size satisfies $\alpha \leq \frac{1}{10LK}$, we then have the following bound on the local drift

$$\begin{aligned} \mathbb{E}\|w_{t,j}^i - w_t\|^2 &\leq 8(K_i - 1)\alpha^2 \mathbb{E}\|\nabla f_i(w_t) - \nabla f_i(w^*)\|^2 \\ &\quad + 2(K_i - 1)\alpha^2 \sigma_i^2 + 8(K_i - 1)\alpha^2 \|\nabla f_i(w^*)\|^2, \end{aligned}$$

for any $0 \leq j \leq K_i - 1$.

Proof. For $j = 0$, we have

$$\|w_{t,j}^i - w_t\|^2 = 0,$$

hence the bound trivially holds. Using the local update rule (5), for any $j > 0$, we have

$$\begin{aligned} \mathbb{E}\|w_{t,j}^i - w_t\|^2 &= \mathbb{E}\|w_{t,j-1}^i - w_t - \alpha g_{t,j-1}^i\|^2 \\ &= \mathbb{E}\|w_{t,j-1}^i - w_t - \alpha \nabla f_i(w_{t,j-1}^i)\|^2 \\ &\quad + \alpha^2 \mathbb{E}\|g_{t,j-1}^i - \nabla f_i(w_{t,j-1}^i)\|^2 - 2\alpha \mathbb{E}\Gamma_{t,j}^i, \end{aligned}$$

where $\Gamma_{t,j}^i$ is given by

$$\Gamma_{t,j}^i = \langle w_{t,j-1}^i - w_t - \alpha \nabla f_i(w_{t,j-1}^i), g_{t,j-1}^i - \nabla f_i(w_{t,j-1}^i) \rangle.$$

Noticing that $\mathbb{E}\Gamma_{t,j}^i = 0$, applying the variance bound and Lemma 1 with $\tau = K_i - 1$, we get

$$\begin{aligned} \mathbb{E}\|w_{t,j}^i - w_t\|^2 &\leq \left(1 + \frac{1}{K_i - 1}\right) \mathbb{E}\|w_{t,j-1}^i - w_t\|^2 \\ &\quad + \alpha^2 K_i \mathbb{E}\|\nabla f_i(w_{t,j-1}^i)\|^2 + \alpha^2 \sigma_i^2 \\ &\leq \left(1 + \frac{1}{K_i - 1} + 2\alpha^2 L^2 K_i\right) \mathbb{E}\|w_{t,j-1}^i - w_t\|^2 \\ &\quad + 2\alpha^2 \mathbb{E}\|\nabla f_i(w_t)\|^2 + \alpha^2 \sigma_i^2. \end{aligned}$$

For $\alpha \leq \frac{1}{10LK}$, we have $2\alpha^2 L^2 K_i \leq \frac{1}{50(K_i - 1)}$. It then follows that

$$\begin{aligned} \mathbb{E}\|w_{t,j}^i - w_t\|^2 &\leq \left(1 + \frac{51}{50(K_i - 1)}\right) \mathbb{E}\|w_{t,j-1}^i - w_t\|^2 \\ &\quad + 2\alpha^2 \mathbb{E}\|\nabla f_i(w_t)\|^2 + \alpha^2 \sigma_i^2. \end{aligned}$$

Set $A_j = \mathbb{E}\|w_{t,j}^i - w_t\|^2$, $B = 2\alpha^2 \mathbb{E}\|\nabla f_i(w_t)\|^2 + \alpha^2 \sigma_i^2$ and $C = 1 + \frac{51}{50(K_i - 1)}$, to obtain the following recursion

$$A_j \leq CA_{j-1} + B.$$

Unrolling the expression, and noting that $A_0 = w_{t,0}^i - w_t = 0$, we get

$$A_j \leq B \sum_{l=0}^{j-1} C^l = B \frac{C^j - 1}{C - 1} \leq B \frac{C^{K_i-1} - 1}{C - 1}.$$

The expression C^{K_i-1} can be further simplified as

$$\begin{aligned} C^{K_i-1} &= \left(1 + \frac{51}{50(K_i - 1)}\right)^{K_i-1} \\ &= \left(1 + \frac{51}{50(K_i - 1)}\right)^{\frac{50(K_i-1)}{51} \frac{51}{50}} \leq e^{\frac{51}{50}} < 3. \end{aligned}$$

Hence, we get the following bound

$$\begin{aligned}\mathbb{E}\|w_{t,j}^i - w_t\|^2 &\leq 2(K_i - 1) \left(2\alpha^2 \mathbb{E}\|\nabla f_i(w_t)\|^2 + \alpha^2 \sigma_i^2 \right) \\ &\leq 2(K_i - 1) \alpha^2 \sigma_i^2 + 8(K_i - 1) \alpha^2 \|\nabla f_i(w^*)\|^2 \\ &\quad + 8(K_i - 1) \alpha^2 \mathbb{E}\|\nabla f_i(w_t) - \nabla f_i(w^*)\|^2,\end{aligned}$$

which completes the proof. \square

D. Completing the proof of Theorem 1

Proof. Note that so far we have

$$\mathbb{E}\|w_{t+1} - w^*\|^2 = \mathbb{E}\|w_{t+1} - v_{t+1}\|^2 + \mathbb{E}\|v_{t+1} - w^*\|^2.$$

Combining Lemmas 8 and 9, we get the following upper-bound

$$\begin{aligned}\mathbb{E}\|w_{t+1} - w^*\|^2 &\leq \left(1 - \frac{\mu\alpha}{2}\right) \mathbb{E}\|w_t - w^*\|^2 + \frac{2\alpha^2}{N^2} \sum_{i=1}^N \frac{\sigma_i^2}{K_i} \\ &\quad + \frac{4\alpha L}{N} \sum_{i,j} \frac{1}{K_i} \mathbb{E}\|w_{t,j}^i - w_t\|^2 - \alpha \mathbb{E}(f(w_t) - f(w^*)) \\ &\quad + \frac{8\alpha^2 L^2}{N^2} \sum_{i,j} \frac{\bar{q}_i}{K_i} \mathbb{E}\|w_{t,j}^i - w_t\|^2 + \frac{4\alpha^2}{N^2} \sum_{i=1}^N \bar{q}_i \|\nabla f_i(w^*)\|^2 \\ &\quad + \frac{2\alpha^2}{N^2} \sum_{i=1}^N \frac{\bar{q}_i \sigma_i^2}{K_i} + \frac{8\alpha^2}{N^2} \sum_{i=1}^N \bar{q}_i \mathbb{E}\|\nabla f_i(w_t) - \nabla f_i(w^*)\|^2.\end{aligned}$$

For $\alpha \leq \frac{1}{2L}$, we get

$$\begin{aligned}\mathbb{E}\|w_{t+1} - w^*\|^2 &\leq \left(1 - \frac{\mu\alpha}{2}\right) \mathbb{E}\|w_t - w^*\|^2 + \frac{2\alpha^2}{N^2} \sum_{i=1}^N \frac{\sigma_i^2}{q_i K_i} \\ &\quad + \frac{4\alpha L}{N} \sum_{i,j} \frac{1}{q_i K_i} \mathbb{E}\|w_{t,j}^i - w_t\|^2 \\ &\quad - \alpha \mathbb{E}(f(w_t) - f(w^*)) + \frac{4\alpha^2}{N^2} \sum_{i=1}^N \bar{q}_i \|\nabla f_i(w^*)\|^2 \\ &\quad + \frac{8\alpha^2}{N^2} \sum_{i=1}^N \bar{q}_i \mathbb{E}\|\nabla f_i(w_t) - \nabla f_i(w^*)\|^2.\end{aligned}$$

Applying Lemma 10, we have

$$\begin{aligned}\mathbb{E}\|w_{t+1} - w^*\|^2 &\leq \left(1 - \frac{\mu\alpha}{2}\right) \mathbb{E}\|w_t - w^*\|^2 + \frac{2\alpha^2}{N^2} \sum_{i=1}^N \frac{\sigma_i^2}{q_i K_i} \\ &\quad - \alpha \mathbb{E}(f(w_t) - f(w^*)) + \frac{4\alpha^2}{N^2} \sum_{i=1}^N \bar{q}_i \|\nabla f_i(w^*)\|^2 \\ &\quad + \frac{8\alpha^3 L}{N} \sum_{i=1}^N \frac{(K_i - 1) \sigma_i^2}{q_i} + \frac{32\alpha^3 L}{N} \sum_{i=1}^N \frac{(K_i - 1) \|\nabla f_i(w^*)\|^2}{q_i} \\ &\quad + \frac{32\alpha^3 L}{N} \sum_{i=1}^N \frac{(K_i - 1) \mathbb{E}\|\nabla f_i(w_t) - \nabla f_i(w^*)\|^2}{q_i} \\ &\quad + \frac{8\alpha^2}{N^2} \sum_{i=1}^N \mathbb{E} \bar{q}_i \|\nabla f_i(w_t) - \nabla f_i(w^*)\|^2.\end{aligned}$$

Using the fact that $1 - \bar{q}_i \leq \frac{1}{q_i}$ and choosing $\alpha \leq \frac{1}{4KL}$, we get

$$\begin{aligned}\mathbb{E}\|w_{t+1} - w^*\|^2 &\leq \left(1 - \frac{\mu\alpha}{2}\right) \mathbb{E}\|w_t - w^*\|^2 + \frac{2\alpha^2}{N^2} \sum_{i=1}^N \frac{\sigma_i^2}{q_i K_i} \\ &\quad - \alpha \mathbb{E}(f(w_t) - f(w^*)) + \frac{4\alpha^2}{N^2} \sum_{i=1}^N \bar{q}_i \|\nabla f_i(w^*)\|^2 \\ &\quad + \frac{8\alpha^3 L}{N} \sum_{i=1}^N \frac{(K_i - 1) \sigma_i^2}{q_i} + \frac{32\alpha^3 L}{N} \sum_{i=1}^N \frac{(K_i - 1) \|\nabla f_i(w^*)\|^2}{q_i} \\ &\quad + \frac{16\alpha^2}{q_{\min} N} \sum_{i=1}^N \mathbb{E}\|\nabla f_i(w_t) - \nabla f_i(w^*)\|^2.\end{aligned}$$

Applying Lemma 2, we get

$$\begin{aligned}\mathbb{E}\|w_{t+1} - w^*\|^2 &\leq \left(1 - \frac{\mu\alpha}{2}\right) \mathbb{E}\|w_t - w^*\|^2 + \frac{2\alpha^2}{N^2} \sum_{i=1}^N \frac{\sigma_i^2}{q_i K_i} \\ &\quad - \alpha \left(1 - \frac{32\alpha L}{q_{\min}}\right) \mathbb{E}(f(w_t) - f(w^*)) + \frac{4\alpha^2}{N^2} \sum_{i=1}^N \bar{q}_i \|\nabla f_i(w^*)\|^2 \\ &\quad + \frac{8\alpha^3 L}{N} \sum_{i=1}^N \frac{(K_i - 1) \sigma_i^2}{q_i} + \frac{32\alpha^3 L}{N} \sum_{i=1}^N \frac{(K_i - 1) \|\nabla f_i(w^*)\|^2}{q_i}.\end{aligned}$$

For $\alpha \leq \frac{q_{\min}}{64L}$, we get

$$\begin{aligned}\mathbb{E}\|w_{t+1} - w^*\|^2 &\leq \left(1 - \frac{\mu\alpha}{2}\right) \mathbb{E}\|w_t - w^*\|^2 + \frac{2\alpha^2}{N^2} \sum_{i=1}^N \frac{\sigma_i^2}{q_i K_i} \\ &\quad - \frac{\alpha}{2} \mathbb{E}(f(w_t) - f(w^*)) + \frac{4\alpha^2}{N^2} \sum_{i=1}^N \bar{q}_i \|\nabla f_i(w^*)\|^2 \\ &\quad + \frac{8\alpha^3 L}{N} \sum_{i=1}^N \frac{(K_i - 1) \sigma_i^2}{q_i} + \frac{32\alpha^3 L}{N} \sum_{i=1}^N \frac{(K_i - 1) \|\nabla f_i(w^*)\|^2}{q_i}.\end{aligned}$$

Rearranging and defining $\Delta_t := \mathbb{E}\|w_t - w^*\|^2$, gives

$$\begin{aligned}\mathbb{E}(f(w_t) - f^*) &\leq 2\alpha^{-1} \left[\left(1 - \frac{\mu\alpha}{2}\right) \mathbb{E}\Delta_t - \mathbb{E}\Delta_{t+1} \right] \\ &\quad + \frac{8\alpha}{N^2} \sum_{i=1}^N \bar{q}_i \|\nabla f_i(w^*)\|^2 + \frac{16\alpha^2 L}{N} \sum_{i=1}^N \frac{(K_i - 1) \sigma_i^2}{q_i} \\ &\quad + \frac{4\alpha}{N^2} \sum_{i=1}^N \frac{\sigma_i^2}{q_i K_i} + \frac{64\alpha^2 L}{N} \sum_{i=1}^N \frac{(K_i - 1) \|\nabla f_i(w^*)\|^2}{q_i}.\end{aligned}$$

If f is general convex, i.e., $\mu = 0$, we can directly apply Lemma 3. Otherwise, using weights $v_r = (1 - \frac{\mu\alpha}{2})^{1-r}$, for some $\frac{1}{\mu R} \leq \alpha \leq \min\left\{\frac{1}{10KL}, \frac{q_{\min}}{64L}\right\}$, we can apply Lemma 4 to obtain the upper bound

$$\mathbb{E}(f(\bar{w}_R) - f^*) = \mathcal{O}\left(\mu \Delta_0 \exp(-\beta R) + \frac{B_1}{NR} + \frac{B_2}{R^2}\right),$$

where β, B_1, B_2 are given in the statement of Theorem 1. \square

APPENDIX B
OPTIMAL COMMUNICATION PROBABILITY

The convergence rates established in Theorem 1 depend on the parameters q_{\min} , $\frac{1}{q_i}$ and \bar{q}_i . Naturally, the terms mostly affected by the communication infrequency are the lower-order terms, namely, the terms of order $\frac{1}{\sqrt{R}}$ and $\frac{1}{R}$, for convex and strongly convex functions, respectively. Therefore, to mitigate the effect of imperfect communication, we can optimize with respect to the communication probabilities associated with lower-order terms. In particular, those terms are of the form

$$\sum_{i=1}^N \frac{c_i}{q_i}, \quad (16)$$

where $c_i := \|\nabla f_i(w^*)\|^2 + \frac{\sigma_i^2}{K_i}$. In order to obtain the best possible rate, we would like to minimize (16), by solving the following problem:

$$\begin{aligned} \min_{q_1, \dots, q_N} \quad & \sum_{i=1}^N \frac{c_i}{q_i} \\ \text{s.t.} \quad & 0 \leq q_i \leq 1, \forall i \in [N]. \end{aligned} \quad (17)$$

The problem defined above has a trivial solution, namely $q_i = 1$. Hence, we impose some additional constraints, like limiting the maximum expected number of agents participating, by defining the following problem

$$\begin{aligned} \min_{q_1, \dots, q_N} \quad & \sum_{i=1}^N \frac{1}{q_i} \left(\|\nabla f_i(w^*)\|^2 + \frac{\sigma_i^2}{K_i} \right) \\ \text{s.t.} \quad & 0 \leq q_i \leq 1, \forall i \in [N], \\ & \sum_{i=1}^N q_i = S. \end{aligned} \quad (18)$$

We obtain the optimal probabilities in the following lemma.

Lemma 11. *The optimal probabilities obtained by solving (18) are given by*

$$q_i^* = \begin{cases} 1, & i \in \mathcal{M} \\ (S - m) \frac{\sqrt{c_i}}{\sum_{j \notin \mathcal{M}} \sqrt{c_j}}, & i \notin \mathcal{M} \end{cases}, \quad (19)$$

where \mathcal{M} is the set of $m := |\mathcal{M}|$ indices corresponding to the largest values of c_i , and $0 \leq m < S$ is either the largest positive integer satisfying

$$S - m \geq \frac{\sum_{i \notin \mathcal{M}} \sqrt{c_i}}{\sqrt{\tilde{c}_m}}, \quad (20)$$

or $m = 0$. Here, \tilde{c}_m represents the m -th largest value of c_i 's.

Remark 1. *Note that, in the case $S = N$, the solution of (18) is trivial, i.e. $q_i = 1$. From the constraint $m < S = N$, one might think that the solution is not captured by (19). However, that is not the case. For the sake of simplicity, assume c_i 's are sorted in a non-decreasing order. Plugging $m = N - 1$ into (20), we have*

$$1 \geq \sqrt{\frac{c_N}{c_{N-1}}},$$

which is true, as the c_i 's are non-decreasing. Hence, we have

$$q_i = 1, \quad i = 1, \dots, N - 1,$$

and

$$q_N = (N - (N - 1)) \sqrt{\frac{c_N}{c_N}} = 1.$$

Remark 2. *Another interesting case happens when $c_i = c$, $\forall i \in [N]$. In that case, the problem (18) becomes*

$$\begin{aligned} \min_{q_1, \dots, q_N} \quad & \sum_{i=1}^N \frac{c}{q_i} \\ \text{s.t.} \quad & 0 \leq q_i \leq 1, \forall i \in [N], \\ & \sum_{i=1}^N q_i = S. \end{aligned}$$

Using the substitution $x_i = \frac{c}{q_i}$ and the arithmetic-harmonic mean inequality, we get the following lower-bound on the solution

$$\sum_{i=1}^N x_i \geq \frac{N^2}{\sum_{i=1}^N \frac{1}{x_i}} = \frac{cN^2}{\sum_{i=1}^N q_i} = c \frac{N^2}{S}. \quad (21)$$

Noting that, for any $m > 0$

$$\frac{\sum_{i=m+1}^N \sqrt{c_i}}{\sqrt{c_m}} = \frac{\sqrt{c}(N - m)}{\sqrt{c}} = N - m \geq S - m,$$

we can conclude that $m = 0$ (the equality only holds for the case $S = N$ and that was covered in the previous remark) and hence we obtain the following probabilities

$$q_i = S \frac{\sqrt{c_i}}{\sum_{i=1}^N \sqrt{c_i}} = S \frac{\sqrt{c}}{N\sqrt{c}} = \frac{S}{N}, \quad \forall i \in [N].$$

Plugging them into the cost function, we get

$$\sum_{i=1}^N \frac{c}{q_i} = c \frac{N^2}{S},$$

which shows the lower bound established in (21) is attained by our solution.

Proof. Defining the cost

$$f(q) = \sum_{i=1}^N \frac{c_i}{q_i},$$

we can reformulate the problem (18) as

$$\begin{aligned} \min_q \quad & \sum_{i=1}^N f(q) \\ \text{s.t.} \quad & g_j(q) \leq 0, \forall j \in [2N], \\ & h(q) = 0, \end{aligned}$$

where

$$\begin{aligned} g_i(q) &= q_i - 1, \quad i \in \{1, \dots, N\}, \\ g_j(q) &= -q_{j-N}, \quad i \in \{N + 1, \dots, 2N\}, \\ h(q) &= \sum_{i=1}^N q_i - S. \end{aligned}$$

Using the KKT conditions, we get the following expressions:

1) **Stationarity:**

$$-\frac{c_i}{q_i^{*2}} + \mu_i - \mu_{i+N} + \lambda = 0, \quad \forall i \in [N]. \quad (22)$$

2) **Primal Feasibility:**

$$\begin{aligned} 0 &\leq q_i^* \leq 1, \\ \sum_{i=1}^N q_i^* &= S. \end{aligned} \quad (23)$$

3) **Dual Feasibility:**

$$\mu_i \geq 0, \quad \forall i \in [2N].$$

4) **Complementary slackness:**

$$\begin{aligned} \mu_i(q_i^* - 1) &= 0, \\ q_i^* \mu_{i+N} &= 0, \end{aligned} \quad i \in \{1, \dots, N\} \quad (24)$$

From the formulation of the problem, it is clear that $q_i > 0$, hence, from (24) we get

$$\mu_{i+N} = 0, \quad \forall i \in [N].$$

From (22), we get

$$q_i^* = \sqrt{\frac{c_i}{\mu_i + \lambda}}. \quad (25)$$

We now differentiate between the following two cases.

Case 1: for all $i \in [N]$, $q_i < 1$. In this case, from (24), we get

$$\mu_i = 0, \quad \forall i \in [N],$$

and hence, from (23) and (25), we get

$$S = \sum_{i=1}^N q_i^* = \sum_{i=1}^N \sqrt{\frac{c_i}{\lambda}} \implies \sqrt{\lambda} = \frac{1}{S} \sum_{i=1}^N \sqrt{c_i},$$

which in turns gives

$$q_i^* = S \frac{\sqrt{c_i}}{\sum_{i=1}^N \sqrt{c_i}}.$$

Case 2: for some $i \in [N]$, $q_i = 1$. By the formulation of the problem, we know that in this case, the $q_i = 1$ will correspond to the largest values of c_i . Without the loss of generality, assume that the c_i 's are sorted in a non-increasing order, i.e.

$$c_1 \geq c_2 \geq \dots \geq c_N.$$

Let m represent the number of agents with the corresponding probabilities $q_i = 1$, i.e. $q_i = 1, i = 1, \dots, m$. We then know, by complementary slackness, that for the remaining agents

$$\mu_i = 0, \quad i = m + 1, \dots, N.$$

Combining (25), primal feasibility condition and the equation above, we get

$$\begin{aligned} S &= \sum_{i=1}^N q_i^* = m + \sum_{i=m+1}^N \sqrt{\frac{c_i}{\lambda}} \\ \implies \sqrt{\lambda} &= \frac{1}{S-m} \sum_{i=m+1}^N \sqrt{c_i}, \end{aligned} \quad (26)$$

and hence

$$q_i^* = (S-m) \frac{\sqrt{c_i}}{\sum_{j=m+1}^N \sqrt{c_j}}, \quad i = m+1, \dots, N.$$

Combining, we get

$$q_i^* = \begin{cases} 1, & i = 1, \dots, m, \\ (S-m) \frac{\sqrt{c_i}}{\sum_{j=m+1}^N \sqrt{c_j}}, & i = m+1, \dots, N. \end{cases}$$

Note that $S-m > 0$, by the primal feasibility conditions. Additionally, using the closed form solution (25), we get

$$1 = q_i^* = \sqrt{\frac{c_i}{\mu_i + \lambda}}, \quad i = 1, \dots, m,$$

and hence

$$\lambda = c_i - \mu_i \leq c_i, \quad \forall i = 1, \dots, m,$$

as $\mu_i \geq 0, i = 1, \dots, m$. By the assumption that c_i 's are sorted in a non-decreasing order, we can equivalently state the equation above as

$$\lambda \leq c_m.$$

Applying (26), we get

$$\begin{aligned} \frac{1}{(S-m)^2} \left(\sum_{i=m+1}^N \sqrt{c_i} \right)^2 &\leq c_m \\ \implies S-m &\geq \frac{\sum_{i=m+1}^N \sqrt{c_i}}{\sqrt{c_m}}. \end{aligned}$$

Unifying both cases, we have

$$q_i^* = \begin{cases} 1, & i = 1, \dots, m, \\ (S-m) \frac{\sqrt{c_i}}{\sum_{j=m+1}^N \sqrt{c_j}}, & i = m+1, \dots, N, \end{cases}$$

where $m < S$ is either the largest positive integer satisfying

$$S-m \geq \frac{\sum_{i=m+1}^N \sqrt{c_i}}{\sqrt{c_m}},$$

or $m = 0$. \square

A. Optimal probabilities with channel constraints

Let k_i quantify the reliability of agent i 's communication channel, with $0 < k_i \leq 1$. We then modify our problem (18) to account for this information, as follows

$$\begin{aligned} \min_{q_1, \dots, q_N} \quad & \sum_{i=1}^N \frac{c_i}{q_i} \\ \text{s.t.} \quad & 0 \leq q_i \leq k_i, \quad \forall i \in [N], \\ & \sum_{i=1}^N q_i \leq S \end{aligned}$$

where modify the final constraint to an inequality one, in order to capture the case where $\sum_{i=1}^N k_i < S$. Repeating the steps

from lemma 11, it can be shown that the optimal probabilities are given by

$$q_i^* = \begin{cases} k_i, & i \in \mathcal{M} \\ (S - k(\mathcal{M})) \frac{\sqrt{c_i}}{\sum_{j \notin \mathcal{M}} \sqrt{c_j}}, & i \notin \mathcal{M} \end{cases},$$

where \mathcal{M} is the set of $m := |\mathcal{M}|$ indices corresponding to the largest values of c_i , $k(\mathcal{M}) = \sum_{i \in \mathcal{M}} k_i$, and m is either the largest positive integer satisfying

$$S - k(\mathcal{M}) \geq \frac{\sum_{i \notin \mathcal{M}} \sqrt{c_i}}{\sqrt{\tilde{c}_m}}, \quad (27)$$

or $m = 0$. Here, \tilde{c}_m represents the m -th largest value of c_i 's. Note that for the case where $\sum_{i=1}^N k_i \leq S$, the condition (27) evaluates to

$$S - \sum_{i=1}^N k_i \geq 0,$$

which is true by default. Hence, the optimal solution of $q_i^* = k_i$, $\forall i \in [N]$ is recovered.