

Decentralized Event-Triggered Federated Learning with Heterogeneous Communication Thresholds

Shahryar Zehtabi, Seyyedali Hosseinalipour, Christopher G. Brinton

Abstract—A recent emphasis of distributed learning research has been on federated learning (FL), in which model training is conducted by the data-collecting devices. Existing research on FL has mostly focused on a star topology learning architecture with synchronized (time-triggered) model training rounds, where the local models of the devices are periodically aggregated by a centralized coordinating node. However, in many settings, such a coordinating node may not exist, motivating efforts to fully decentralize FL. In this work, we propose a novel methodology for distributed model aggregations via asynchronous, event-triggered consensus iterations over the network graph topology. We consider heterogeneous communication event thresholds at each device that weigh the change in local model parameters against the available local resources in deciding the benefit of aggregations at each iteration. Through theoretical analysis, we demonstrate that our methodology achieves asymptotic convergence to the globally optimal learning model under standard assumptions in distributed learning and graph consensus literature, and without restrictive connectivity requirements on the underlying topology. Subsequent numerical results demonstrate that our methodology obtains substantial improvements in communication requirements compared with FL baselines.

I. INTRODUCTION

Federated learning (FL) has emerged as a popular technique for distributing machine learning model training across network devices [1]. In the conventional FL architecture, a set of devices are connected to a central coordinating node (e.g., an edge server) in a star topology configuration. Devices conduct local model updates based on their individual datasets, and the coordinator periodically aggregates these into a global model, synchronizing the devices to begin the next round of training. Several works in the past few years have built functionality into this architecture to manage different types of network heterogeneity, including varying communication and computation abilities of devices and statistical properties of local datasets [2], [3].

However, access to a central coordinating node is not always feasible/desirable. For instance, ad-hoc wireless networks serve as an efficient alternative for communication among devices in settings where device-to-server connectivity is energy intensive or unavailable [4]. The proliferation of such settings motivates consideration of *fully-decentralized FL*, where the model aggregation step, in addition to the data processing step, is distributed across devices. In this paper, we propose a cooperative learning approach for achieving

this via consensus iterations over the available distributed graph topology, and analyze its convergence characteristics.

The central coordinator in FL is also typically employed for timing synchronization, i.e., determining the time between global aggregations. To overcome this, we consider an *asynchronous, event-triggered communication framework* for distributed model consensus. Event-triggered communications can offer several benefits in this context. For one, the amount of redundant communications can be reduced by defining event triggering conditions based on the significance of each device’s model update. Also, removing the assumption of devices communicating at every iteration opens the possibility of alleviating straggler issues [5]. Third, we can improve computational efficiency at each device by limiting aggregations to only when new parameters are received.

A. Related Work

1) *Consensus-based distributed optimization*: There is a rich literature on distributed optimization over graphs via consensus algorithms, e.g., [6]–[11]. For connected, undirected graph topologies, symmetric and doubly-stochastic transition matrices can be constructed for consensus iterations. In typical approaches [6]–[8], each device maintains a local gradient of the target system objective (e.g., error minimization), with the consensus matrices designed to satisfy additional convergence criteria outlined in [12], [13]. More recently, gradient tracking optimization techniques have been developed where the global gradient is simultaneously learned alongside local parameters [9].

In this work, our focus is on decentralized FL, which adds two unique aspects to the distributed optimization problem. First, the local data distributions across devices for machine learning tasks are in general not independent and identically distributed (non-i.i.d.), which can have significant impacts on convergence [5]. Second, we consider the realistic scenario in which the devices have heterogeneous resources [2].

2) *Resource-efficient federated learning*: Several recent works in FL have investigated techniques for improving communication and computation efficiency. A popular line of research has aimed to adaptively control the FL process based on device capabilities, e.g., [14]–[18]. [17] studies FL convergence under a total network resource budget, in which the server adapts the frequency of global aggregation iterations. Others [14], [15], [18] have considered FL under partial device participation, where the communication and processing capabilities of devices are taken into account when assessing which clients will participate in each training round. [16] remove the necessity that every local client needs

S. Zehtabi, S. Hosseinalipour and C. Brinton are with the School of Electrical and Computer Engineering, Purdue University, West Lafayette, IN, 47906 USA email: {szehtabi, hosseina, cgb}@purdue.edu

This work was supported in part by the Office of Naval Research under grant N00014-21-1-2472

to share the same global model as the server, allowing weaker clients to take smaller subsets of the model to optimize.

Different from these works, we focus on novel learning topologies for decentralized FL. In this respect, a few recent works [19]–[22] have proposed peer-to-peer (P2P) communication approaches for collaborative learning over local device topologies. [21], [22] investigated a semi-decentralized FL methodology across hierarchical networks, where local model aggregations are conducted via P2P-based cooperative consensus formation to reduce the frequency of global aggregations by the coordinating node. In our work, we consider the fully decentralized setting, where a central node is not available, as in [19], [20]: alongside local model updates, devices conduct consensus iterations with their neighbors in order to gradually minimize the global machine learning loss in a distributed manner. Different from [19], [20], our methodology incorporates asynchronous event-triggered communications among devices, where local resource levels are factored in to the event thresholds to account for device heterogeneity. We will see that this approach leads to substantial improvements in model convergence time compared with non-heterogeneous thresholding.

B. Outline and Summary of Contributions

- We develop a novel methodology for fully decentralizing FL, with model aggregations occurring via cooperative model consensus iterations (Sec. II). In our methodology, communications are asynchronous and event-driven. With event thresholds defined to incorporate local model evolution and resource availability, our methodology adapts to device communication and processing limitations in heterogeneous networks.
- We provide a convergence analysis of our methodology, which shows that each device arrives at the globally optimal learning model asymptotically under standard assumptions for distributed learning (Sec. III). This result is obtained without overly restrictive connectivity assumptions on the underlying communication graph. Our analysis also leads to guardrails for the event-triggering conditions to ensure convergence.
- We conduct numerical experiments comparing our methodology to baselines in decentralized FL and a randomized gossip algorithm on a real-world machine learning dataset (Sec. V). We show that our method is able to reduce model training communication time substantially compared to FL baselines. Additionally, we find that the convergence rate of our method scales well with consensus graph connectivity.

II. METHODOLOGY AND ALGORITHM

In this section, we develop our methodology for decentralized FL with event-triggered communications. After discussing preliminaries of the learning model in FL (Sec. II-A), we present our cooperative consensus algorithm for distributed model aggregations (Sec. II-B). Finally, we remark on hyperparameters introduced in our algorithm (Sec. IV).

A. Device and Learning Model

We consider a system of m devices/nodes, collected via the set \mathcal{M} , which are engaged in distributed training of a machine learning model. Under the FL framework, each device $i \in \mathcal{M}$ trains a local model \mathbf{w}_i using its own generated dataset \mathcal{D}_i . Each data point $\xi \triangleq (\mathbf{x}_\xi, y_\xi) \in \mathcal{D}_i$ consists of a feature vector \mathbf{x}_ξ and (optionally, in the case of supervised learning) a target label y_ξ . The performance of the local model is measured via the local loss $F_i(\cdot)$:

$$F_i(\mathbf{w}) = \sum_{\xi \in \mathcal{D}_i} \ell_\xi(\mathbf{w}), \quad (1)$$

where $\ell_\xi(\mathbf{w}_i)$ is the loss of the model on datapoint ξ (e.g., squared prediction error) under parameter realization $\mathbf{w} \in \mathbb{R}^n$, with n denoting the dimension of the target model. The global loss is defined in terms of these local losses as

$$F(\mathbf{w}) = \frac{1}{|\mathcal{M}|} \sum_{i \in \mathcal{M}} F_i(\mathbf{w}). \quad (2)$$

The goal of the training process is to find an optimal parameter vector \mathbf{w}^* which minimizes the global loss function, i.e., $\mathbf{w}^* = \arg \min_{\mathbf{w} \in \mathbb{R}^n} F(\mathbf{w})$. In the distributed setting, we desire $\mathbf{w}_1 = \dots = \mathbf{w}_m = \mathbf{w}^*$, which requires a synchronization mechanism. In conventional FL, synchronization is conducted periodically by a central coordinator globally aggregating the local models. However, in this work, we are interested in settings where no such central node exists. Thus, alongside using optimization techniques to minimize local loss functions, we must desire a technique to reach consensus over the parameters in our decentralized FL scheme.

To accomplish this, we propose event-triggered FL with heterogeneous communication thresholds (EF-HC). In EF-HC, devices conduct peer-to-peer (P2P) communications during the model training period to synchronize their locally trained models and avoid overfitting to their local datasets. The overall EF-HC algorithm is given in Alg. 1. Two model parameter vectors are kept at each device i : (i) its instantaneous *main* model parameters \mathbf{w}_i , and (ii) the *auxiliary* model parameters $\widehat{\mathbf{w}}_i$, which is the outdated version of its main parameters that had been broadcast to the neighbors. Decentralized ML is conducted over the (time-varying, undirected) device graph through a sequence of four events detailed in Sec. II-B. Although in our distributed setup there is no physical notion of a global iteration, we introduce the iteration variable k for analysis purposes.

B. Model Updating and Event-Triggered Communications

We consider the physical network graph $\mathcal{G}^{(k)} = (\mathcal{M}, \mathcal{E}^{(k)})$ among devices, where $\mathcal{E}^{(k)}$ is the set of edges available at iteration k in the underlying time-varying communication graph. We assume that link availability varies over time according to the underlying communication protocol in place [5]. In each iteration, some of these edges are employed for transmission/reception of model parameters between devices. To represent this process, we define the information flow graph $\mathcal{G}'^{(k)} = (\mathcal{M}, \mathcal{E}'^{(k)})$, which is a subgraph of $\mathcal{G}^{(k)}$. $\mathcal{E}'^{(k)}$ only contains those links in $\mathcal{E}^{(k)}$ that are being used at

iteration k to exchange parameters. Based on this, we denote the neighbors of device i at iteration k as $\mathcal{N}_i^{(k)} = \{j : (i, j) \in \mathcal{E}^{(k)}, j \in \mathcal{M}\}$, with node degree $d_i^{(k)} = |\mathcal{N}_i^{(k)}|$. We also denote neighbors of i which are directly communicating with it at iteration k as $\mathcal{N}'_i{}^{(k)} = \{j : (i, j) \in \mathcal{E}'^{(k)}, j \in \mathcal{M}\}$. Additionally, the aggregation weight associated with the link $(i, j) \in \mathcal{E}^{(k)}$ and $(i, j) \in \mathcal{E}'^{(k)}$ are defined as $\beta_{ij}^{(k)}$ and $p_{ij}^{(k)}$ respectively, with $p_{ij}^{(k)} = \beta_{ij}^{(k)}$ if the link (i, j) is used for aggregation at iteration k , and $p_{ij}^{(k)} = 0$ otherwise.

In EF-HC, there are four types of communication events:

Event 1: Neighbor connection. The first event (lines 2-8 of Alg. 1) is triggered at device i if new devices connect to it or existing devices disconnect from it due to the time-varying nature of the graph. In this event, model parameters $\mathbf{w}_i^{(k)}$ and the degree of device i at that time $d_i^{(k)}$ are exchanged with this new neighbor. Consequently, this results in an aggregation event (Event 3) at both devices.

Event 2: Broadcast. Second, if the normalized difference between $\mathbf{w}_i^{(k)}$ and $\widehat{\mathbf{w}}_i^{(k)}$ at device i is greater than a *threshold* value $r\rho_i\gamma^{(k)}$, i.e., $(\frac{1}{n})^{\frac{1}{q}} \|\mathbf{w}_i^{(k)} - \widehat{\mathbf{w}}_i^{(k)}\|_q \geq r\rho_i\gamma^{(k)}$, then a broadcast event (lines 9-13 of Alg. 1) is triggered at that device. In other words, communication at a device is triggered once the instantaneous local model is sufficiently different than the outdated local model. When this event triggers, device i broadcasts its parameters $\mathbf{w}_i^{(k)}$ and its instantaneous degree $d_i^{(k)}$ to all of its neighbors and receives the same information from them.

The threshold $r\rho_i\gamma^{(k)}$ is treated as heterogeneous across devices $i \in \mathcal{M}$, to assess whether the gain from a consensus iteration on the instantaneous main models at the devices will be worth the induced network resource utilization. Specifically: (i) $r > 0$ is a scaling hyperparameter value; (ii) $\gamma^{(k)} > 0$ is a decaying factor that accounts for smaller expected variations in the local models over time; and (iii) ρ_i quantifies the resource availability of device i . Developing the threshold measure $(\frac{1}{n})^{\frac{1}{q}} \|\mathbf{w}_i^{(k)} - \widehat{\mathbf{w}}_i^{(k)}\|_q$ and the condition $r\rho_i\gamma^{(k)}$ is one of the contributions of this paper relative to existing event-triggered schemes [23]. For example, in a bandwidth limited environment, the transmission delay of model transfer will be inversely proportional to the bandwidth among two devices. Thus, to decrease the latency of model training, ρ_i can be defined inversely proportional to the bandwidth, promoting lower frequency of communications at the devices with less available bandwidth. In EF-HC, we set $\rho_i \propto \frac{1}{b_i}$, where b_i is the average bandwidth on outgoing links of device i . Further details on choosing the broadcast threshold are given in Sec. IV.

Event 3: Aggregation. Following a broadcast event (Event 2) or a neighbor connection event (Event 1) at device i , an aggregation event (lines 14-16 of Alg. 1) is triggered at device i and all of its neighbors. This aggregation is carried out through a distributed weighted averaging consensus method [13] as $\mathbf{w}_i^{(k+1)} = \mathbf{w}_i^{(k)} + \sum_{j \in \mathcal{N}'_i{}^{(k)}} \beta_{ij}^{(k)} (\mathbf{w}_j^{(k)} - \mathbf{w}_i^{(k)})$, where $\beta_{ij}^{(k)}$ is the aggregation weight that device i will assign to parameters received from

Algorithm 1 EF-HC procedure for device i .

Input: K, q

Initialize $k = 0, \mathbf{w}_i^{(0)} = \widehat{\mathbf{w}}_i^{(0)}$

1: **while** $k \leq K$ **do**

▷ **Event 1.** Neighbor Connection Event

2: **if** device j is connected to device i **then**

3: device i appends device j to its list of neighbors

4: device i sends $\mathbf{w}_i^{(k)}$ and $d_i^{(k)}$ to device j

5: device i receives $\mathbf{w}_j^{(k)}$ and $d_j^{(k)}$ from device j

6: **else if** device j is disconnected from device i **then**

7: device i removes device j from its list of neighbors

8: **end if**

▷ **Event 2.** Broadcast Event

9: **if** $(\frac{1}{n})^{\frac{1}{q}} \|\mathbf{w}_i^{(k)} - \widehat{\mathbf{w}}_i^{(k)}\|_q \geq r\rho_i\gamma^{(k)}$ **then**

10: device i broadcasts $\mathbf{w}_i^{(k)}, d_i^{(k)}$ to all neighbors $j \in \mathcal{N}_i^{(k)}$

11: device i receives $\mathbf{w}_j^{(k)}, d_j^{(k)}$ from all neighbors $j \in \mathcal{N}_i^{(k)}$

12: $\widehat{\mathbf{w}}_i^{(k+1)} = \mathbf{w}_i^{(k)}$

13: **end if**

▷ **Event 3.** Aggregation Event

14: **if** updated parameters $\mathbf{w}_j^{(k)}$ and $d_j^{(k)}$ received from neighbor j **then**

15: $\mathbf{w}_i^{(k+1)} = \mathbf{w}_i^{(k)} + \sum_{j \in \mathcal{N}'_i{}^{(k)}} \beta_{ij}^{(k)} (\mathbf{w}_j^{(k)} - \mathbf{w}_i^{(k)})$

16: **end if**

▷ **Event 4.** Gradient Descent Event

17: device i conducts SGD iteration $\mathbf{w}_i^{(k+1)} = \mathbf{w}_i^{(k)} -$

$\alpha^{(k)} \mathbf{g}_i(\mathbf{w}_i^{(k)})$

18: $k \leftarrow k + 1$

19: **end while**

device j at iteration k . The aggregation weights $\{\beta_{ij}^{(k)}\}$ for graph $\mathcal{G}^{(k)}$ can be selected based on degree information of the neighbors, as will be discussed in Sec. III-B.

Event 4: Gradient descent. Each device i conducts stochastic gradient descent (SGD) iterations for local model training. Formally, device i obtains $\mathbf{w}_i^{(k+1)} = \mathbf{w}_i^{(k)} - \alpha^{(k)} \mathbf{g}_i(\mathbf{w}_i^{(k)})$, where $\alpha^{(k)}$ is the learning rate and $\mathbf{g}_i(\mathbf{w}_i^{(k)})$ is the stochastic gradient approximation defined as $\mathbf{g}_i(\mathbf{w}_i^{(k)}) = \frac{1}{|\mathcal{S}_i^{(k)}|} \sum_{\xi \in \mathcal{S}_i^{(k)}} \nabla \ell_{\xi}(\mathbf{w}_i^{(k)})$. Here, $\mathcal{S}_i^{(k)}$ denotes the set of data points (mini-batch) used to compute the gradient, chosen uniformly at random from the local dataset.

III. CONVERGENCE ANALYSIS

In this section, we first present our main theoretical result in this paper (Sec. III-A). We then enumerate and discuss the assumptions needed to obtain the main result (Sec. III-B).

A. Main Convergence Result

We first obtain the convergence characteristics of EF-HC. We reveal that (a) all devices reach consensus asymptotically, i.e., each device i 's model $\mathbf{w}_i^{(k)}$ converges to $\bar{\mathbf{w}}^{(k)} = \frac{1}{m} \sum_{i=1}^m \mathbf{w}_i^{(k)}$ as $k \rightarrow \infty$, and (b) the final model across the devices (i.e., $\bar{\mathbf{w}}^{(k)}, k \rightarrow \infty$) minimizes the global loss.

Theorem 1. *Under the standard distributed learning assumptions in Sec. III-B, model training under EF-HC satisfies the following convergence behaviors:*

(a) $\lim_{k \rightarrow \infty} \|\mathbf{w}_i^{(k)} - \bar{\mathbf{w}}^{(k)}\|_2 = 0$ for all i , where $\bar{\mathbf{w}}^{(k)} = \frac{1}{m} \sum_{i=1}^m \mathbf{w}_i^{(k)}$, and

(b) $\lim_{k \rightarrow \infty} F(\bar{\mathbf{w}}^{(k)}) - F^* = 0$.

Proof. See Appendix A. \blacksquare

B. Assumptions for Theorem 1

Assumption 1 (Simultaneous information exchange). *The devices exchange information simultaneously: if device j communicates with device i at some time, device i also communicates with device j at that same time.*

Assumption 2 (Transition weights). *Let $\{p_{ij}^{(k)}\}$ be the set of aggregation weights in the information graph $\mathcal{G}'(k)$. $p_{ij}^{(k)}$ is the transition weight that device i utilizes to aggregate device j 's parameters at iteration k :*

$$p_{ij}^{(k)} = \begin{cases} \beta_{ij}^{(k)} v_{ij}^{(k)} & i \neq j \\ 1 - \sum_{j=1}^m \beta_{ij}^{(k)} v_{ij}^{(k)} & i = j \end{cases}, \quad (3)$$

where $v_i^{(k)}$ indicates whether a broadcast event has occurred at device i at iteration k :

$$\begin{aligned} v_i^{(k)} &= \begin{cases} 1 & \left(\frac{1}{n}\right)^{\frac{1}{q}} \|\mathbf{e}_i^{(k)}\|_q > r\rho_i\gamma^{(k)} \\ 0 & \text{o.w.} \end{cases}, \\ \mathbf{e}_i^{(k)} &= \mathbf{w}_i^{(k)} - \widehat{\mathbf{w}}_i^{(k)}, \quad \rho_i = \frac{1}{b_i}, \\ v_{ij}^{(k)} &= \begin{cases} \max\{v_i^{(k)}, v_j^{(k)}\} & j \in \mathcal{N}_i^{(k)} \\ 0 & \text{o.w.} \end{cases}. \end{aligned} \quad (4)$$

The following conditions must hold:

- (a) (Non-negative weights) There exists a scalar η , $0 < \eta < 1$, such that $\forall i \in \mathcal{M}$, we have
 - (i) $p_{ii}^{(k)} \geq \eta$ and $p_{ij}^{(k)} \geq \eta$ for all $k \geq 0$ and all neighbor devices $j \in \mathcal{N}_i^{(k)}$.
 - (ii) $p_{ij}^{(k)} = 0$, if $j \notin \mathcal{N}_i^{(k)}$.
- (b) (Doubly-stochastic weights) The rows and columns of matrix $\mathbf{P}^{(k)} = [p_{ij}^{(k)}]$ are both stochastic, i.e., $\sum_{j=1}^m p_{ij}^{(k)} = 1$, $\forall i$, and $\sum_{i=1}^m p_{ij}^{(k)} = 1$, $\forall j$.
- (c) (Symmetric weights) $p_{ij}^{(k)} = p_{ji}^{(k)}$, $\forall i, k$ and $p_{ii}^{(k)} = 1 - \sum_{j \neq i} p_{ij}^{(k)}$.

Considering the conditions mentioned in Assumption 2, and the definition of $p_{ij}^{(k)}$ in (3), a choice of parameters $\beta_{ij}^{(k)}$ that satisfy these assumptions are as follows:

$$\beta_{ij}^{(k)} = \min \left\{ \frac{1}{1 + d_i^{(k)}}, \frac{1}{1 + d_j^{(k)}} \right\}, \quad (5)$$

which is inspired by the Metropolis-Hastings algorithm [12]. Note that $p_{ij}^{(k)}$ also depends on $v_{ij}^{(k)}$, which was defined in (4).

Assumption 3 (Convexity). (a) *The local objective function at each device i , i.e., F_i , is convex:*

$$F_i(\mathbf{w}') \geq F_i(\mathbf{w}) + \nabla F_i(\mathbf{w})^\top (\mathbf{w}' - \mathbf{w}),$$

$$\forall (\mathbf{w}', \mathbf{w}) \in \mathbb{R}^n \times \mathbb{R}^n.$$

- (b) *The global objective function $F(\mathbf{w}) = \frac{1}{m} \sum_{i=1}^m F_i(\mathbf{w})$ is convex, and thus has a non-empty minimizer set*

denoted by $\mathbf{W}^* = \text{Arg min}_{\mathbf{w} \in \mathbb{R}^n} F(\mathbf{w})$, such that $F^* = F(\mathbf{w}^*)$ for any $\mathbf{w}^* \in \mathbf{W}^*$.

Assumption 4 (Bounded gradients). *The gradient of each loss function F_i is bounded, i.e., there exists a scalar $L_i > 0$ such that $\forall i \in \mathcal{M}$, $\mathbf{w} \in \mathbb{R}^n$,*

$$\|\nabla F_i(\mathbf{w})\|_2 \leq L_i \leq L,$$

where $L = \max_{i \in \mathcal{M}} L_i$. We define L_∞ as the bound for the infinity norm of all F_i 's, i.e., $\|\nabla F_i(\mathbf{w})\|_\infty \leq L_\infty$, $\forall i$.

Assumption 5 (Step sizes). *All devices use the same step size for model training, which is diminishing over time and satisfies the following conditions:*

$$\lim_{k \rightarrow \infty} \alpha^{(k)} = 0, \quad \sum_{k=0}^{\infty} \alpha^{(k)} = \infty, \quad \sum_{k=0}^{\infty} \left(\alpha^{(k)}\right)^2 < \infty.$$

In particular, setting $\alpha^{(k)} = \frac{a}{(b+k)^c}$ meets the criteria of the above assumption if $c \in (0.5, 1]$.

The previous assumptions are common in literature [17], [21]. In the next assumption, we introduce a relaxed version of graph connectivity requirements relative to existing work in distributed learning, which underscores the difference of our decentralized event-triggered FL method compared with traditional distributed optimization algorithms.

Assumption 6 (Network graph connectivity). *The physical network graph $\mathcal{G}^{(k)} = (\mathcal{M}, \mathcal{E}^{(k)})$ satisfies the following:*

- (a) *There exists an integer $B_1 \geq 1$ such that the graph union of $\mathcal{G}^{(k)}$ from any arbitrary iteration k to $k + B_1 - 1$, i.e., $\mathcal{G}^{(k:k+B_1-1)} = (\mathcal{M}, \cup_{s=0}^{B_1-1} \mathcal{E}^{(k+s)})$, is connected for any $k \geq 0$.*
- (b) *There exists an integer $B_2 \geq 1$ such that for every device i , triggering conditions for the broadcasting event occurs at least once every B_2 consecutive iterations $\forall k \geq 0$. This is equivalent to the following condition:*

$$\exists B_2 \geq 1, \forall i : \max\{v_i^{(k)}, v_i^{(k+1)}, \dots, v_i^{(k+B_2-1)}\} = 1.$$

Together, (a) and (b) imply that each device i broadcasts its information to its neighboring devices at least once every B consecutive iterations, where $B = (l+2)B_1$ in which $lB_1 < B_2 \leq (l+1)B_1$.¹ Hence, the information flow graph $\mathcal{G}'(k)$ is B -connected, i.e., $\mathcal{G}'(k:k+B-1) = (\mathcal{M}, \cup_{s=0}^{B-1} \mathcal{E}'(k+s))$ is connected, for any $k \geq 0$. It is important to note that we use B only for convergence analysis, and it can have any arbitrarily large integer value. Therefore, we are not making strict connectivity assumptions on the underlying graph.

IV. REMARKS ON HYPERPARAMETERS

We make a few remarks on the hyperparameters used in Alg. 1. Remark 1 elaborates on the choice of q when calculating $\|\mathbf{w}_i^{(k)} - \widehat{\mathbf{w}}_i^{(k)}\|_q$, Remark 2 discusses the choice of r in the threshold $r\rho_i\gamma^{(k)}$, and Remark 3 elaborates on the

¹Note that in Algorithm 1, once two unconnected devices become connected, they exchange their parameters regardless of triggering conditions.

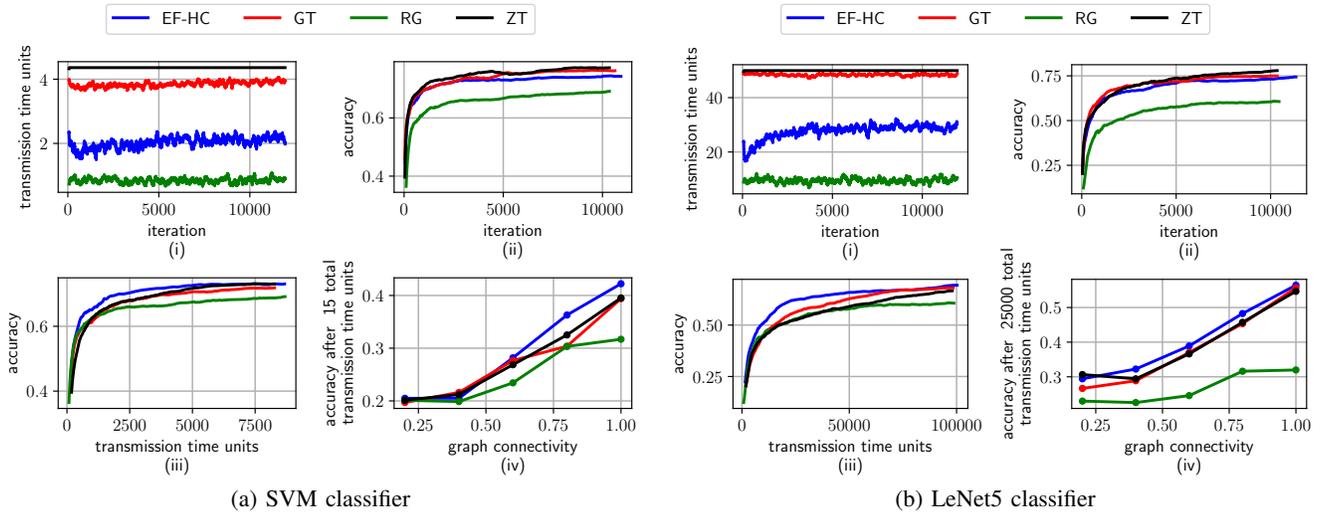


Fig. 1: Performance comparison between our method (EF-HC), global threshold (GT), zero threshold (ZT), and randomized gossip (RG) algorithms. The plots show (i) transmission time per iteration, (ii) accuracy per iteration, (iii) accuracy per transmission time, and (iv) accuracy after a certain number of transmissions with respect to graph connectivity.

threshold decay rate $\gamma^{(k)}$ and the learning rate $\alpha^{(k)}$. More explanations of these remarks are deferred to Appendix B.

Remark 1. Factor $(\frac{1}{n})^{\frac{1}{q}}$ in the event-triggering condition is a normalization factor, making the conditions independent of the model dimension n and the norm $q \geq 1$ used.

Remark 2. The constant r in the threshold of event-triggering condition is a hyperparameter to set the threshold $r\rho_i\gamma^{(k)}$ to a value comparable with $(\frac{1}{n})^{\frac{1}{q}}\|\mathbf{w}_i^{(k)} - \widehat{\mathbf{w}}_i^{(k)}\|_q$. The value of this constant can be chosen as follows:

$$r = \frac{\alpha^{(0)}}{\gamma^{(0)}} \frac{1}{\rho} KL_{\infty} = \frac{1}{\rho} KL_{\infty} \quad \text{if } \alpha^{(0)} = \gamma^{(0)}, \quad (6)$$

in which K is an approximation on the number of local iterations between aggregation events, $\frac{1}{\rho}$ is an approximation of $\frac{1}{m} \sum_{i \in \mathcal{M}} \frac{1}{\rho_i}$, and L_{∞} is an upper bound obtained via the relation $\|\nabla F_i(x)\|_{\infty} \leq L_{\infty}$ for all i from Assumption 4.

Remark 3. To ensure sporadic aggregations at each device in EF-HC, the learning rate $\alpha^{(k)}$ and threshold decay rate $\gamma^{(k)}$ should satisfy the following conditions:

- (a) $\lim_{k \rightarrow \infty} \frac{\gamma^{(k)}}{\alpha^{(k)}} = \Omega$, where Ω is a finite positive constant.
- (b) $\lim_{k \rightarrow \infty} \frac{\gamma^{(k)}/\gamma^{(0)}}{\alpha^{(k)}/\alpha^{(0)}} = 1$, i.e., $\Omega = \frac{\gamma^{(0)}}{\alpha^{(0)}}$.

The conditions in Remark 3 ensure that aggregation events (Event 3) neither cease completely nor occur continuously after a while, but instead are executed only when our proposed triggering condition is met (see Appendix B for details). To satisfy these conditions, first the learning rate $\alpha^{(k)}$ should be chosen to meet the criteria of Assumption 5, and then $\gamma^{(k)}$ should be chosen to satisfy the above conditions. One choice that satisfies these conditions is $\gamma^{(k)} = \alpha^{(k)}$.

V. NUMERICAL RESULTS

We now conduct numerical experiments to validate our methodology. We explain our simulation setup in Sec. V-A and provide the results and discussion in Sec. V-B.

A. Simulation Setup

We evaluate our proposed methodology classification tasks on the Fashion-MNIST image recognition dataset [24]. We employ support vector machine (SVM) and the LeNet5 neural network model as classifiers; SVM satisfies Assumption 3 while LeNet5 (and deep learning models in general) does not.

We consider a network of $m = 10$ devices, where the topology is generated according to a random geometric graph with connectivity 0.4 [21]. To generate non-i.i.d. data distributions across devices, each device only contains samples of Fashion-MNIST from a fraction of the 10 labels. For SVM and LeNet5, we consider 1 and 2 labels/device, respectively.

We set the average link bandwidth to 5000. We introduce a resource heterogeneity measure H , $0 \leq H < 1$, which we use to generate networks with two types of devices: (i) “weak,” which have outgoing links with an average bandwidth of 1000, and (ii) “powerful,” which have an average outgoing link bandwidth of $\frac{5000-1000H}{1-H}$. We set $H = 0.8$ for LeNet5 and $H = 0.4$ for SVM.

In each experiment, the learning rate is selected as $\alpha^{(k)} = \frac{1}{\sqrt{1+k}}$, and threshold decay rate is set to $\gamma^{(k)} = \alpha^{(k)}$, satisfying the conditions in Remark 3. The 2-norm is used for the event-triggering conditions (see Remark 1), and $r = 5000 \times 10^{-2}$ following the guidelines of Remark 2.

At iteration k , we define a resource utilization score as $\frac{1}{m} \sum_{i=1}^m \frac{\sum_{j=1}^m v_{ij}^{(k)}}{d_i^{(k)}} \rho_i n$. The term $\frac{\sum_{j=1}^m v_{ij}^{(k)}}{d_i^{(k)}}$ is the outgoing link utilization, and therefore this score is a weighted average of link utilization, penalizing devices with larger ρ_i . For our proposed method where $\rho_i = \frac{1}{b_i}$, this score is the same as the average transmission time, i.e., $\frac{1}{m} \sum_{i=1}^m \frac{\sum_{j=1}^m v_{ij}^{(k)}}{d_i^{(k)}} \frac{n}{b_i}$.

B. Results and Discussion

We compare the performance of our method EF-HC against three baseline methods: (i) distributed learning with

aggregations at every iteration, i.e., zero thresholds (denoted by ZT), (ii) decentralized event-triggered FL with the same global threshold across all devices (denoted by GT), and (iii) randomized gossip algorithm where each device engages in communication with probability of $\frac{1}{m}$ at each iteration [9] (denoted by RG). The performance of our method against these baselines for each classifier is depicted in Fig. 1.

Figs. 1a-(i) and 1b-(i) show the average transmission time units each algorithm requires per training iteration. As can be observed, $EF-HC$ results in less transmission delay compared to ZT and GT , which helps to resolve the impact of stragglers by not requiring the same amount of communications and aggregations from devices with less available bandwidth. Note that although less transmission delay per iteration is desirable for a decentralized FL algorithm, this runs the risk of degrading the performance of the classification task in terms of accuracy when the data distribution across devices is non-i.i.d. Thus, a good comparison between multiple decentralized algorithms is to consider the accuracy reached per transmission time units. In this regard, although RG achieves less transmission delay per iteration compared to our method, Figs. 1a-(iii) and 1b-(iii) reveal that it achieves substantially lower model performance, indicating that our method strikes an effective balance between these objectives.

The average accuracy of devices per iteration is plotted in Figs. 1a-(ii) and 1b-(ii). These plots are indicative of processing efficiency since they evaluate the accuracy of algorithms per number of gradient descent computations. As expected, the baseline algorithm ZT is able to achieve the highest accuracy per iteration since it does not take resource efficiency into account, and thus sacrifices network resources to reach a better accuracy. In these plots, we show that unlike RG , the performance of our proposed method $EF-HC$ as well as GT do not considerably degrade although they use less communication resources, as will be discussed next.

Figs. 1a-(iii) and 1b-(iii) are perhaps the most critical results, as they assess the accuracy vs. communication time tradeoff. We see that our algorithm $EF-HC$ can achieve a higher accuracy while using less transmission time compared to all the baselines, both for the SVM classifier and LeNet5, i.e., with and without the model convexity assumption from our convergence analysis. These plots reveal that our method can adapt to non-i.i.d data distributions across the devices, which is an important characteristic for FL algorithms [1], and achieve a better accuracy as compared to the baselines given a fixed transmission time, i.e., under a fixed network resource consumption.

Finally, we evaluate the effect of network connectivity on our method and baseline methods in Figs. 1a-(iv) and 1b-(iv)². Since the graphs are generated randomly in our simulations, we have taken the average performance of all four algorithms over 5 Monte Carlo instances to reduce the effect of random initialization on the results. It can be observed that higher network connectivity improves the

convergence speed of our method and most of the baselines, as expected. Importantly, however, we see that our method has the highest improvement per increase in connectivity.

VI. CONCLUSION AND FUTURE WORK

In this paper, we developed a novel methodology for event-triggered FL with heterogeneous communication thresholds ($EF-HC$). $EF-HC$ introduces a scenario where the conventional centralized model aggregations in FL are carried out in a decentralized manner via P2P communications among the devices. To further alleviate the burden of a centralized scheduler and take into account resources heterogeneity across the devices, it considers event-triggered communications with heterogeneous communication thresholds. We conducted a theoretical analysis of $EF-HC$ and demonstrated that model training under $EF-HC$ asymptotically achieves the global optimal model for standard assumptions in distributed learning. Future work can focus on deriving optimal/data-driven algorithms for setting the event-triggering communication conditions under different network settings.

REFERENCES

- [1] P. Kairouz, H. B. McMahan, B. Avent, A. Bellet, M. Bennis, A. N. Bhagoji, K. Bonawitz, Z. Charles, G. Cormode, R. Cummings, *et al.*, "Advances and open problems in federated learning," *Found. Trends® ML*, vol. 14, no. 1–2, pp. 1–210, 2021.
- [2] 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," *Proc. Machine Learn. Sys.*, vol. 1, pp. 374–388, 2019.
- [3] T. Li, A. K. Sahu, A. Talwalkar, and V. Smith, "Federated learning: Challenges, methods, and future directions," *IEEE Signal Proc. Mag.*, vol. 37, no. 3, pp. 50–60, 2020.
- [4] M. Chiang and T. Zhang, "Fog and IoT: An overview of research opportunities," *IEEE IoT J.*, vol. 3, no. 6, pp. 854–864, 2016.
- [5] S. Hosseinalipour, C. G. Brinton, V. Aggarwal, H. Dai, and M. Chiang, "From federated to fog learning: Distributed machine learning over heterogeneous wireless networks," *IEEE Commun. Mag.*, vol. 58, no. 12, pp. 41–47, 2020.
- [6] J. Tsitsiklis, D. Bertsekas, and M. Athans, "Distributed asynchronous deterministic and stochastic gradient optimization algorithms," *IEEE Trans. Auto. Control*, vol. 31, no. 9, pp. 803–812, 1986.
- [7] A. Nedic and A. Ozdaglar, "Distributed subgradient methods for multi-agent optimization," *IEEE Trans. Auto. Control*, vol. 54, no. 1, pp. 48–61, 2009.
- [8] A. Nedic, A. Ozdaglar, and P. A. Parrilo, "Constrained consensus and optimization in multi-agent networks," *IEEE Trans. Auto. Control*, vol. 55, no. 4, pp. 922–938, 2010.
- [9] S. Pu and A. Nedić, "Distributed stochastic gradient tracking methods," *Math. Program.*, vol. 187, no. 1, pp. 409–457, 2021.
- [10] A. Nedić and A. Olshevsky, "Distributed optimization over time-varying directed graphs," *IEEE Trans. Auto. Control*, vol. 60, no. 3, pp. 601–615, 2014.
- [11] R. Xin and U. A. Khan, "A linear algorithm for optimization over directed graphs with geometric convergence," *IEEE Control Sys. Lett.*, vol. 2, no. 3, pp. 315–320, 2018.
- [12] S. Boyd, P. Diaconis, and L. Xiao, "Fastest mixing markov chain on a graph," *SIAM review*, vol. 46, no. 4, pp. 667–689, 2004.
- [13] L. Xiao and S. Boyd, "Fast linear iterations for distributed averaging," *Sys. & Control Lett.*, vol. 53, no. 1, pp. 65–78, 2004.
- [14] T. Nishio and R. Yonetani, "Client selection for federated learning with heterogeneous resources in mobile edge," in *IEEE Int. Conf. Commun. (ICC)*, pp. 1–7, 2019.
- [15] H. T. Nguyen, V. Schwag, S. Hosseinalipour, C. G. Brinton, M. Chiang, and H. V. Poor, "Fast-convergent federated learning," *IEEE J. Sel. Areas Commun.*, vol. 39, no. 1, pp. 201–218, 2020.
- [16] E. Diao, J. Ding, and V. Tarokh, "Heteroff: Computation and communication efficient federated learning for heterogeneous clients," in *Int. Conf. Learn. Represent. (ICLR)*, 2020.

²For the LeNet5 classifier, we change the simulation setup and set $r = 5000 \times 10^{-3}$, and let the devices to have samples from only 1 labels/device.

- [17] S. Wang, T. Tuor, T. Salonidis, K. K. Leung, C. Makaya, T. He, and K. Chan, "Adaptive federated learning in resource constrained edge computing systems," *IEEE J. Sel. Areas in Commun.*, vol. 37, no. 6, pp. 1205–1221, 2019.
- [18] X. Gu, K. Huang, J. Zhang, and L. Huang, "Fast federated learning in the presence of arbitrary device unavailability," *Advances Neur. Info. Process. Sys. (NeurIPS)*, vol. 34, 2021.
- [19] S. Savazzi, M. Nicoli, and V. Rampa, "Federated learning with cooperating devices: A consensus approach for massive IoT networks," *IEEE Internet Things J.*, vol. 7, no. 5, pp. 4641–4654, 2020.
- [20] A. Lalitha, O. C. Kilinc, T. Javidi, and F. Koushanfar, "Peer-to-peer federated learning on graphs," *arXiv preprint arXiv:1901.11173*, 2019.
- [21] S. Hosseinalipour, S. S. Azam, C. G. Brinton, N. Michelusi, V. Aggarwal, D. J. Love, and H. Dai, "Multi-stage hybrid federated learning over large-scale D2D-enabled fog networks," *IEEE/ACM Trans. Network.*, 2022.
- [22] F. P.-C. Lin, S. Hosseinalipour, S. S. Azam, C. G. Brinton, and N. Michelusi, "Semi-decentralized federated learning with cooperative D2D local model aggregations," *IEEE J. Sel. Areas Commun.*, vol. 39, no. 12, pp. 3851–3869, 2021.
- [23] J. George and P. Gurrarn, "Distributed stochastic gradient descent with event-triggered communication," in *Proc. AAAI Conf. Artif. Intell.*, vol. 34, pp. 7169–7178, 2020.
- [24] H. Xiao, K. Rasul, and R. Vollgraf, "Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms," *arXiv preprint arXiv:1708.07747*, 2017.

APPENDIX

A. Proof of Theorem 1

Rewriting the event-based updates of Algorithm 1, we get

$$\begin{aligned} \mathbf{w}_i^{(k+1)} &= \mathbf{w}_i^{(k)} + \sum_{j \in \mathcal{N}'_i(k)} \beta_{ij}^{(k)} \left(\mathbf{w}_j^{(k)} - \mathbf{w}_i^{(k)} \right) v_{ij}^{(k)} \\ &\quad - \alpha^{(k)} \mathbf{g}_i \left(\mathbf{w}_i^{(k)} \right), \\ \widehat{\mathbf{w}}_i^{(k+1)} &= \widehat{\mathbf{w}}_i^{(k)} \left(1 - v_i^{(k)} \right) + \mathbf{w}_i^{(k)} v_i^{(k)}. \end{aligned} \quad (7)$$

Rearranging the relations in (7), we have

$$\begin{aligned} \mathbf{w}_i^{(k+1)} &= \left(1 - \sum_{j=1}^m \beta_{ij}^{(k)} v_{ij}^{(k)} \right) \mathbf{w}_i^{(k)} \\ &\quad + \sum_{j=1}^m \beta_{ij}^{(k)} v_{ij}^{(k)} \mathbf{w}_j^{(k)} - \alpha^{(k)} \mathbf{g}_i \left(\mathbf{w}_i^{(k)} \right) \\ &= \sum_{j=1}^m p_{ij}^{(k)} \mathbf{w}_j^{(k)} - \alpha^{(k)} \mathbf{g}_i \left(\mathbf{w}_i^{(k)} \right). \end{aligned} \quad (8)$$

Next, we collect the vectors of all devices that were previously introduced into matrix form as follows: $\mathbf{W}^{(k)} = \begin{bmatrix} \mathbf{w}_1^{(k)} & \cdots & \mathbf{w}_m^{(k)} \end{bmatrix}^\top$, $\widehat{\mathbf{W}}^{(k)} = \begin{bmatrix} \widehat{\mathbf{w}}_1^{(k)} & \cdots & \widehat{\mathbf{w}}_m^{(k)} \end{bmatrix}^\top$, $\mathbf{G}^{(k)} = \begin{bmatrix} \mathbf{g}_1 \left(\mathbf{w}_1^{(k)} \right) & \cdots & \mathbf{g}_m \left(\mathbf{w}_m^{(k)} \right) \end{bmatrix}^\top$, $\mathbf{P}^{(k)} = [p_{ij}^{(k)}]_{1 \leq i, j \leq m}$.

Now, we transform the recursive update rules of (8) into matrix form to get the following relationship:

$$\mathbf{W}^{(k+1)} = \mathbf{P}^{(k)} \mathbf{W}^{(k)} - \alpha^{(k)} \mathbf{G}^{(k)}. \quad (9)$$

The recursive expression in (9) has been investigated before [7]. In the following, we build upon some lemmas from prior work given our assumptions in Sec. III-B to obtain the final result of the theorem.

Starting from iteration s , where $s \leq k$, we have

$$\begin{aligned} \mathbf{W}^{(k+1)} &= \mathbf{P}^{(k:s)} \mathbf{W}^{(s)} - \sum_{r=s+1}^k \alpha^{(r-1)} \mathbf{P}^{(k:r)} \mathbf{G}^{(r-1)} \\ &\quad - \alpha^{(k)} \mathbf{G}^{(k)}, \\ \mathbf{P}^{(k:s)} &= \mathbf{P}^{(k)} \mathbf{P}^{(k-1)} \cdots \mathbf{P}^{(s+1)} \mathbf{P}^{(s)}. \end{aligned} \quad (10)$$

If we let $s = 0$ in (10), we get an explicit relationship for the model parameters at iteration k with respect to their initial values. Focusing on the parameters of each device i (row i of $\mathbf{W}^{(k+1)}$), we get

$$\begin{aligned} \mathbf{w}_i^{(k+1)} &= \sum_{j=1}^m p_{ij}^{(k:0)} \mathbf{w}_j^{(0)} \\ &\quad - \sum_{r=1}^k \alpha^{(r-1)} \sum_{j=1}^m p_{ij}^{(k:r)} \mathbf{g}_j \left(\mathbf{w}_j^{(r-1)} \right) - \alpha^{(k)} \mathbf{g}_i \left(\mathbf{w}_i^{(k)} \right). \end{aligned} \quad (11)$$

To analyze the local model consensus, we define the average model $\bar{\mathbf{w}}^{(k)}$ as

$$\bar{\mathbf{w}}^{(k)} = \frac{1}{m} \sum_{i=1}^m \mathbf{w}_i^{(k)} = \frac{1}{m} \mathbf{1}_m^\top \mathbf{W}^{(k)}.$$

The recursive relation for $\bar{\mathbf{w}}^{(k)}$ using (9) and the stochasticity of $\mathbf{P}^{(k)}$ is

$$\bar{\mathbf{w}}^{(k+1)} = \frac{1}{m} \mathbf{1}_m^\top \mathbf{W}^{(k+1)} = \bar{\mathbf{w}}^{(k)} - \frac{\alpha^{(k)}}{m} \mathbf{1}_m^\top \mathbf{G}^{(k)}.$$

Also, the explicit relationship connecting $\bar{\mathbf{w}}^{(k+1)}$ to its corresponding value at iteration 0 can be calculated using (11) together with the stochasticity of $\mathbf{P}^{(k:0)}$:

$$\bar{\mathbf{w}}^{(k+1)} = \bar{\mathbf{w}}^{(0)} - \frac{1}{m} \sum_{r=1}^{k+1} \alpha^{(r-1)} \sum_{j=1}^m \mathbf{g}_j \left(\mathbf{w}_j^{(r-1)} \right). \quad (12)$$

Part (a) of the following lemma shows that model parameters $\mathbf{w}_i^{(k)}$ of each device i asymptotically converge to $\bar{\mathbf{w}}^{(k)}$, thus reaching consensus as $k \rightarrow \infty$.

Lemma 1 (Follows from Lemma 8 of [8]). *Let the sequence $\{\mathbf{w}_i^{(k)}\}$ be generated by iteration (11) and the sequence $\{\bar{\mathbf{w}}^{(k)}\}$ be generated by (12). Then $\forall i \in \mathcal{M}$ we have*

- $\lim_{k \rightarrow \infty} \left\| \mathbf{w}_i^{(k)} - \bar{\mathbf{w}}^{(k)} \right\|_2 = 0$, if the step size satisfies $\lim_{k \rightarrow \infty} \alpha^{(k)} = 0$, and
- $\sum_{k=1}^{\infty} \alpha^{(k)} \left\| \mathbf{w}_i^{(k)} - \bar{\mathbf{w}}^{(k)} \right\|_2 < \infty$, if the step size satisfies $\sum_{k=0}^{\infty} (\alpha^{(k)})^2 < \infty$.

We next move on to show that $\bar{\mathbf{w}}^{(k)}$ under our method asymptotically converges to the optimizer of the global loss. First, we provide the following lemma, which reveals the relationship between $F(\cdot)$ evaluated at $\bar{\mathbf{w}}^{(k)}$ and $\mathbf{w}_i^{(k)}$.

Lemma 2 (Follows from Lemma 6 in [8]). *Let the sequence $\{\mathbf{w}_i^{(k)}\}$ be generated by iteration (11) $\forall i \in \mathcal{M}$ and the sequence $\{\bar{\mathbf{w}}^{(k)}\}$ be generated by iteration (12). If Assumptions 3 and 4 hold, we have*

$$\begin{aligned} &\frac{2\alpha^{(k)}}{m} \left(F \left(\bar{\mathbf{w}}^{(k)} \right) - F \left(\mathbf{w}_i^{(k)} \right) \right) \\ &\leq \left\| \bar{\mathbf{w}}^{(k)} - \mathbf{w}_i^{(k)} \right\|_2^2 - \left\| \bar{\mathbf{w}}^{(k+1)} - \mathbf{w}_i^{(k)} \right\|_2^2 \\ &\quad + \frac{L^2}{m} \left(\alpha^{(k)} \right)^2 + \frac{4L}{m} \alpha^{(k)} \sum_{j=1}^m \left\| \bar{\mathbf{w}}^{(k)} - \mathbf{w}_j^{(k)} \right\|_2. \end{aligned}$$

Finally, we only need to show that the average of models $\bar{\mathbf{w}}^{(k)}$ asymptotically optimizes the global loss. To prove this, we take the summation of the relation in Lemma 2 from $k = 0$ to ∞ , and then use the results of Lemma 1-(b) alongside the step size conditions $\lim_{k \rightarrow \infty} \alpha^{(k)} = 0$ and $\sum_{k=0}^{\infty} (\alpha^{(k)})^2 < \infty$. It follows that $\lim_{k \rightarrow \infty} F(\bar{\mathbf{w}}^{(k)}) - F^* = 0$.

B. Further Explanation of Remarks 1-3

Remark 1. For the q -norm of a vector $\mathbf{w} \in \mathbb{R}^n$, we have

$$\|\mathbf{w}\|_u \leq \|\mathbf{w}\|_q \leq n^{\frac{1}{q} - \frac{1}{u}} \|\mathbf{w}\|_u, \quad (13)$$

where $1 \leq q < u$. Also note that based on the way Algorithm 1 defines the event-triggering conditions, the relation $C \|\mathbf{w}_i^{(k)} - \widehat{\mathbf{w}}_i^{(k)}\|_q < r \rho_i \gamma^{(k)}$ holds at every iteration, since otherwise an event will be triggered to ensure it. C is a normalization factor to be derived here.

- (a) Considering all the norms that can be used for a vector, it is only the ∞ -norm that does not depend on the dimension of the vector. Thus, a model-invariant event-triggering condition would result in the relation $\|\mathbf{w}_i^{(k)} - \widehat{\mathbf{w}}_i^{(k)}\|_{\infty} < r \rho_i \gamma^{(k)}$ holding with $C = 1$.
- (b) To not be constrained by the ∞ -norm over the choice of q in $C \|\mathbf{w}_i^{(k)} - \widehat{\mathbf{w}}_i^{(k)}\|_q < r \rho_i \gamma^{(k)}$, and to still ensure invariance over the model dimension n , we can write the following by letting $u \rightarrow \infty$ in (13):

$$\left(\frac{1}{n}\right)^{\frac{1}{q}} \left\| \mathbf{w}_i^{(k)} - \widehat{\mathbf{w}}_i^{(k)} \right\|_q \leq \left\| \mathbf{w}_i^{(k)} - \widehat{\mathbf{w}}_i^{(k)} \right\|_{\infty} < r \rho_i \gamma^{(k)}.$$

Remark 2. Based on (10), we can obtain the following relationship between the state of device i from iteration s to $k \geq s$:

$$\begin{aligned} \mathbf{w}_i^{(k)} &= \sum_{j=1}^m p_{ij}^{(k-1:s)} \mathbf{w}_j^{(s)} - \alpha^{(k-1)} \mathbf{g}_i(\mathbf{w}_i^{(k-1)}) \\ &\quad - \sum_{r=s+1}^{k-1} \alpha^{(r-1)} \sum_{j=1}^m p_{ij}^{(k-1:r)} \mathbf{g}_j(\mathbf{w}_j^{(r-1)}). \end{aligned}$$

Assuming no aggregation events occur at device i or its neighbors from iteration s to iteration k , we will have: (i) $\widehat{\mathbf{w}}_i^{(k)} = \mathbf{w}_i^{(s)}$; and (ii) $p_{ii}^{(k-1:r)} = 1$ and $p_{ij}^{(k-1:r)} = 0$ for all $j \neq i$ and $s \leq r \leq k-1$. As a result, we get

$$\mathbf{w}_i^{(k)} = \mathbf{w}_i^{(s)} - \sum_{r=s}^{k-1} \alpha^{(r)} \mathbf{g}_i(\mathbf{w}_i^{(r)}).$$

In other words, device i solely conducts SGD from iteration s to k , and thus

$$\begin{aligned} \left\| \mathbf{w}_i^{(k)} - \widehat{\mathbf{w}}_i^{(k)} \right\|_q &\leq \sum_{r=s}^{k-1} \alpha^{(r)} n^{\frac{1}{q}} \left\| \mathbf{g}_i(\mathbf{w}_i^{(r)}) \right\|_{\infty} \\ &\leq \alpha^{(s)} n^{\frac{1}{q}} (k-s) L_{\infty}. \end{aligned} \quad (14)$$

The expression above gives us a guideline on selecting the hyperparameter r . Since r has a constant value throughout the training process, we select it in a way to have our desired behavior from the early iterations, i.e., $s = 0$. Considering

the extreme case where maximum steps are taken to update $\mathbf{w}_i^{(k)}$, i.e., steps of size $\alpha^{(0)} \|\mathbf{g}_i(\mathbf{w}_i^{(0)})\|_q \approx \alpha^{(0)} n^{\frac{1}{q}} L_{\infty}$, we set r to a value such that it would take approximately K iterations with maximum steps before the event-triggering condition is reached. Note that for the threshold decay rate $\gamma^{(k)}$, its extreme case value $\gamma^{(0)}$ is considered as well:

$$\begin{aligned} \left(\frac{1}{n}\right)^{\frac{1}{q}} \alpha^{(0)} n^{\frac{1}{q}} K L_{\infty} &= r \rho_i \gamma^{(0)}, \\ r &= \frac{\alpha^{(0)}}{\gamma^{(0)}} \frac{1}{\rho_i} K L_{\infty}. \end{aligned} \quad (15)$$

However, r should be a global variable that has the same value across all devices. Thus, we take the average of the relation above across all devices

$$r = \frac{\alpha^{(0)}}{\gamma^{(0)}} \left(\frac{1}{m} \sum_{i=1}^m \frac{1}{\rho_i} \right) K L_{\infty}.$$

Since in our fully-decentralized setting there is no central server with the knowledge of each ρ_i , calculating $\frac{1}{m} \sum_{i=1}^m \frac{1}{\rho_i}$ exactly is not possible. Thus, we replace that term with an estimate $\frac{1}{\rho}$ to get (6).

Remark 3. Expression (6) in Remark 2 was used to derive a value for the constant r . We use similar arguments to find a relationship between the learning rate $\alpha^{(k)}$ and the threshold decay rate $\gamma^{(k)}$. Using (14) and (15) and solving for $\Delta k_i^u = k_i^{u+1} - k_i^u$, where k_i^u denotes the iteration where the $\{u\}$ -th aggregation event occurs at device i , gives us

$$\Delta k_i^u = \frac{r \rho_i \gamma^{(k)}}{L_{\infty} \alpha^{(k)}}.$$

We are interested in the asymptotic behavior of Δk_i^u as $k \rightarrow \infty$. $\lim_{k \rightarrow \infty} \Delta k_i^u = \infty$ implies that aggregation events become less frequent as time goes by and stop after a while. This contradicts Assumption 6-(b) (bounded intercommunication intervals) and hence should be avoided. There is no particular issue when $\lim_{k \rightarrow \infty} \Delta k_i^u = 0$ in terms of consensus, as it implies an aggregation occurs at every iteration after a while. However, we avoid this situation as it defeats our purpose of sporadic event-triggered communications. Therefore, we aim for having a finite constant value for $\lim_{k \rightarrow \infty} \Delta k_i^u$ (this constant is equal to K , which is the approximate number of iterations between aggregation events), and thus

$$K = \frac{r \rho_i}{L_{\infty}} \lim_{k \rightarrow \infty} \frac{\gamma^{(k)}}{\alpha^{(k)}} \Rightarrow \lim_{k \rightarrow \infty} \frac{\gamma^{(k)}}{\alpha^{(k)}} = \frac{K L_{\infty}}{r \rho_i}.$$

So, the decay rate of $\gamma^{(k)}$ and $\alpha^{(k)}$ should be the same. We next substitute the value of r derived in (6) to get

$$\lim_{k \rightarrow \infty} \frac{\gamma^{(k)}/\gamma^{(0)}}{\alpha^{(k)}/\alpha^{(0)}} = \frac{\rho}{\rho_i}.$$

Similar to the argument made in Remark 2, since both $\gamma^{(k)}$ and $\alpha^{(k)}$ are global variables, we take the average of the relationship above to obtain

$$\lim_{k \rightarrow \infty} \frac{\gamma^{(k)}/\gamma^{(0)}}{\alpha^{(k)}/\alpha^{(0)}} = \rho \left(\frac{1}{m} \sum_{i=1}^m \frac{1}{\rho_i} \right) = 1.$$