# Stragglers-Aware Low-Latency Synchronous Federated Learning via Layer-Wise Model Updates

Natalie Lang, Alejandro Cohen, and Nir Shlezinger

Abstract-Synchronous federated learning (FL) is a popular paradigm for collaborative edge learning. It typically involves a set of heterogeneous devices locally training neural network (NN) models in parallel with periodic centralized aggregations. As some of the devices may have limited computational resources and varying availability, FL latency is highly sensitive to stragglers. Conventional approaches discard incomplete intra-model updates done by stragglers, alter the amount of local workload and architecture, or resort to asynchronous settings; which all affect the trained model performance under tight training latency constraints. In this work, we propose straggler-aware layerwise federated learning (SALF) that leverages the optimization procedure of NNs via backpropagation to update the global model in a layer-wise fashion. SALF allows stragglers to synchronously convey partial gradients, having each layer of the global model be updated independently with a different contributing set of users. We provide a theoretical analysis, establishing convergence guarantees for the global model under mild assumptions on the distribution of the participating devices, revealing that SALF converges at the same asymptotic rate as FL with no timing limitations. This insight is matched with empirical observations, demonstrating the performance gains of SALF compared to alternative mechanisms mitigating the device heterogeneity gap in FL.

### I. INTRODUCTION

**D** EEP learning algorithms require large volumes of data. In practice, data is often gathered by edge devices such as mobile phones, sensors, and vehicles, which may be limited in their ability to share this data, due to, e.g., privacy or regulation constraints [1]. *Federated learning (FL)* [2]–[5] is an emerging paradigm that allows multiple devices to learn a model collaboratively. To avoid data sharing, FL exploits the local computational capabilities of edge users [6], training locally with periodic aggregations orchestrated by a server.

Prevalent FL protocols are synchronous, operating in multiple rounds. In each round, the server transmits the latest update of the global model to all participating clients. Each client then locally trains the model using its local compute power and data set, and subsequently sends the updated version back to the server for aggregation. This distributed operation of FL inherently induces several core challenges that are not encountered in conventional centralized learning [4], [5]. A notable challenge is associated with the *latency* of learning in a federated manner, which is a dominant factor, particularly in FL applications that operate continuously in dynamic environments and must thus learn rapidly. Such applications include, e.g., intelligent transportation systems [7], and wireless networks adaptation [8]–[10].

The excessive latency of FL is mainly due to two main factors: (i) the time it takes to communicate the model updates between the users and the server; and (ii) the local computation time at the users' side. The former is typically addressed by different forms of sparsification [11]–[14] and compression [15]-[18]. Thus, the dominant factor is often the latter, i.e., the time it takes each user to update the local model. This issue is most significant in FL settings operating under system heterogeneity, arising from the existence of devices with low computational capabilities [19], [20], which may even be temporally unavailable while the local training takes place [21]. The heterogeneous nature of edge devices induces possibly substantial variations between different clients in their local update latency. This in turn affects the time it takes the server to update the global model on each round. As a result, FL is sensitive to *stragglers*, as each round takes as long as the local training time of the slowest user, dictating latency and throughput implications [3]. The presence of heterogeneous users thus makes synchronous FL abortive for applications with tight latency constraints.

Various schemes were proposed to provide robustness against stragglers in FL [22]-[28], see also survey [20]. Conventional synchronous FL limits local computation latency by imposing a deadline, while discarding a fixed-size set of delayed stragglers and their contributions [29], [30]. This can be extended to support varying deadlines via user selection, where only a subset of the users participate in each round [31], by identifying and grouping potential stragglers [22], [23]. Deadline-based synchronous FL facilitates the incorporation of latency constraints without altering the learning procedure, and can be combined with additional FL latency reduction techniques based on scheduling [8] and resource allocation [10]. However, the fact that stragglers are discarded affects performance when operating under low latency requirements, where a large portion of the users may not meet the local computation deadline.

Alternative approaches to handle stragglers involve deviating from the conventional operation of FL by either altering the learning procedure, or by switching to an asynchronous operation [20]. Several examples for the former are dedicated aggregation [26]; introducing redundancy on the devices' data via distributed codes [24], [32]; and altering the amount of local workload and architecture [19], [25]. However, these approaches assume a limited portion of straggling users and

This work was partially supported by the Israeli Ministry of Science and Technology. N. Lang and N. Shlezinger are with the School of ECE, Ben-Gurion University of the Negev, Be'er-Sheva, Israel (e-mails: langn@post.bgu.ac.il; nirshl@bgu.ac.il). A. Cohen is with the Faculty of ECE, Technion – Israel Institute of Technology, Haifa, Israel (e-mail: alecohen@technion.ac.il).

are thus unsuitable for low-latency settings, where many users might struggle to locally compute sufficient local training iterations in time. The asynchronous approach operates FL without requiring participating users to aggregate on each round [27], [28], [33]. For instance, TimelyFL proposed in [34] was experimentally shown to mitigate stragglers in asynchronous settings by having the server and the clients share calculation times to dictate a latency and local workload while utilizing layer-wise aggregations. Whereas asynchronous FL allows slower clients to continue the local training and contribute to future aggregation rounds, it also requires the number of slow computing users to be small for stable learning [20], limiting its applicability under tight latency constraints. Moreover, having different nodes operate on different versions of the global model leads to staleness [35], complicating the orchestration of the FL procedure compared to synchronous FL, while often ending up with an inferior performance.

In this work, we propose straggler-aware layer-wise federated learning (SALF), for FL with time varying system heterogeneity, that enables synchronous deadline-based high-performance low-latency operation. SALF is particularly geared for learning deep neural networks (DNNs), being the common family of machine learning models considered in FL, while exploiting their gradient-based training operation. We leverage the inherent recursive nature in which DNN are being optimized, i.e., the fact that the empirical risk gradient with respect to the model weights is computed from the last layer to the first via backpropagation [36], and that gradient computation typically induces a dominant portion of the latency. This operation indicates that stragglers may still compute lastlayers gradients that can be utilized rather than discarded under strict deadlines. SALF is thus based on this layer-wise approach, allowing users to convey (possibly partial) gradients when local training expires. Then, to update the global model, SALF averages the local updates *per layer*, while preserving the simplicity of conventional federated averaging (FedAvg) [2], and having different contributing devices for each.

To capture the realistic dynamic heterogeneity of the users, we analyze SALF assuming a probabilistic model on the straggling clients IDs per round. Specifically, in each round, the stragglers' set ranges from being empty to contain all FL users, neither limiting its cardinality nor its items. In our analysis we rigorously prove that SALF converges to the optimal model in the same asymptotic rate as local stochastic gradient descent (SGD) [30], while characterizing an upper bound on the gap in its learning objective value compared to the optimal model in the non-asymptotic regime.

We extensively evaluate SALF for the federated training of different DNN architectures, considering various latency constraints. Our numerical results show that SALF allows reliable training under tight latency constraints where a large bulk of the users become stragglers, while achieving similar accuracy to conventional FedAvg with no latency requirements.

The rest of this paper is organized as follows: Section II briefly reviews the FL system and the heterogeneity models. Section III presents SALF along with its convergence analysis. We numerically evaluate SALF in Section IV, and provide concluding remarks in Section V.

Throughout this paper, we use boldface lower-case letters for vectors, e.g.,  $\boldsymbol{x}$ , and calligraphic letters for sets, e.g.,  $\mathcal{X}$ , with  $|\mathcal{X}|$  being the cardinality of  $\mathcal{X}$ . The stochastic expectation, probability operator, and  $\ell_2$  norms are denoted by  $\mathbb{E}[\cdot]$ ,  $\mathbb{P}[\cdot]$ , and  $\|\cdot\|$ , respectively, while  $\mathbb{R}$  is the set of real numbers.

## II. SYSTEM MODEL

In this section, we set the ground for the derivation of SALF. We commence by presenting the system model of synchronous FL in Subsection II-A. Then, we provide a description of the system heterogeneity model and its effects on local computation latency in Subsection II-B.

#### A. Federated Learning

We consider a central server training a model with parameters  $\boldsymbol{w} \in \mathbb{R}^m$  using data available at U users, where each is indexed by  $u \in \{1, \ldots, U\}$ . Unlike conventional centralized learning, these datasets, denoted  $\mathcal{D}_1, \ldots, \mathcal{D}_U$ , cannot be shared with the server. Thus, by letting  $F_u(\boldsymbol{w})$  be the empirical risk of a model  $\boldsymbol{w}$  evaluated with dataset  $\mathcal{D}_u$ , FL aims to recover the  $m \times 1$  optimal weights vector,  $\boldsymbol{w}_{opt}$ , satisfying

$$\boldsymbol{w}_{\mathrm{opt}} = \operatorname*{arg\,min}_{\boldsymbol{w}} \left\{ F(\boldsymbol{w}) \triangleq \sum_{u=1}^{U} \frac{1}{U} F_u(\boldsymbol{w}) \right\},$$
 (1)

where it is implicitly assumed that the local datasets are of balanced (equal) cardinality.

Generally speaking, FL operates in rounds where for each time step t, the server distributes the global model  $w_t$  to the users, who each locally trains it, and sends back the model updates [5]. In conventional synchronous FL, the server collects the model updates from all the participating users, aggregates the models into an updated global model, and the overall procedure repeats iteratively.

We focus on settings where w represents a DNN with gradient-based local training. Here, each user of index ucomputes the stochastic gradient of its local empirical risk  $F_u$  evaluated on the global model at time  $t \ge 1$ ,  $w_t$ ; i.e.,  $\nabla F_u(w_t; i_t^u)$ , where  $i_t^u$  denotes the data sample index, chosen uniformly from  $\mathcal{D}_u$ . Then, for a step-size  $\eta_t$ , the user shares its local update (gradients), i.e,

$$\boldsymbol{w}_{u,t} \triangleq \boldsymbol{w}_t - \eta_t \nabla F_u \left( \boldsymbol{w}_t; i_t^u \right), \tag{2}$$

with the server who updates the global model. Conventional aggregation of the local updates is based on FedAvg [2], in which the server sets the global model to be

$$\boldsymbol{w}_{t+1} \triangleq \sum_{u=1}^{U} \frac{1}{U} \boldsymbol{w}_{u,t} = \boldsymbol{w}_t - \eta_t \sum_{u=1}^{U} \frac{1}{U} \nabla F_u\left(\boldsymbol{w}_t; i_t^u\right). \quad (3)$$

The updated global model is again distributed to the users, and the learning procedure continues until convergence is reached.

FL of DNNs involves edge users training, where each user obtains its local gradients in (3) using its local computational capabilities. The users' devices can notably vary in their computational resources based on their hardware and instantaneous operation. This property gives rise to the core challenge of device heterogeneity gap, discussed next.



Fig. 1. A device-heterogeneous FL-aided system learning object recognition that is expected to operate under tight latency and edge power constraints. Note that  $w_t$  and  $w_{u,t}$  denote the global and local models, respectively; where  $T_t^u$  it the local computational time of user u in FL round t.

#### B. System Heterogeneity

Edge devices participating in FL can widely differ in their computational powers, leading to varying processing times for calculating gradients [23]. Specifically, let  $T_t^u$  denote the time it takes the *u*th user to compute the gradient at the *t*th round, i.e.,  $\nabla F_u(w_t; i_t^u)$ . Higher values of  $T_t^u$  are attributed with clients who are slower at the *t*th round due to, e.g., limited hardware or additional external computations being carried out. These slower users are termed *stragglers*. Hence, the local computation latency of the *t*th FL round is given by  $\max_u T_t^u$ , as the server has to wait for the slowest user. This limits FL applications with tight latency constraints, e.g., intelligent transportation systems [37], as illustrated in Fig. 1.

A requirement for executing an FL round with a fixed latency is satisfied by setting a deadline  $T_{\text{max}}$ . Deadline-based synchronous FL typically discards the stragglers at round t for which  $T_t^u > T_{\text{max}}$ , i.e., those not meeting the threshold [29]. In such cases, denoting the set of users meeting the deadline at round t as  $U_t \triangleq \{u : T_t^u \leq T_{\text{max}}\}$ , the aggregation rule in (3) is replaced with

$$\boldsymbol{w}_{t+1} = \boldsymbol{w}_t - \eta_t \sum_{u \in \mathcal{U}_t} \frac{1}{|\mathcal{U}_t|} \nabla F_u\left(\boldsymbol{w}_t; \boldsymbol{i}_t^u\right).$$
(4)

Aggregation via (4) guarantees that each round does not surpass the desired deadline. However, for low-latency, i.e., small  $T_{\rm max}$ , this could result in a few users participating in each round, degrading the training procedure. Such a scenario motivates the deriviation of a scheme which allows FL to operate with small  $T_{\rm max}$  without fully discarding stragglers, as proposed next.

#### III. STRAGGLER-AWARE LAYER-WISE FL

In this section, we introduce SALF, formulating its operation in Subsection III-A. Then, in Subsection III-B we analyze the convergence properties of DNNs trained in a federated manner using SALF, and provide a discussion in Subsection III-C. For clarity, we summarize the symbols and notations used throughout this section in Table I.

#### A. SALF Algorithm

Consider the gradient-based federated training of an *L*-layered DNN with layer-wise parameters

$$\boldsymbol{w}_t = \begin{bmatrix} \boldsymbol{w}_t^1, \dots, \boldsymbol{w}_t^L \end{bmatrix}^T, \tag{5}$$

where layers 1 and L correspond to the input and output layers, respectively. The local gradients  $\nabla F_u$  are typically computed using backpropagation [36]. Namely, gradients are recursively calculated from the last-to-first layer, constructing

$$\left[\nabla F_u^1(\boldsymbol{w}_t; i_t^u), \dots, \nabla F_u^l(\boldsymbol{w}_t; i_t^u), \dots, \nabla F_u^L(\boldsymbol{w}_t; i_t^u)\right]^T, \quad (6)$$

each component at a time, where  $F_u^l(\boldsymbol{w}_t; i_t^u)$  denotes the (stochastic) gradient computed with respect to the parameters of the *l*th layer. This operation suggests that when  $T_{\max}$  expires, stragglers are likely to evaluate partial gradients, corresponding to the last DNN layers.

To exploit this expected behavior, we design SALF not to discard stragglers, but to have them contribute to the aggregation of these intermediate layers. This is done by aggregating via FedAvg in (3) over each layer separately. Since a device's computational power availability for an FL round is assumed to diverse along the training, the depth reached in backpropagation when the deadline expires is expected to vary between users and rounds. In our analysis provided in the sequel, it is modeled as a discrete random variable obeying a uniform distribution.

SALF, illustrated in Fig. 2 and summarized as Algorithm 1, affects two aspects of conventional synchronous FL: (*i*) *local training* at the users side; and (*ii*) *aggregation* of the model updates by the server. The resulting operation is formalized below.

1) Users Local Training: As in conventional DNN training, each user computes  $\nabla F_u(\boldsymbol{w}_t; i_t^u)$  via backpropagation, sequentially, from the last layer to the first. When the deadline  $T_{\text{max}}$  expires, the *u*th user calculates the gradients up to layer  $d_t^u \in \{1, \dots, L, L+1\}$ , denoting its associated depth; where



Fig. 2. Illustrative overview of SALF for training a deep CNN. The left dashed-box represents local training, where colored layers correspond to gradients calculated within  $T_{max}$ ; the left dashed-box shows the layer-wise aggregation with updated colored global model layers.

Al	Algorithm 1: SALF at round t									
1 I	1 Initialization: FL round running time $T_{\text{max}}$ ;									
2 Users side:										
3	do in parallel for each $u$ , until deadline $T_{\max}$ :									
4	Compute $\nabla F_u(\boldsymbol{w}_t, i_t^u)$ up to layer $d_t^u$ ;									
5	Convey partial gradients to server;									
0										
6 8	6 Server side:									
7	for $1 \le l \le L$ do									
8	Recover $\mathcal{U}_t^l$ ; $\triangleright$ the users updating the <i>l</i> th layer									
9	Compute $\tilde{\boldsymbol{w}}_{t+1}^l$ via (9); $\triangleright$ layer-wise update									
F	<b>Result:</b> The updated global model, $\tilde{w}_{t+1}$ :									

 $d_t^u = L + 1$  stands for a straggler that does not compute any layer gradient.

Accordingly, if  $d_t^u \leq L$ , the *u*th user conveys to the server the  $L - d_t^u + 1$  sub-vectors of its empirical risk stochastic gradient evaluated on  $w_t$ , i.e.,

$$\left[\nabla F_u^{d_t^u}(\boldsymbol{w}_t; i_t^u), \dots, \nabla F_u^L(\boldsymbol{w}_t; i_t^u)\right]^T.$$
(7)

This operation implies that users that are typically viewed as stragglers also convey (partial) model updates to the server.

2) Server Aggregation: Let  $\mathcal{U}_t^l \subseteq \{1, \ldots, U\}$  be the set of users that managed to compute the *l*th layer gradient on round *t*. Using the above notations, this set is given by

$$\mathcal{U}_t^l \triangleq \{ u : d_t^u \le l \},\tag{8}$$

where  $\mathcal{U}_t^1 \subseteq \cdots \subseteq \mathcal{U}_t^l \subseteq \cdots \subseteq \mathcal{U}_t^L$ . This follows by the backpropagation operation, as if a gradient was calculated up to layer  $d_t^u < L$ , so do all the gradients of layers  $d_t^u + 1, \ldots, L$ .

The server recovers  $\{\mathcal{U}_{t}^{l}\}_{l=1}^{L}$  from the received gradients, which are then aggregated via a form of *layer-wise* FedAvg. Rather than using FedAvg over the full DNN as in (3), the model updates obtained by SALF, denoted by  $\tilde{w}_{t+1}$ , are aggregated in a layer-wise fashion. Specifically, the *l*th layer of the global model, according to the decomposition in (5), is updated via

$$\tilde{\boldsymbol{w}}_{t+1}^{l} = \begin{cases} \boldsymbol{w}_{t}^{l} & |\mathcal{U}_{t}^{l}| = 0, \\ \frac{1}{1-p_{l}} \left( \sum_{u \in \mathcal{U}_{t}^{l}} \frac{1}{|\mathcal{U}_{t}^{l}|} \boldsymbol{w}_{u,t}^{l} - p_{l} \boldsymbol{w}_{t}^{l} \right) & |\mathcal{U}_{t}^{l}| > 0; \end{cases}$$
(9)

TABLE I SUMMERY OF NOTATIONS



where  $p_l \in [0, 1)$  is a hyperparameter set such that  $\tilde{w}_{t+1}^l$  is an unbiased estimator of the corresponding stragglers-free FedAvg update  $w_{t+1}^l$ . For instance, when the  $\{d_t^u\}$  are i.i.d with the uniform distribution,  $p_l$  is given by

$$p_l = \left(1 - \frac{l}{L+1}\right)^U,\tag{10}$$

see further details in Subsection III-B.

Note that in (9), as opposed to conventional FedAvg (3), the set of participating users can vary between layers. This is because, when operating under a deadline that imposes a fixed latency, different users compute their gradients up to different layers (see Fig. 2), depending on their computational resources and availability at that round, encapsulated in the stochastic  $d_t^u$ . As each layer of the global model is updated by different users at different rounds, the convergence analysis of SALF deviates from traditional FL [30]. Yet, due to the formulation of the model updates in (9), we are still able to rigorously prove convergence under conventional modelling assumptions, as reviewed next.

#### B. Analysis

As described in the previous subsection, SALF tackles system heterogeneity using a layer-wise aggregation rule to update the global model in each iteration. Here, we theoretically characterize the convergence profile of SALF. Note that the layer-wise operation of SALF and the non-zero probability of having the first layers not updated in given rounds indicate that existing FL analysis with partial device participation, e.g., [30, Thm. 3], do not apply here. We first elaborate bellow the assumptions introduced in our analysis as well as the chosen statistical characteristics of the stragglers; from which the convergence bound is subsequently derived.

1) Analysis Assumptions: We carry out our analysis of SALF subject to the following assumptions, that are commonly employed in FL convergence studies [15], [18], [38], [39]:

ASI The local objectives  $\{F_u(\cdot)\}_{u=1}^U$  are  $\rho_c$  strongly convex and  $\rho_s$ -smooth. That is, for all  $w_1, w_2 \in \mathbb{R}^m$ , it holds that

$$\begin{aligned} (\boldsymbol{w}_1 - \boldsymbol{w}_2)^T \nabla F_u(\boldsymbol{w}_2) + \frac{1}{2} \rho_c \|\boldsymbol{w}_1 - \boldsymbol{w}_2\|^2 \\ &\leq F_u(\boldsymbol{w}_1) - F_u(\boldsymbol{w}_2) \leq \\ (\boldsymbol{w}_1 - \boldsymbol{w}_2)^T \nabla F_u(\boldsymbol{w}_2) + \frac{1}{2} \rho_s \|\boldsymbol{w}_1 - \boldsymbol{w}_2\|^2. \end{aligned}$$

AS2 For each user u and round index t, the variance of the stochastic gradients  $\nabla F_u(\boldsymbol{w}; i_t^u)$  is bounded by some  $\sigma_u^2$  for all  $\boldsymbol{w} \in \mathbb{R}^m$ , i.e.,

$$\mathbb{E}\left[\left\|\nabla F_u(\boldsymbol{w}; i_t^u) - \nabla F_u(\boldsymbol{w})\right\|^2\right] \leq \sigma_u^2.$$

AS3 For each user u and round index t, the expected squared  $\ell_2$  norm of the stochastic gradients  $\nabla F_u(\boldsymbol{w}; i_t^u)$  is uniformly bounded by some  $G^2$  for all  $\boldsymbol{w} \in \mathbb{R}^m$ , i.e.,

$$\mathbb{E}\left[\left\|\nabla F_u(\boldsymbol{w}; i_t^u)\right\|^2\right] \leq G^2.$$

AS4 For each user u and round index t, the depth reached in backpropagation, denoted  $d_t^u$ , is i.i.d. (in t and u), and uniformly distributed over the set  $\{1, \ldots, L+1\}$ .

Notice that smoothness, as assumed in ASI, holds for a range of objective functions used in FL, including  $\ell_2$ -norm regularized linear/logistic regression [15]. The heterogeneity between the users is also reflected in AS2 via the dependence on u, which implies that the local objectives differ between users, as different datasets can be statistically heterogeneous, i.e., arise from different distributions. Accordingly, we follow [30] and define the (datasets) heterogeneity gap as

$$\Gamma \triangleq F(\boldsymbol{w}_{\text{opt}}) - \frac{1}{U} \sum_{u=1}^{U} \min_{\boldsymbol{w}} F_u(\boldsymbol{w}), \qquad (11)$$

where  $w_{\text{opt}}$  is given in (1).

2) Stragglers Statistical Modelling: The dynamic nature of the system's devices heterogeneity (see Subsection II-B) is statistically modelled in AS4, as the set of random variables  $\{d_t^u\}$  in AS4 form the stragglers sets  $\{\mathcal{U}_t^l\}$  (8). This formulation subsumes the extreme cases of both  $|\mathcal{U}_t^l| = 0$  (a layer which is not updated by any user) and  $|\mathcal{U}_t^l| = U$  (a layer that is updated by all users). In particular, the marginal distribution of  $|\mathcal{U}_t^l|$  is obtained in the following lemma:

**Lemma 1.** When AS4 holds, then for every FL round t and DNN layer l, and the cardinality of the set  $\mathcal{U}_t^l$  is distributed as

$$|\mathcal{U}_t^l| \sim \operatorname{Bin}\left(U, \frac{l}{L+1}\right),$$
 (12)

where Bin denotes the Binomial distribution.

*Proof:* Eq. (12) follows from the definition of  $\mathcal{U}_t^l$  in (8), which implies that for every  $M \in \{0, \ldots, U\}$ , it holds that

$$\mathbb{P}[|\mathcal{U}_t^l| = M] = \mathbb{P}[|\{u : d_t^u \le l\}| = M]$$
$$= \mathbb{P}\left[\left(\sum_{u=1}^U \mathbf{1}_{d_t^u \le l}\right) = M\right]$$

where  $\mathbf{1}_{\{\cdot\}}$  is the indicator function. By *AS4*, it follows that  $\{\mathbf{1}_{d_t^u \leq l}\}_{u=1}^U$  are i.i.d. Bernoulli random variables with parameter  $\mathbb{P}[\mathbf{1}_{d_t^u \leq l} = 1] = \frac{l}{L+1}$ , proving (12) by the definition of the binomial distribution [40, Ch. 4].

According to Lemma 1, for a given global iteration, the chance that a particular layer would not be updated at all decreases for either deeper layers or growing number of total users U. Lemma 1 is thereby used to establish two auxiliary lemmas: The first identifies that the SALF update rule, which yields a random vector due to the stochastic nature of the stragglers set, is an unbiased estimate of the full FedAvg rule, as stated next:

**Lemma 2** (Unbiasedness). When AS4 holds, for every FL round of index t and a given set of training data samples  $\{i_t^u\}$ , the global model aggregated via SALF (9) with (10) is an unbiased estimator of the one obtained via vanilla FedAvg (3), namely,

$$\mathbb{E}[\tilde{\boldsymbol{w}}_{t+1}] = \boldsymbol{w}_{t+1} \tag{13}$$

*Proof:* The proof is given in Appendix A.

The second auxiliary lemma which follows from Lemma 1 bounds the variance of the SALF DNN parameters.

**Lemma 3** (Bounded variance). Consider SALF where AS4 holds and  $p_l$  is set via (10), and for every FL round t and a given set of training data samples  $\{i_t^u\}$ , the learning rate  $\eta_t$ is set to be non-increasing and satisfying  $\eta_t \leq 2\eta_{t+1}$ . Then, the expected difference between  $\tilde{w}_{t+1}$  and  $w_{t+1}$ , which is the variance of  $\tilde{w}_{t+1}$ , is bounded by

$$\mathbb{E}[\|\tilde{\boldsymbol{w}}_{t+1} - \boldsymbol{w}_{t+1}\|^2] \le \eta_t^2 \frac{4ULG^2}{(U-1)} \frac{1 + \left(1 - \frac{1}{L+1}\right)^U}{1 - \left(1 - \frac{1}{L+1}\right)^U}.$$
 (14)

*Proof:* The proof is given in Appendix B.

*3) Convergence Bound:* The characterization of the DNN parameters produced by SALF as a bounded-variance unbiased stochastic estimate of the parameters produced with full FedAvg, allows to characterize the convergence of the learning algorithm. Using the above notations, the following theorem establishes the convergence bound of SALF.

**Theorem III.1.** Consider SALF-aided FL with (10) satisfying AS1-AS4; define  $\kappa = \frac{\rho_s}{\rho_c}, \gamma = \max\{8\kappa, 1\}$ , and set the learning rate to  $\eta_t = \frac{2}{\rho_c(\gamma+t)}$ . Then, it holds that

$$\mathbb{E}\left[F(\tilde{\boldsymbol{w}}_{t})\right] - F(\boldsymbol{w}_{\text{opt}}) \leq \frac{\kappa}{\gamma + t - 1} \times \left(\frac{2(B+C)}{\rho_{c}} + \frac{\rho_{c}\gamma}{2} \mathbb{E}\left[\|\boldsymbol{w}_{1} - \boldsymbol{w}_{\text{opt}}\|^{2}\right]\right), \quad (15)$$

where

$$B = \sum_{u=1}^{U} \frac{1}{U^2} \sigma_u^2 + 6\rho_s \Gamma; \qquad (16a)$$

$$C = \frac{4ULG^2}{(U-1)} \frac{1 + \left(1 - \frac{1}{L+1}\right)^U}{1 - \left(1 - \frac{1}{L+1}\right)^U}.$$
 (16b)

# *Proof:* The proof is given in Appendix C.

Theorem III.1 rigorously bounds the difference between the objective value of a model learned by SALF at round t to the optimal model  $w_{opt}$ . By setting the step size  $\eta_t$  to decrease gradually, which is also known to contribute to the convergence of FL [30], [38], Theorem III.1 indicates that SALF converges at a rate of  $\mathcal{O}(1/t)$ . This asymptotic rate is identical to FL schemes with no latency limitations [30], [38], stressing the ability of SALF to carry out low latency FL in the presence of dynamic system heterogeneity while mitigating its harmful effects with hardly affecting the learning procedure compared to conventional FL.

Nonetheless, in the non-asymptotic regime, the integration of a latency deadline influences model convergence. This is revealed in (15) by the constant C in (16), that depends on (growing at least linearly in) the number of layers L, resulting from SALF's layers-wise aggregation technique. This implies that synchronous FL with deadline  $T_{\text{max}}$ , whose values results in some of the users being occasionally stragglers, is expected to converge slower for deeper architectures; for which C is larger compared with shallower DNNs.

Such behaviour is also experimentally demonstrated in Section IV, stemming from the fact that for DNNs with many layers, the first few layers will be trained by a small portion of the participating users when operating under a fixed deadline compared with shallow layer. This is in line with DNN typical behaviour, as deeper model architectures aim to learn more complex mappings, often require lengthy learning and are slower to converge [41], [42]. It is emphasized though that under conventional (non layer-wise) synchronous FL, training deep DNNs under similar low latency constraints is expected to often be infeasible, as effectively all users become stragglers.

## C. Discussion

SALF is particularly designed for FL systems that are constrained to operate with low latency while learning over networks comprised of heterogeneous edge devices. It allows synchronous FL operation with small deadline  $T_{\rm max}$ , which typically results in a large number of stragglers. This is achieved without notably affecting the conventional FedAvgbased FL flow by exploiting the recursive nature of the backpropagation to leverage partial gradients for updating the global model. Due to its simple layer-wise aggregation, SALF asymptotically converges at the same rate as unconstrained FedAvg, for a random set of stragglers, while supporting extremely low-latency FL. For instance, we numerically show in Section IV that SALF can learn accurate models in settings where as much as 90% of the users are stragglers that cannot finish computing their gradients in time.

Our convergence analysis of SALF assumes that the local computations of each user behave randomly. The distribution imposed in *AS4* accounts for the fact that a device can become a straggler in a given round not only due to its fundamental hardware, but also due to its availability on that particular round, as edge devices may be occupied also for other tasks than FL. Conversely, if we were to utilize SALF's layer-wise aggregation using *deterministic* sets of stragglers in which the same users update the same layers in each training round, an inevitable bias is expected to arise as some of the overall data would not affect all the layers. In that sense, the stochastic nature of system heterogeneity is leveraged as a contributing factor, eliminating the probability of such a scenario to occur.

The outline of SALF is based on a generic formulation of synchronous gradient-based FL. It only requires the learned model to be a DNN trained using backpropagation. While this form of gradient-based learning by far dominates DNN training methods to-date, alternative ones such as Kalmanbased learning [43] and zero-order optimization [44] were also proposed in the literature, for which SALF would require a dedicated adaptation. Additionally, SALF operates without increasing the complexity at the clients and/or the server compared to conventional FedAvg. Hence, its methodology can be combined with other schemes for decreasing FL latency via model update compression [15], [18], [39]. While our design considers a single local iteration, being tailored to tight deadlines, it can be extended to multiple iterations and possibly combined with proximal-aided aggregation to account for user-varying iterations, e.g., [26], on a layer-wise basis. These extensions of SALF are left for future study.

#### IV. EXPERIMENTAL STUDY

In this section we numerically evaluate SALF, and compare it to existing approaches [2], [19], [20], [29] tackling system heterogeneity in low-latency FL<sup>1</sup>. Our aim is to experimentally validate that the layer-wise approach of SALF allows to learn reliable DNN models for various architectures in a synchronous low-latency manner. We focus on the training of DNNs for image processing tasks, where we first consider a simple handwritten digit recognition task (Subsection IV-A), which allows us to evaluate SALF in different terms of performance as convergence, accuracy, and latency; in controlled settings. Then, we proceed to a more challenging image classification task (Subsection IV-B), where we evaluate performance as well as suitability for different deep architectures.

#### A. Handwritten Digit Recognition

1) Setup: We first consider the federated training of a handwritten digit classification model using the MNIST dataset [45]. The data, comprised of  $28 \times 28$  gray-scale images divided into 60,000 training examples and 10,000 test examples, is uniformly distributed among U = 30 users.

<sup>1</sup>The source code used in our experimental study, including all the hyperparameters, is available online at https://github.com/langnatalie/SALF.



Fig. 3. FL convergence profile, MLP trained on MNIST.

Architectures: We train two different DNN architectures:

- A multi-layer perceptron (MLP) with two hidden layers, intermediate ReLU activations and a softmax output layer.
- A CNN composed of two convolutional layers and two fully-connected ones, with intermediate ReLU activations, max-pooling layers, and a softmax output layer.

**FL Training**: In each FL iteration, the users train the MLP/CNN model using local mini-batch SGD (3) with learning rate 0.05/0.1; and the global model is learned using 250/150 FL iterations, respectively. To simulate controllable latency-constrained system heterogeneity for different deadlines in a hardware-invariant manner, in each iteration we randomly set a predefined ratio of the users to be stragglers, where for each straggler the depth reached in backpropagation is randomized uniformly (in line with the distribution assumed in *AS4*).

- FL Algorithms: We evaluate the following FL methods:
- *Vanilla FL*, which implements full FedAvg without latency constraints [2] and thus without any stragglers. This approach constitutes the desired performance for the methods that operate in the presence of stragglers.
- Drop-stragglers FedAvg, that discards stragglers [29].
- *HetroFL* [19] which addresses heterogeneous clients by equipping them with corresponding heterogeneous local models with varying computational complexities. To guarantee fair comparison with SALF in the sense of average number of gradient computations, we set the straggling users to shrink their local models with ratio 0.5. Comparison to HetroFL is considered only for the CNN, as this architecture was covered in [19].
- AsyncFL [20] which has the stragglers participate in the FedAvg once they finish the local training rather than at each round. Here, to achieve comparable computations to synchronous settings, a straggler that finished the update up to the *l*th layer within  $T_{\rm max}$ , is set to asynchronously participate in the FedAvg every 2*l* global iterations.
- Our proposed *SALF* (Algorithm 1).

2) Results: We compare the performance of the above FL methods for the considered task in their convergence profile



Fig. 4. FL convergence profile, CNN trained on MNIST

(namely, during training); the accuracy of the trained model (i.e., after training is concluded); and their overall latency.

**Convergence Results**: We begin by evaluating the convergence profiles of the learning methods, i.e., the resulting validation accuracy achieved over the training procedure. To emphasize the gains of SALF compared with conventional drop-stragglers in facilitating federated training under tight latency constraints, we focus on the extreme case in which 90% of the users are stragglers. This exemplifies a tight deadline-based FL application, where the major bulk of the users effectively become stragglers that do not complete the whole model update. Fig. 3-4 illustrate the convergence profile of the methods for the MLP and CNN models, respectively. There, it is systematically demonstrated that SALF is the closet to vanilla FL, whereas AsyncFL is the second best (while requiring an asynchronous operation). For comparison, drop-stragglers and HetroFL show an inferior performance, as a result of ignoring stragglers contributions for the former, and having the majority of local architectures differ from the archietcture of the global model for the latter.

Accuracy Results: We proceed with examining the performance of the trained models in terms of their test accuracy. Table II summarizes the test accuracy result for both the MLP and CNN models trained on the MNIST datasets for either 30/50/70/90 percent of the users being stragglers. Table II reveals that for all the considered techniques, as expected, the higher the percentage of stragglers, the lower the test accuracy is. This monotonic behaviour results with dramatic degradation for growing percentages (corresponding to tight deadlines) for the synchronous drop-stragglers and HetroFL. Despite that, the degradation of SALF is notably lower compared to all other baselines, maintaining a minor gap from vanilla FL (which operates without latency constraint), and consistently achieving the best performance of the trained model among all stragglers-constrained methods.

**Latency Results**: As discussed in Subsection II-B, for timing-based FL deployments, the deadline  $T_{\text{max}}$  determines the latency. To translate the improved resilancy to stragglers of FL reported above into concrete timings, we evaluate the overall test performance of all considered FL methods under

 TABLE II

 Test accuracy results for different stragglers' percents, MNIST dataset.

	Vanilla FL	Drop-stragglers				SALF				HetroFL				AsyncFL			
Straggler %	N/A	0.3	0.5	0.7	0.9	0.3	0.5	0.7	0.9	0.3	0.5	0.7	0.9	0.3	0.5	0.7	0.9
MLP	0.9	0.87	0.84	0.77	0.49	0.88	0.85	0.85	0.81	N/A	N/A	N/A	N/A	0.86	0.84	0.82	0.73
CNN	0.95	0.93	0.9	0.83	0.28	0.94	0.93	0.92	0.90	0.88	0.84	0.61	0.43	0.94	0.92	0.91	0.86



Fig. 5. Test set accuracy vs. latency constrains, CNN trained on MNIST.

different timing constraints  $T_{\text{max}}$ . We focus on the 2-layer CNN model, for which the full backward pass for a single user was empirically evaluated as taking the maximal value of 7.5 microseconds ( $\mu$ sec), when computed using a Quadro RTX 6000 GPU. The resulting test accuracy versus  $T_{\text{max}}$  in the range of  $[1, 9]\mu$ sec are reported in Fig. 5 for the considered FL methods, in comparison with vanilla FL (which operates without latency constraints and is thus invariant of  $T_{\text{max}}$ ).

We first clearly observe in Fig. 5 that, as expected, all FL methods coincide with full FedAvg when the deadline surpasses the maximal local computation latency of 7.5  $\mu$ sec. This follows since in such latency regimes, none of the users are stragglers. However, in the more interesting regimes of  $T_{\rm max} < 7.5 \ \mu {
m sec}$ , we observe that SALF yields a minor desegregation in performance which hardly grows when the latency is decreased; and outperforms the (asynchronous) AsyncFL counterpart. This is in contrast to either drop-stragglers, which is left with hardly any devices updating the global model under tight latency constraints, or HetroFL, where in this case tighter latency implies that the majority of the users train a different model than the global one. Specifically, for low  $T_{\text{max}}$ , i.e., below 1.5  $\mu$ sec, which also corresponds to high percentages of straggling clients, HetroFL is superior to drop-stragglers in performance (where the opposite holds for moderate values of  $T_{\rm max}$ , i.e., for  $1.5 < T_{\rm max} < 7 \ \mu {\rm sec}$ ), and AsyncFL is superior to both; aligned with similar findings in Fig. 4 and Table II.

Finally, as evidenced in Fig. 5, for the tightest latency value, SALF realizes a drop of merely 5% from the accuracy of vanilla FL, compared to 15%, 70% and 90% in the case of the AsyncFL, HetroFL, and drop-stragglers, respectively. Consequently, the gains of SALF in performance are persistent, and most dominant in the low-latency regime, where stragglers

mostly fail to meet the deadline calculation time  $T_{\text{max}}$  and their partial updates are harnessed for modifying the global model by the layer-wise approach of SALF.

#### B. Image Classification

1) Setup: We proceed to evaluating SALF in tasks typically requiring deeper DNNs compared to the ones used in the previous subsection. Here, FL is implemented for the distributed training of natural image classification model using the CIFAR-10 dataset [46]. This set is comprised of  $32 \times 32$  RGB images divided into 50,000 training examples and 10,000 test examples.

**Architecture**: We explore whether the learning profile of a model trained via the layer-wise aggregation of SALF changes with varying the depth of its given architecture. To that aim, we use the VGG model architecture [47] with four depths, namely, number of convolutional layers, that are 11, 13, 16, and 19. FL training is similar to the one described in Subsection IV-A, while using a learning rate of 0.05 and set the amount of global iterations to be 1,500 for each architecture.

2) Results: We first evaluate SALF for different timing constraints. Fig. 6 shows the convergence profile of the VGG models for vanilla FL and SALF operating with either 30/50/70 percents of straggling clients. It can be generally observed that, SALF converges for all considered architectures and amount of stragglers. However, as discussed in Subsection III-C, the training of deeper architectures exhibit different profiles compared with shallower ones, which are consistent for all considered stragglers percentages. A notable phenomenon is observed when training the deepest model of VGG19; There, SALF with 70% or 50% stragglers yields an improved learning procedure compared with 30%. This can be associated with the fact that when training deep models, adding minor levels of distortion, which in our case result from the growth in the variance of the stochastic estimate in Lemma 3 for deeper networks, can lead to improving the converged model, in accordance with similar findings in [48], [49].

To better highlight the interplay between our layer-wise FL and the overall DNN depth, we conclude our study by focusing merely on a challenging setup with tight deadlines. In Fig. 7, we depict SALF and drop-stragglers for 90% straggles. The results in Fig. 7 stress the power of SALF in mitigating their harmful effect significantly better than drop-stragglers, also for various deep architectures, similarly to the findings evidenced in Subsection IV-A. In addition, the convergence profiles of Fig. 7, once observed in comparison with Fig. 3-4, numerically support Theorem III.1 under the non-asymptotic regime; indeed indicating that low-latency FL



Fig. 6. Convergence profile of FL schemes training VGG models using the CIFAR-10 dataset for different stragglers percentages.

converges slower with deeper DNNs, while systematically exceeding drop-stragglers and still approaching the performance achieved with stragglers-free full FedAvg.

## V. CONCLUSIONS

In this work we proposed SALF, which is an FL algorithm that implements layer-wise global aggregation to incorporate stochastically determined straggling users in tight timings synchronous settings. SALF utilizes the last-to-first layer update policy of backpropagation-based DNN training to exploit partial gradients and update each of the model layers separately, with possibly different amount of users in each. We analyzed the convergence profile of SALF accompanied by numerical evaluations, demonstrating that it operates reliably under tight latency constraints and approaches the performance achieved by FL with no stragglers.

#### ACKNOWLEDGMENTS

The authors are grateful to Dor Elimelech for fruitful discussions and helpful ideas regarding the probabilistic modelling.

#### Appendix

## A. Proof of Lemma 2

In order to show that  $\tilde{w}_{t+1}$  (9) is an unbiased estimator of  $w_{t+1}$  (3), it suffices to show that this holds for the *l*th subvector of both, according to the decomposition of a vector into its *L* sub-vectors in (5). Assuming that the random data sample indexes  $\{i_t^u\}$  are given, the only source of randomness in  $w_{t+1}^l$  is encapsulated in  $\mathcal{U}_t^l$ . Yet, instead of calculating the expectation of  $\tilde{w}_{t+1}^l$  with respect to  $\mathcal{U}_t^l$ , we leverage the observation that  $|\mathcal{U}_t^l|$  is a binomial random variable (12), and utilize the law of total expectation, yielding

$$\begin{split} \mathbb{E}[\tilde{\boldsymbol{w}}_{t+1}^{l}] &= \mathbb{E}\left[\mathbb{E}\left[\tilde{\boldsymbol{w}}_{t+1}^{l} \middle| |\mathcal{U}_{t}^{l}|\right]\right] \\ &= \sum_{K=0}^{U} \mathbb{P}\left[|\mathcal{U}_{t}^{l}| = K\right] \cdot \mathbb{E}\left[\tilde{\boldsymbol{w}}_{t+1}^{l} \middle| |\mathcal{U}_{t}^{l}| = K\right] \\ &= \mathbb{P}\left[|\mathcal{U}_{t}^{l}| = 0\right] \cdot \mathbb{E}\left[\tilde{\boldsymbol{w}}_{t+1}^{l} \middle| |\mathcal{U}_{t}^{l}| = 0\right] \\ &+ \sum_{K=1}^{U} \mathbb{P}\left[|\mathcal{U}_{t}^{l}| = K\right] \cdot \mathbb{E}\left[\tilde{\boldsymbol{w}}_{t}^{l} \middle| |\mathcal{U}_{t}^{l}| = K\right]. \quad (A.1) \end{split}$$

10



Fig. 7. Convergence profile of FL schemes training VGG models using the CIFAR-10 dataset at 90% stragglers percentage.

Next, we note that for the distribution of  $|\mathcal{U}_t^l|$  in Lemma 1, alternative. Equivalently, the definition of  $p_l$  in (10) holds that

$$p_l = \mathbb{P}\left[|\mathcal{U}_t^l| = 0\right]. \tag{A.2}$$

Substituting (A.2) combined with SALF's aggregation rule in (9) into (A.1) results in

$$\mathbb{E}[\tilde{\boldsymbol{w}}_{t+1}^{l}] = p_{l} \cdot \boldsymbol{w}_{t}^{l} + \sum_{K=1}^{U} \mathbb{P}\left[|\mathcal{U}_{t}^{l}| = K\right] \frac{1}{1 - p_{l}} \times \left(\frac{1}{K} \mathbb{E}\left[\sum_{u \in \mathcal{U}_{t}^{l}} \boldsymbol{w}_{u,t}^{l} \middle| |\mathcal{U}_{t}^{l}| = K\right] - p_{l} \boldsymbol{w}_{t}^{l}\right). \quad (A.3)$$

Notice that, according to AS4, the distribution of  $\mathcal{U}_t^l$  given that  $|\mathcal{U}_t^l| = K$  is uniform over all size-K-out-of-size-U subsets. In [30, Lem. 4], the authors proved that FedAvg with partial device participation, such that a fixed-size amount of users are sampled uniformly in a without replacement fashion, is an unbiased estimator of the full device participation

$$\frac{1}{K}\mathbb{E}\left[\sum_{u\in\mathcal{U}_{t}^{l}}\boldsymbol{w}_{u,t}^{l}\middle||\mathcal{U}_{t}^{l}|=K\right]=\frac{1}{K}K\sum_{u=1}^{U}\frac{1}{U}\boldsymbol{w}_{u,t}^{l}=\boldsymbol{w}_{t+1}^{l}.$$
(A.4)

Plugging (A.4) into (A.3) and the fact that

$$\sum_{K=1}^{U} \mathbb{P}\left[ |\mathcal{U}_t^l| = K \right] = 1 - \mathbb{P}\left[ |\mathcal{U}_t^l| = 0 \right] = 1 - p_l$$

results with

$$\mathbb{E}[\tilde{\boldsymbol{w}}_{t+1}^{l}] = p_l \cdot \boldsymbol{w}_t^{l} + \left(\boldsymbol{w}_{t+1}^{l} - p_l \cdot \boldsymbol{w}_t^{l}\right) \frac{1}{1 - p_l} (1 - p_l)$$
$$= \boldsymbol{w}_{t+1}^{l},$$

thus proving (13).

#### B. Proof of Lemma 3

By the definition of the  $\ell_2$  norm, it follows that the variance of  $\tilde{w}_{t+1}$  is the sum over all its  $1, \ldots, L$  sub vectors (layers) variances; and, similarly to Appendix A, we have that

$$\mathbb{E}\left[\|\tilde{\boldsymbol{w}}_{t+1}^{l} - \boldsymbol{w}_{t+1}^{l}\|^{2}\right] = \mathbb{E}\left[\mathbb{E}\left[\|\tilde{\boldsymbol{w}}_{t+1}^{l} - \boldsymbol{w}_{t+1}^{l}\|^{2}||\mathcal{U}_{t}^{l}|\right]\right] \\
= \mathbb{P}\left[|\mathcal{U}_{t}^{l}| = 0\right] \mathbb{E}\left[\|\tilde{\boldsymbol{w}}_{t+1}^{l} - \boldsymbol{w}_{t+1}^{l}\|^{2}||\mathcal{U}_{t}^{l}| = 0\right] + \\
\sum_{K=1}^{U} \mathbb{P}\left[|\mathcal{U}_{t}^{l}| = K\right] \cdot \mathbb{E}\left[\|\tilde{\boldsymbol{w}}_{t+1}^{l} - \boldsymbol{w}_{t+1}^{l}\|^{2}||\mathcal{U}_{t}^{l}| = K\right]. \quad (B.1)$$

For the first summoned, it holds that

$$\mathbb{E}\left[\|\tilde{\boldsymbol{w}}_{t+1}^{l} - \boldsymbol{w}_{t+1}^{l}\|^{2}||\mathcal{U}_{t}^{l}| = 0\right] = \\\mathbb{E}\left[\|\boldsymbol{w}_{t}^{l} - \boldsymbol{w}_{t+1}^{l}\|^{2}\right] \stackrel{(a)}{\leq} \sum_{u=1}^{U} \frac{1}{U} \mathbb{E}\left[\|(\boldsymbol{w}_{u,t}^{l} - \boldsymbol{w}_{t}^{l})\|^{2}\right] = \\\sum_{u=1}^{U} \frac{1}{U} \mathbb{E}\left[\|\eta_{t} \nabla F_{u}(\boldsymbol{w}_{t}, i_{t}^{u})\|^{2}\right] \stackrel{(b)}{=} \eta_{t}^{2} G^{2}, \qquad (B.2)$$

where (a) follows by the convexity of  $\|\cdot\|^2$  and (b) stems from the bounded norm assumed by *AS3*.

As for the second summoned (B.1), it can be written as

$$\begin{split} & \mathbb{E}\left[\left\|\tilde{\boldsymbol{w}}_{t+1}^{l} - \boldsymbol{w}_{t+1}^{l}\right\|^{2} \left||\mathcal{U}_{t}^{l}| = K\right] = \\ & \mathbb{E}\left[\left\|\frac{1}{1-p_{l}}\left(\sum_{u\in\mathcal{U}_{t}^{l}}\frac{1}{|\mathcal{U}_{t}^{l}|} \; \boldsymbol{w}_{u,t}^{l} - p_{l}\boldsymbol{w}_{t}^{l}\right) - \boldsymbol{w}_{t+1}^{l}\right\|^{2} \left||\mathcal{U}_{t}^{l}| = K\right] \\ & = \frac{1}{(1-p_{l})^{2}} \times \\ & \mathbb{E}\left[\left\|\sum_{u\in\mathcal{U}_{t}^{l}}\frac{1}{|\mathcal{U}_{t}^{l}|} \boldsymbol{w}_{u,t}^{l} - \boldsymbol{w}_{t+1}^{l} - p_{l}(\boldsymbol{w}_{t}^{l} - \boldsymbol{w}_{t+1}^{l})\right\|^{2} \left||\mathcal{U}_{t}^{l}| = K\right] \\ & = \frac{1}{(1-p_{l})^{2}} \times \\ & \left(\mathbb{E}\left[\left\|\sum_{u\in\mathcal{U}_{t}^{l}}\frac{1}{|\mathcal{U}_{t}^{l}|} \boldsymbol{w}_{u,t}^{l} - \boldsymbol{w}_{t+1}^{l}\right\|^{2} \left||\mathcal{U}_{t}^{l}| = K\right] + \\ & \quad (B.3) \end{split}$$

$$p_{l}^{2} \mathbb{E}\left[\left\|\boldsymbol{w}_{t}^{l} - \boldsymbol{w}_{t+1}^{l}\right\|^{2} \middle| |\mathcal{U}_{t}^{l}| = K\right] +$$
(B.4)

$$-p_{l}(\boldsymbol{w}_{t}^{l}-\boldsymbol{w}_{t+1}^{l})^{T}\mathbb{E}\left[\sum_{u\in\mathcal{U}_{t}^{l}}\frac{1}{|\mathcal{U}_{t}^{l}|}\boldsymbol{w}_{u,t}^{l}-\boldsymbol{w}_{t+1}^{l}\middle||\mathcal{U}_{t}^{l}|=K\right]\right).$$
(B.5)

Now, (B.5) = 0 by (A.4); (B.4)  $\leq p_l^2 \eta_t^2 G^2$  by AS3; and finally, (B.3) can be bounded using the result obtained in [30, Lem. 5] due to the same reasons mentioned in Appendix A. That is, by [30, Lem. 5] it holds that

$$\mathbb{E}\left[\left\|\sum_{u\in\mathcal{U}_{t}^{l}}\frac{1}{|\mathcal{U}_{t}^{l}|}\boldsymbol{w}_{u,t}^{l}-\boldsymbol{w}_{t+1}^{l}\right\|^{2}\left||\mathcal{U}_{t}^{l}|=K\right]\leq\frac{U}{K(U-1)}\left(1-\frac{K}{U}\right)4\eta_{t}^{2}G^{2}.$$

Overall, (B.1) is thus given by

$$\begin{aligned} (\mathbf{B}.\mathbf{1}) &\leq \frac{1}{(1-p_l)^2} \sum_{K=1}^U \mathbb{P}\left[ |\mathcal{U}_t^l| = K \right] \times \\ &\left( \frac{U}{K(U-1)} \left( 1 - \frac{K}{U} \right) 4\eta_t^2 G^2 + p_l^2 \eta_t^2 G^2 \right) \\ &\leq \frac{1}{(1-p_l)^2} \sum_{K=1}^U \mathbb{P}\left[ |\mathcal{U}_t^l| = K \right] \cdot \left( \frac{U}{(U-1)} 4\eta_t^2 G^2 + p_l^2 \eta_t^2 G^2 \right) \\ &= \eta_t^2 \frac{G^2 \left( \frac{4U}{(U-1)} + p_l^2 \right)}{1-p_l}. \end{aligned}$$
(B.6)

Adding both bounds of (B.2) and (B.6) while summing over all *L* layers results with

$$\mathbb{E}[\|\tilde{\boldsymbol{w}}_{t+1} - \boldsymbol{w}_{t+1}\|^2] = \sum_{l=1}^{L} \mathbb{E}\left[\|\tilde{\boldsymbol{w}}_{t+1}^l - \boldsymbol{w}_{t+1}^l\|^2\right]$$
$$\leq \sum_{l=1}^{L} \eta_t^2 \frac{G^2\left(\frac{4U}{(U-1)} + p_l\right)}{1 - p_l} \leq \eta_t^2 G^2 \frac{4U}{(U-1)} \sum_{l=1}^{L} \frac{1 + p_l}{1 - p_l}$$
$$\leq \eta_t^2 G^2 \frac{4U}{(U-1)} L \frac{1 + \left(1 - \frac{1}{L+1}\right)^U}{1 - \left(1 - \frac{1}{L+1}\right)^U},$$

where the last inequality follows by bounding each fraction in the summation with the one obtained from setting the largest numerator and the smallest denominator; proving (14).

## C. Proof of Theorem III.1

By modelling the model updates of SALF as unbiased stochastic estimates of FedAvg (Lemma 2) with bounded variance (Lemma 3), we recast our setting of FL with random layer-wise computations as FL with random partial participation as considered in [30, Thm. 3]. Accordingly, by replacing Lemmas 4, 5 of [30] with our Lemmas 2, 3, respectively, we obtain that the derivation of [30, Thm. 3] applies for our setting of SALF with the corresponding coefficients in (16), thus proving the theorem.

#### REFERENCES

- E. Horvitz and D. Mulligan, "Data, privacy, and the greater good," Science, vol. 349, no. 6245, pp. 253–255, 2015.
- [2] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas, "Communication-efficient learning of deep networks from decentralized data," in *Artificial intelligence and statistics*. PMLR, 2017, pp. 1273– 1282.
- [3] P. Kairouz et al., "Advances and open problems in federated learning," Foundations and Trends<sup>®</sup> in Machine Learning, vol. 14, no. 1–2, pp. 1–210, 2021.
- [4] T. Li, A. K. Sahu, A. Talwalkar, and V. Smith, "Federated learning: Challenges, methods, and future directions," *IEEE Signal Process. Mag.*, vol. 37, no. 3, pp. 50–60, 2020.
- [5] T. Gafni, N. Shlezinger, K. Cohen, Y. C. Eldar, and H. V. Poor, "Federated learning: A signal processing perspective," *IEEE Signal Process. Mag.*, vol. 39, no. 3, pp. 14–41, 2022.
- [6] J. Chen and X. Ran, "Deep learning with edge computing: A review." Proc. IEEE, vol. 107, no. 8, pp. 1655–1674, 2019.
- [7] S. Samarakoon, M. Bennis, W. Saad, and M. Debbah, "Distributed federated learning for ultra-reliable low-latency vehicular communications," *IEEE Trans. Commun.*, vol. 68, no. 2, pp. 1146–1159, 2019.

- [8] W. Xia, W. Wen, K.-K. Wong, T. Q. Quek, J. Zhang, and H. Zhu, "Federated-learning-based client scheduling for low-latency wireless communications," *IEEE Wireless Commun. Lett.*, vol. 28, no. 2, pp. 32– 38, 2021.
- [9] Y. Lu, X. Huang, K. Zhang, S. Maharjan, and Y. Zhang, "Low-latency federated learning and blockchain for edge association in digital twin empowered 6G networks," *IEEE Trans. Ind. Informat.*, vol. 17, no. 7, pp. 5098–5107, 2020.
- [10] W. Shi, S. Zhou, Z. Niu, M. Jiang, and L. Geng, "Joint device scheduling and resource allocation for latency constrained wireless federated learning," *IEEE Trans. Wireless Commun.*, vol. 20, no. 1, pp. 453–467, 2020.
- [11] P. Han, S. Wang, and K. K. Leung, "Adaptive gradient sparsification for efficient federated learning: An online learning approach," in *IEEE International Conference on Distributed Computing Systems (ICDCS)*, 2020, pp. 300–310.
- [12] C. Hardy, E. Le Merrer, and B. Sericola, "Distributed deep learning on edge-devices in the parameter server model," in *Workshop on Decentralized Machine Learning, Optimization and Privacy*, 2017.
- [13] A. F. Aji and K. Heafield, "Sparse communication for distributed gradient descent," arXiv preprint arXiv:1704.05021, 2017.
- [14] D. Alistarh et al., "The convergence of sparsified gradient methods," Advances in Neural Information Processing Systems, vol. 31, 2018.
- [15] N. Shlezinger, M. Chen, Y. C. Eldar, H. V. Poor, and S. Cui, "UVeQFed: Universal vector quantization for federated learning," *IEEE Trans. Signal Process.*, vol. 69, pp. 500–514, 2020.
- [16] D. Alistarh, D. Grubic, J. Li, R. Tomioka, and M. Vojnovic, "QSGD: Communication-efficient SGD via gradient quantization and encoding," *Advances in Neural Information Processing Systems*, vol. 30, pp. 1709– 1720, 2017.
- [17] A. Reisizadeh, A. Mokhtari, H. Hassani, A. Jadbabaie, and R. Pedarsani, "Fedpaq: A communication-efficient federated learning method with periodic averaging and quantization," in *International Conference on Artificial Intelligence and Statistics*. PMLR, 2020, pp. 2021–2031.
- [18] N. Lang, N. Shlezinger, R. G. D'Oliveira, and S. E. Rouayheb, "Compressed private aggregation for scalable and robust federated learning over massive networks," arXiv preprint arXiv:2308.00540, 2023.
- [19] E. Diao, J. Ding, and V. Tarokh, "Heterofl: Computation and communication efficient federated learning for heterogeneous clients," in *International Conference on Learning Representations*, 2020.
- [20] K. Pfeiffer, M. Rapp, R. Khalili, and J. Henkel, "Federated learning for computationally-constrained heterogeneous devices: A survey," ACM Computing Surveys, 2023.
- [21] S. Vahidian, S. Kadaveru, W. Baek, W. Wang, V. Kungurtsev, C. Chen, M. Shah, and B. Lin, "When do curricula work in federated learning?" in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 5084–5094.
- [22] T. Nishio and R. Yonetani, "Client selection for federated learning with heterogeneous resources in mobile edge," in *IEEE international* conference on communications (ICC), 2019.
- [23] A. Reisizadeh, I. Tziotis, H. Hassani, A. Mokhtari, and R. Pedarsani, "Straggler-resilient federated learning: Leveraging the interplay between statistical accuracy and system heterogeneity," *IEEE J. Sel. Areas Inf. Theory*, vol. 3, no. 2, pp. 197–205, 2022.
- [24] R. Schlegel, S. Kumar, E. Rosnes, and A. G. i Amat, "CodedPaddedFL and CodedSecAgg: Straggler mitigation and secure aggregation in federated learning," *IEEE Trans. Commun.*, 2023.
- [25] I. Wang, P. J. Nair, and D. Mahajan, "Fluid: Mitigating stragglers in federated learning using invariant dropout," in *Thirty-seventh Conference* on Neural Information Processing Systems, 2023.
- [26] T. Li, A. K. Sahu, M. Zaheer, M. Sanjabi, A. Talwalkar, and V. Smith, "Federated optimization in heterogeneous networks," *Machine Learning* and Systems (MLSys), vol. 2, pp. 429–450, 2020.
- [27] J. Park, D.-J. Han, M. Choi, and J. Moon, "Sageflow: Robust federated learning against both stragglers and adversaries," *Advances in neural information processing systems*, vol. 34, pp. 840–851, 2021.
- [28] Y. Chen, Y. Ning, M. Slawski, and H. Rangwala, "Asynchronous online federated learning for edge devices with non-iid data," in *IEEE International Conference on Big Data*, 2020, pp. 15–24.
- [29] K. Bonawitz, H. Eichner, W. Grieskamp, D. Huba, A. Ingerman, V. Ivanov, C. Kiddon, J. Konečný, S. Mazzocchi, B. McMahan *et al.*, "Towards federated learning at scale: System design," *Machine Learning* and Systems (MLSys), vol. 1, pp. 374–388, 2019.
- [30] X. Li, K. Huang, W. Yang, S. Wang, and Z. Zhang, "On the convergence of fedavg on non-iid data," in *International Conference on Learning Representations*, 2019.

- [31] M. Chen, N. Shlezinger, H. V. Poor, Y. C. Eldar, and S. Cui, "Communication-efficient federated learning," *Proceedings of the National Academy of Sciences*, vol. 118, no. 17, 2021.
- [32] H. Esfahanizadeh, A. Cohen, and M. Médard, "Stream iterative distributed coded computing for learning applications in heterogeneous systems," in *IEEE INFOCOM 2022-IEEE Conference on Computer Communications*. IEEE, 2022, pp. 230–239.
- [33] J. Nguyen, K. Malik, H. Zhan, A. Yousefpour, M. Rabbat, M. Malek, and D. Huba, "Federated learning with buffered asynchronous aggregation," in *International Conference on Artificial Intelligence and Statistics*. PMLR, 2022, pp. 3581–3607.
- [34] T. Zhang, L. Gao, S. Lee, M. Zhang, and S. Avestimehr, "Timelyfl: Heterogeneity-aware asynchronous federated learning with adaptive partial training," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 5063–5072.
- [35] T. Ortega and H. Jafarkhani, "Asynchronous federated learning with bidirectional quantized communications and buffered aggregation," in International Conference on Machine Learning (ICML), Workshop on Federated Learning and Analytics, 2023.
- [36] D. E. Rumelhart, G. E. Hinton, and R. J. Williams, "Learning internal representations by error propagation, parallel distributed processing, explorations in the microstructure of cognition, ed. de rumelhart and j. mcclelland. vol. 1. 1986," *Biometrika*, vol. 71, pp. 599–607, 1986.
- [37] S. Zhang, J. Li, L. Shi, M. Ding, D. C. Nguyen, W. Tan, J. Weng, and Z. Han, "Federated learning in intelligent transportation systems: Recent applications and open problems," *IEEE Transactions on Intelligent Transportation Systems*, 2023.
- [38] S. U. Stich, "Local SGD converges fast and communicates little," in *International Conference on Learning Representations*, 2018.
- [39] N. Lang, E. Sofer, T. Shaked, and N. Shlezinger, "Joint privacy enhancement and quantization in federated learning," *IEEE Trans. Signal Process.*, vol. 71, pp. 295–310, 2023.
- [40] F. Edition, A. Papoulis, and S. U. Pillai, *Probability, random variables, and stochastic processes*. McGraw-Hill Europe: New York, NY, USA, 2002.
- [41] X. Glorot and Y. Bengio, "Understanding the difficulty of training deep feedforward neural networks," in *Proceedings of the thirteenth international conference on artificial intelligence and statistics*. JMLR Workshop and Conference Proceedings, 2010, pp. 249–256.
- [42] S. Zagoruyko and N. Komodakis, "Wide residual networks," arXiv preprint arXiv:1605.07146, 2016.
- [43] P. Chang, G. Duràn-Martín, A. Y. Shestopaloff, M. Jones, and K. Murphy, "Low-rank extended Kalman filtering for online learning of neural networks from streaming data," arXiv preprint arXiv:2305.19535, 2023.
- [44] S. Liu, P.-Y. Chen, B. Kailkhura, G. Zhang, A. O. Hero III, and P. K. Varshney, "A primer on zeroth-order optimization in signal processing and machine learning: Principals, recent advances, and applications," *IEEE Signal Process. Mag.*, vol. 37, no. 5, pp. 43–54, 2020.
- [45] L. Deng, "The MNIST database of handwritten digit images for machine learning research," *IEEE Signal Process. Mag.*, vol. 29, no. 6, pp. 141– 142, 2012.
- [46] A. Krizhevsky, V. Nair, and G. Hinton, "Cifar-10 (canadian institute for advanced research)." [Online]. Available: http://www.cs.toronto.edu/ ~kriz/cifar.html
- [47] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," arXiv preprint arXiv:1409.1556, 2014.
- [48] G. An, "The effects of adding noise during backpropagation training on a generalization performance," *Neural computation*, vol. 8, no. 3, pp. 643–674, 1996.
- [49] T. Sery, N. Shlezinger, K. Cohen, and Y. C. Eldar, "Over-the-air federated learning from heterogeneous data," *IEEE Trans. Signal Process.*, vol. 69, pp. 3796–3811, 2021.