

REWAFL: Residual Energy and Wireless Aware Participant Selection for Efficient Federated Learning over Mobile Devices

Yixuan Li, *Student Member, IEEE*, Xiaoqi Qin, *Member, IEEE*, Jiayang Geng, *Student Member, IEEE*, Rui Chen, Yanzhao Hou, Yanmin Gong, *Member, IEEE*, Miao Pan, *Senior Member, IEEE*, and Ping Zhang, *Fellow, IEEE*

Abstract—Participant selection (PS) helps to accelerate federated learning (FL) convergence, which is essential for the practical deployment of FL over mobile devices. While most existing PS approaches focus on improving training accuracy and efficiency rather than residual energy of mobile devices, which fundamentally determines whether the selected devices can participate. Meanwhile, the impacts of mobile devices’ heterogeneous wireless transmission rates on PS and FL training efficiency are largely ignored. Moreover, PS causes the staleness issue. Prior research exploits isolated functions to force long-neglected devices to participate, which is decoupled from original PS designs.

In this paper, we propose a residual energy and wireless aware PS design for efficient FL training over mobile devices (REWAFL). REWAFL introduces a novel PS utility function that jointly considers global FL training utilities and local energy utility, which integrates energy consumption and residual battery energy of candidate mobile devices. Under the proposed PS utility function framework, REWAFL further presents a residual energy and wireless aware local computing policy. Besides, REWAFL buries the staleness solution into its utility function and local computing policy. The experimental results show that REWAFL is effective in improving training accuracy and efficiency, while avoiding “flat battery” of mobile devices.

Index Terms—Federated learning, Mobile devices, Participant selection, Residual energy awareness.

I. INTRODUCTION

NOWADAYS, federated learning (FL) has been migrating from cable-powered datacenters to battery-powered mobile devices [1]. With ever advancing hardware development, mobile devices (e.g., Google Pixel 4a, Galaxy Note20, iPad Pro, MacBook laptop, UAVs, etc.) have increasing on-device computing capability ready for local training. Besides, mobile devices crowdsense the raw training data and provide the foundation for FL beyond datacenters [2], [3]. Following FL principle of training deep neural networks (DNNs) locally

without exposing raw data, FL over mobile devices fosters numerous promising applications in various domains. For example, Google and Apple have deployed FL for computer vision and natural language processing tasks across mobile devices [4]–[7]; NVIDIA exploits FL to create medical imaging AI [8]; many other applications [9]–[11] are emerging in e-Healthcare, autonomous driving, hazard detection in smart home, smart surveillance and monitoring in industrial Internet of Things (IoT), etc. Most of those applications require to train FL models across a crowd of clients, while mobile devices may not all be simultaneously available or suitable for FL training. Therefore, how to appropriately select FL participants is the key to their successful deployment in practice.

The biggest challenges for FL participant selection (PS) stem from the heterogeneity of mobile devices. To improve the performance of FL over mobile devices, research efforts in the literature tried to address the data heterogeneity and system heterogeneity issues, and developed FL PS approaches considering *statistical utility* or/and *system utility* as defined in [12]. For example, to address data heterogeneity/improve statistical utility, the importance-aware scheduling policies are proposed to select mobile devices with significant training improvement. Various importance metrics for mobile devices are defined based on gradient divergence [13], local loss values or gradient norm [14]–[16], and the differences between the local and global models [17]. To address system heterogeneity/reduce FL training latency, the mobile devices with the best instantaneous channel conditions at each round are selected to minimize the communication latency [18], [19]. To improve the efficiency of model aggregation, the central server evaluates the computing capability at model devices and select participants with the “fine” models [11]. Pioneeringly, Oort [12] proposed a utility function that jointly considers statistical and system heterogeneity to guide PS for efficient FL training. While most existing PS approaches focus on improving learning accuracy and reducing training latency from the perspective of global FL system, there are very limited discussions about concerns from participating mobile devices’ side, i.e., the energy consumption of FL training and the residual-battery heterogeneity among candidate mobile devices. It is not trivial since if the selected device cannot participate in FL training, it will backfire FL performance in terms of accuracy and convergence latency. Note that only reducing system-level energy consumption [20]–[22] may not be enough, since the mobile device with very low residual

Y. Li, X. Qin, J. Geng, Y. Hou and P. Zhang are with the State Key Laboratory of Networking and Switching Technology, Beijing University of Posts and Telecommunications, Beijing, 100876, China, (e-mail: lyixuan@bupt.edu.cn; xiaoqiqin@bupt.edu.cn; lelegjx@bupt.edu.cn; houyanzhao@bupt.edu.cn; pzhang@bupt.edu.cn). P. Zhang is also with the Department of Broadband Communication, Peng Cheng Laboratory, Shenzhen 518066, Guangdong, China.

R. chen and M. Pan are with the Department of Electrical and Computer Engineering, University of Houston, Houston, TX 77204 USA (e-mail: rchen19@uh.edu; mpan2@uh.edu).

Y. Gong is with the Department of Electrical and Computer Engineering, University of Texas at San Antonio, San Antonio, TX 78249 USA (e-mail: yanmin.gong@utsa.edu).

battery energy is not feasible for selection, even though its energy consumption for training is small.

Besides, state-of-the-art (SOTA) designs largely ignore the impacts of wireless transmission heterogeneity on local computing policy so as to reduce latency/energy consumption, and its consequent influence on PS in FL training. It is obvious that given a FL task, fast/slow transmissions will decrease/increase the communication latency/energy consumption of uploading local model updates from a mobile device to FL server. However, it is not clear when a mobile device's wireless transmission rate is high/low, how to adjust its own local computing policy over rounds. Additionally, it is not clear how these adjustments affect local computing latency/energy consumption considering device's residual energy, and further PS and FL performance.

Moreover, since not every mobile device participates in every round of FL training, frequently selecting a fixed subset of devices to participate may lead to parameter staleness and biased training, which eventually degrades training performance. Most prior PS designs enforce those high-staleness mobile devices to participate and update their model parameters by a separated function [12], which is decoupled from their original PS designs. In addition, existing staleness solutions ignore the selected participating devices' energy consumption and its residual battery energy.

Motivated by the challenges above, in this paper, we propose a residual energy and wireless aware PS design for efficient FL training over mobile devices (REWAFL). Briefly, REWAFL introduces a novel PS utility function that jointly considers global FL training utilities (statistical utility and global latency utility) and local energy utility, which is aware of energy consumption and residual battery energy of candidate participating mobile devices. Under the proposed PS utility function framework, REWAFL further presents a residual energy and wireless aware local computing policy, which includes a wireless aware local stochastic gradient descent (SGD) computing strategy, and an energy utility aware stopping criterion for increasing local training iterations. Besides, REWAFL buries the staleness solution into its utility function and local computing development policy without incurring any extra or external mechanism.

To demonstrate its effectiveness, we implemented and evaluated a prototype of REWAFL system consisting of FL server (NVIDIA RTX 3090), mobile devices (Android smartphones, tablet computers, and laptops) with heterogeneous battery energy and different wireless transmission techniques (Wi-Fi 5 and 5G) and rates for model updates between them. We conducted extensive experiments with various FL learning tasks, training model and datasets. Our experimental results show that the proposed REWAFL can reduce device dropout ratio, reach the target accuracy with less overall latency and energy consumption, and effectively solve the staleness issue in a self-contained manner compared with SOTA designs.

II. PRELIMINARIES AND MOTIVATION

A. Federated Learning over Mobile Devices

We consider a wireless federated learning system consisting of an edge server (e.g., base station or gNodeB) as the FL

aggregator and a set of \mathcal{S} mobile devices as FL clients. We assume that the computing resources (e.g., GPU and CPU frequencies), wireless transmission rates (i.e., s), data distribution and residual battery energy (i.e., E) among mobile devices are heterogeneous. The edge server and mobile devices jointly execute FL training as follows. The edge server first broadcasts an initialized global model to the participating mobile devices. After receiving the global model, each selected mobile device i utilizes local computing resources to train its local models by performing $H(i)$ local iterations based on its local data of size $|B_i|$. When the local training is complete, mobile devices upload model parameters to the edge server using available wireless accessing technologies, i.e., Wi-Fi 5, 5G, etc. After that, the edge server aggregates the local model to update the global model and selects the participants for the next round based on the defined PS utility function. The procedure above repeats until FL converges, where PS plays an essential role for FL's learning performance and efficiency.

B. Revisiting SOTA PS Designs

Some pioneering FL PS designs in the literature considered not only FL learning performance, but also FL training efficiency in terms of global system latency or local participating mobile devices' energy consumption. For example, Oort in [12] proposed the PS design to associate the statistical utility (i.e., contributions to FL learning performance) and system efficiency (i.e., system latency, which is the sum of local computing latency and communication latency). In Oort, mobile device i 's utility function for PS is formulated as follows.

$$Util(i) = |B_i| \sqrt{\frac{1}{|B_i|} \sum_{k \in B_i} Loss(k)^2} \times \left(\frac{T}{t(i)}\right)^{\mathbb{I}(T < t(i)) \times \alpha}, \quad (1)$$

where B_i denotes the local training samples of mobile device i , $Loss(k)$ denotes the training loss of the k -th sample, $t(i)$ denotes the system latency of mobile device i , T is the developer-preferred duration of each round and α is the penalty factor. Besides, $\mathbb{I}(x)$ is an indicator function, where $\mathbb{I}(x) = 1$ if x is true, and 0 otherwise.

Another example is AutoFL in [20], which proposed a PS design to improve learning performance while reducing total energy consumption. Briefly, AutoFL introduced a reward function to associate learning accuracy and the energy consumption of mobile devices, and exploited reinforcement learning to select participants with high contributions to FL training and low energy consumption.

Although Oort/AutoFL exhibits good learning performance and global latency/local energy consumption efficiency, SOTA PS approaches ignored the battery energy constraints and various transmission conditions confronted by candidate mobile devices in FL training, i.e., (i) the battery-powered mobile devices may have heterogeneous levels of residual energy, and (ii) mobile devices have to face heterogeneous wireless transmission environments. Such ignorance may downgrade the training performance/efficiency of FL over mobile devices,

TABLE I
The Dropout Ratio of SOTA PS Designs

Local Model	Target Acc. (%)	Dropout Ratio (%)	
		Oort	AutoFL
CNN@HAR	89.3	85.0	55.0
CNN@CIFAR10	72.2	46.0	25.0
LSTM@Shakespeare	50.3	66.0	53.0

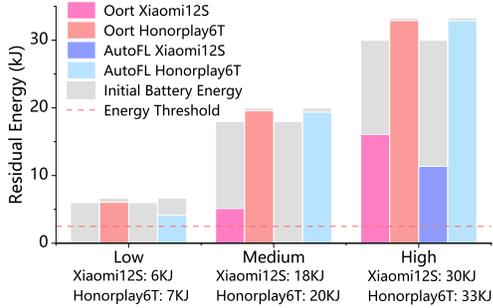


Fig. 1: Energy consumption and residual energy of devices with different initial energy (CNN@MNIST).

or even make the selected mobile devices unable to participate in FL training.

C. Unawareness of Heterogeneous Residual Battery Energy

Mobile devices are powered by the battery, which has a limited energy capacity. If the FL PS design is unaware of mobile devices’ residual energy heterogeneity, it may quickly deplete the energy of frequently selected mobile devices. Since a mobile device has to reserve certain amount of energy for its mandatory/regular operations (e.g., display, voice calls, data services, location services, etc.), such designs may drain out the energy of frequently selected devices and disable them from participation afterwards in FL training. That will have detrimental effects on FL learning performance and efficiency, especially when there is a heterogeneous data distribution among candidate devices.

In order to understand the impact of residual battery energy heterogeneity on training efficiency, we conducted empirical experiments to validate our projection. Here, we follow SOTA designs, i.e., Oort in [12] and AutoFL in [20], to select participants in each training round. The training performances are acquired on mobile devices, i.e., Android smartphones, tablet and laptop. We set NVIDIA RTX 3090 as the FL server. We consider a hybrid communication scenario with 5G and Wi-Fi 5. More details of the experimental setup can be found in Section V.

First, we define the dropout ratio as the ratio of the number of mobile devices that cannot participate in FL training due to depleted battery energy and the total number of candidate mobile devices. From the results in Table I, we observe that both Oort and AutoFL have very high dropout ratio to reach the target accuracy under different learning tasks. The potential reason is that the SOTA PS designs are not aware of candidate mobile devices’ residual energy, and thus may drain out the batteries of frequently selected devices. Then, for CNN@MNIST, we take the high-end mobile devices

(Xiaomi 12S smartphone with Adreno 730 GPU) with high transmission rates (i.e., 79.60Mbps for 5G) and the low-end devices (Honor Play 6T smartphone with Mali-G57 MC2 GPU) with low transmission rates (i.e., 0.64Mbps for 5G) as examples to present how much energy is left for mobile devices with different initial battery energy levels after FL training. From the results in Fig. 1, we find that high-end mobile devices use a lot or even use up its battery energy, while low-end devices don’t use much. It indicates that high-end devices are much more frequently selected for training than low-end ones, in order to reduce global latency [12] or energy consumption [20]. Again, such designs may deplete the energy of frequently selected devices with low initial battery energy at early rounds, and exclude them from participation at later training rounds, which may backfire FL learning performance and latency/energy efficiency.

D. Unawareness of Heterogeneous Wireless Transmission Rates

Mobile devices confront various wireless transmission environments, and thus have heterogeneous transmission rates. During each FL training round, mobile device i ’s wireless transmission rate $s(i)$ explicitly affects its energy consumption, training latency and residual energy, all of which are important metrics for FL PS. As we know, the energy consumption/latency of mobile device i per FL training round consists of two parts: on-device local computing energy consumption/latency, and communication energy consumption/latency for uploading local model updates. Obviously, for a given FL task, the high/low transmission rate $s(i)$ reduces/increases both the energy consumption and latency of mobile device i for uploading the model updates to FL server.

However, it is not clear what the corresponding local computing policy should be, when the wireless transmission rate is high/low. Questions to answer include: (i) How to adjust local computing to be aware of the candidate mobile device’s transmission rate? (ii) How to adjust local computing to be aware of the candidate device’s residual energy? (iii) How does the local computing policy affect PS and FL performance? In order to achieve good learning performance and training efficiency, a fine-grained local computing policy under the FL PS framework is in need.

E. Limitations of Oort’s Staleness Solution

Frequently selecting a fixed subset of mobile devices to participate may result in the staleness issue in FL training. Oort in [12] provided a solution by introducing an extra “temporal uncertainty” mechanism, i.e., if mobile device i has not been selected for a long time, its statistical utility will be forced to increase by adding an independent component characterized by the round counter and mobile device i ’s last involved round. Such a staleness solution is isolated from the PS utility function. Besides, it ignores their participating devices’ energy consumption and its residual battery energy. It would be desirable if (i) the staleness solution can be aware of participating devices’ energy consumption, wireless transmission rate, residual energy levels, etc., and (ii) the

staleness issue can inherently be addressed by the PS design in FL training.

III. ALGORITHM DESIGN

Aiming to address the concerns of both global FL system utility and local mobile devices' energy constraints, we propose a REWAFL scheme that includes: (i) a residual energy aware (REA) PS utility function design, (ii) a fine-grained REWA local computing policy under the proposed PS utility function framework, and (iii) a self-contained REWA staleness solution. The sketch of the REWAFL design is shown in Fig. 2.

A. REAFL: REA PS Utility Function Design

As illustrated in Sec. II-C, for FL over mobile devices, it may not be adequate for PS utility function to just consider the learning performance while reducing energy consumption or/and latency of FL, which may drain out the battery of frequently selected mobile devices. The PS utility function design has to be aware of the residual energy of candidate mobile devices.

Our intuitive idea behind the residual energy aware PS utility function design is to trade-off and jointly improve statistical utility, global system/latency efficiency and energy efficiency in FL training, while satisfying the residual energy constraints of mobile devices. Rooting from Oort's utility function [12] in Eqn. (1), we propose the residual-energy aware PS utility function as follows¹.

$$\begin{aligned}
 Util(i, r) = & \underbrace{|B_i^r| \sqrt{\frac{1}{|B_i^r|} \sum_{k \in B_i^r} Loss(k)^2}}_{\text{Statistical utility}} \times \underbrace{\left(\frac{T^r}{t(i, r)}\right)^{\mathbb{I}(T^r < t(i, r)) \times \alpha}}_{\text{Global latency utility}} \\
 & \times \underbrace{\left(\frac{E_i^r - E_0}{e(i, r)}\right)^{\mathbb{U}(e(i, r) < E_i^r - E_0) \times \beta}}_{\text{Local mobile device's energy utility}},
 \end{aligned} \tag{2}$$

where $|B_i^r|$ is the number of training data samples of mobile device i at the r -th round, $t(i, r)$ is the training latency of mobile device i , T^r is the developer-preferred duration at r th round [12], $e(i, r)$ is the energy consumption of mobile device i at the r -th round, E_i^r is residual battery energy of mobile device i at the r -th round, E_0 is the threshold energy reserved for its mandatory/regular operations, and thus $E_i^r - E_0$ is the available energy to consume. Moreover, $\mathbb{I}(x)$ is an indicator function, where $\mathbb{I}(x) = 1$, if x is true, and 0 otherwise. $\mathbb{U}(x)$ is also an indicator function that takes value 1 if x is true, and ∞ otherwise. Here, α and β are scaling coefficients to balance/associate different metrics/utilities in terms of the model accuracy, latency and energy efficiency.

Different from the existing PS utility function designs (e.g., Oort [12] or AutoFL [20]), the proposed utility function

considers not only the global FL training utility including statistical utility and global latency utility (similar to those in Oort [12]), but also the local mobile device's energy utility, which is aware of heterogeneous residual battery energy of candidate mobile devices. Here is a simple interpretation of the energy utility. During the r -th round, if the energy consumption of mobile device i is less than its residual available battery energy, i.e., $e(i, r) < E_i^r - E_0$, the energy utility of mobile device i is proportional to the residual available battery energy and inversely proportional to the energy consumption. That makes mobile devices with more residual energy and less energy consumption easier to be selected for participation. Otherwise, if the energy consumption of mobile device i is greater than or equal to its residual available energy, i.e., $e(i, r) \geq E_i^r - E_0$, the value of the energy utility sharply drops to zero and mobile device i will not be able to join model training in the r -th round. Therefore, the energy efficiency of the proposed energy utility component has two-fold meanings, i.e., low energy consumption and residual energy awareness.

Following the proposed utility function in Eqn. (3), at the beginning of the r -th FL training round, the edge server collects the reported information (i.e., $Loss, t(i, r)$, and device i 's energy utility)² from mobile devices and calculate their utilities to select the participants for this round. Similar to that in Oort [12], we exploit mobile devices' most recent aggregate training losses to estimate their current ones. Besides, mobile device i needs to estimate its $t(i, r)$ and local energy utility, where $t(i, r)/e(i, r)$ consists of local computing latency/energy consumption and communication latency/energy consumption. Given a FL learning task, assuming device i 's transmission power is fixed, the size of device i 's local model update transmitted to FL server does not change over global training rounds. Thus, the communication latency/energy consumption of mobile devices can be calculated given a known transmission rate $s(i, r)$. Moreover, since device i 's computing capability is fixed, the computing latency/energy consumption can be approximately estimated by the product³ of the latency/energy consumption of each iteration and the number of local iterations $H(i, r)$ at the r -th round. Thus, the estimation of $t(i, r)$ and device i 's energy utility relies on $H(i, r)$. Under the proposed PS utility function framework, in the following subsection, we further illustrate the wireless aware local computing policy (i.e., how to adjust $H(i, r)$) to further improve the energy/latency efficiency for FL training over mobile devices.

B. REWA Local Computing Policy

The proposed REWA local computing policy for PS includes two parts: (i) a wireless aware local stochastic gradient

¹Note that the global latency utility in Eqn. (2) is the same as the global system utility defined in [12]. We rename them to (i) explicitly distinguish global FL training utility from local mobile device's energy utility, and (ii) jointly consider statistical, global latency and device's energy utilities.

²As for the privacy preservation of $Loss$ and $t(i, r)$, we follow the same analysis in Oort [12]. For the privacy of $e(i, r)$ and E_i^r , (i) $e(i, r)$ and E_i^r are not that related to the training data, which is sensitive and worth privacy protection, (ii) the client can self-calculate its local energy utility and report it without disclosing $e(i, r)$ and E_i^r , and (iii) the client can report the estimated energy utility, or even dishonestly report its energy utility, while the benefits of such dishonest reporting and the consequent games between the client and the server or among clients are out of the scope of this paper.

³We neglect the non-linear impacts [21] brought by advanced hardware latency/energy saving techniques (e.g., DVFS) for local on-device training.

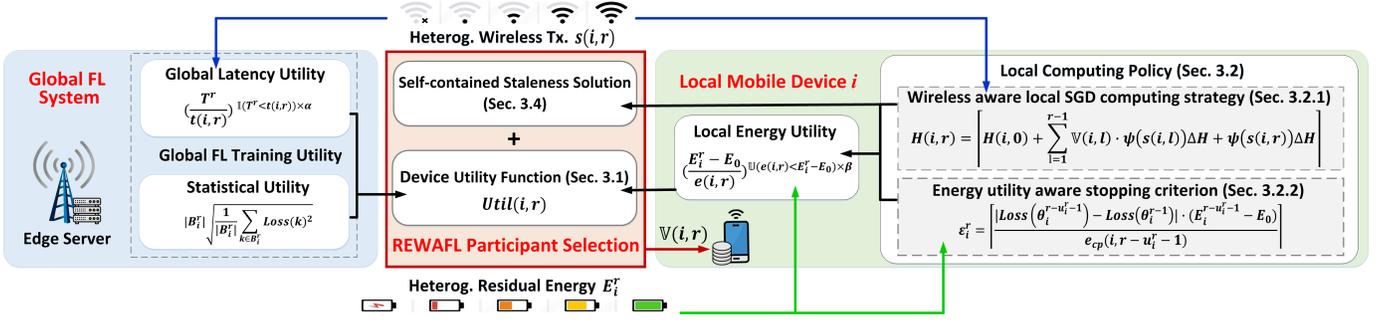


Fig. 2: The design sketch of residual energy and wireless aware federated learning (REWAFL).

descent (SGD) computing strategy, and (ii) an energy utility aware stopping criterion for $H(i)$.

1) *A wireless aware local SGD computing strategy*: Let us start with existing adaptive local SGD computing strategy, AdaH in [23], which gradually increases the number of local iterations over rounds, i.e., $H(i,r) = \left\lceil H(0) + \sum_{l=1}^r \psi \Delta H \right\rceil$. Here, $H(i,r)$ is the number of local training iterations for mobile device i at the r -th round, $H(0)$ is the initial value of the number of local iterations, ψ is the growth-rate coefficient, and ΔH is a non-negative increment unit. The rationale behind it is as follows. At the beginning of the training process, the learning performance is less sensitive to the number of local iterations. Even if the local model is updated with fewer local iterations in the early rounds, learning performance is barely affected. However, as FL training proceeds, it is more sensitive to local model quality, and thus requires more local training iterations for better learning accuracy. By using such an AdaH strategy, the statistical utility in Eqn. (2) can be effectively improved.

To tailor the AdaH strategy [23] for FL PS, we introduce a binary variable $\mathbb{V}(i,r)$ to represent the selection status of mobile device i at the r -th round. If mobile device i is selected, $\mathbb{V}(i,r) = 1$, and 0 otherwise. Considering participant selection, we can update the local AdaH in [23] to $H(i,r) = \left\lceil H(i,0) + \sum_{l=1}^{r-1} \mathbb{V}(i,l) \cdot \psi \Delta H + \psi \Delta H \right\rceil$. Thus, the relationship between the local iteration number of the r -th round and that of the last participation round can be characterized as follows. If the mobile device i is selected to participate in training at the current r -th round, $H(i,r)$ will increase based on $H(i,r - u_i^r - 1)$ at device i 's latest participating round, i.e., $H(i,r) = \lceil H(i,r - u_i^r - 1) + \psi \Delta H \rceil$, where u_i^r is the number of non-participating rounds between the current r -th round and mobile device i 's latest participating round. Otherwise, $H(i,r) = H(i,r - u_i^r - 1)$.

Furthermore, to make the local SGD computing strategy aware of wireless communications, we refine it as follows.

$$\begin{aligned} H(i,r) &= \lceil H(i,r - u_i^r - 1) + \psi(s(i,r))\Delta H \rceil \\ &= \left\lceil H(i,0) + \sum_{l=1}^{r-1} \mathbb{V}(i,l) \cdot \psi(s(i,l))\Delta H + \psi(s(i,r))\Delta H \right\rceil, \end{aligned} \quad (3)$$

where $s(i,r)$ is the averaged wireless transmission rate from device i to the edge server at the r -th round, and $\psi(\cdot)$ is a non-

negative function that decreases with the wireless transmission rate. For the r -th round, if transmission rate $s(i,r)$ is high, the increment of $H(i,r)$ will be small; otherwise, if $s(i,r)$ is low, the increment of $H(i,r)$ will be large.

The rationale behind the strategy in Eqn. (3) is threefold. First, since it inherits ‘‘increasing H over rounds’’ strategy from AdaH [23] and is tailored for PS, it helps to increase the learning accuracy and *statistical utility*. Second, such a wireless aware local SGD computing strategy makes device i with higher transmission rate at the r -th round have smaller increase of local iterations, i.e., less computing workload added. Thus, for a given local learning task, it helps to reduce device i 's local computing latency in the r -th round of FL training. Since the communication latency has already been reduced due to the fast transmission rate, the wireless aware local computing strategy helps to increase the *global latency utility* for mobile device i . Third, once $H(i,r)$ is determined, device i 's local computing energy consumption can be estimated for the r -th round. Similar to computing latency, small increase of local iterations leads to small increase of local computing energy consumption for device i with high transmission rate. Meanwhile, assuming device i 's wireless transmission power is fixed, the energy consumption of local gradient updates transmissions will be reduced for device i , if its transmission rate is high. Therefore, the wireless aware local computing strategy in Eqn. (3) helps to reduce the energy consumption of device i with fast transmission and may increase its *energy utility*, which is jointly determined by both the energy consumption and the residual energy of device i . With such a wireless aware local computing strategy, mobile devices with high transmission rate may have good utility as defined in Eqn. (2), so that they may have good chances to be selected for participation.

Although the proposed wireless aware local computing strategy seems to favor the devices with fast transmissions, it actually also creates participation opportunities for the devices with slow transmissions, and addresses the staleness issue in an inherent manner. The details of staleness analysis and solution will be provided in Sec. III-D.

2) *An energy utility aware stopping criterion for $H(i,r)$* : As $H(i,r)$ keeps increasing over rounds, the computing energy consumption of device i per round will increase in parallel. An unstoppable increase of $H(i,r)$ may ultimately drain out mobile device's battery energy. That prompts us to develop a stopping criterion for $H(i,r)$'s increase, which is

aware of device i 's energy utility (i.e., energy consumption and residual energy). Our idea for the energy utility aware stopping criterion for increasing $H(i)$ is simple. When increase $H(i, r)$ fails to notably reduce the value of the loss function but consumes a lot of computing energy instead, the increasing of $H(i, r)$ should be stopped. Besides, if the residual battery energy of mobile device i is low, we should avoid further increase of $H(i, r)$, allowing the mobile device to retain enough energy for its regular operations. Following this idea, we propose the energy utility aware stopping criterion for increasing $H(i, r)$ as follows.

$$\varepsilon_i^r = \frac{|Loss(\theta_i^{r-u_i^r-1}) - Loss(\theta^{r-1})| \cdot (E_i^{r-u_i^r-1} - E_0)}{e_{cp}(i, r - u_i^r - 1)}, \quad (4)$$

where $\theta_i^{r-u_i^r-1}$ denotes the local model parameters at the last participating round, i.e., the $(r - u_i^r - 1)$ -th round, θ^{r-1} denotes the global model parameters at the $(r - 1)$ -th round, and $e_{cp}(i, r - u_i^r - 1)$ is the computing energy consumption of mobile device i at the $(r - u_i^r - 1)$ -th round. Here, if ε_i^r is large, we will continue increasing $H(i, r)$ following Eqn. (3); otherwise, if ε_i^r is smaller than a predefined threshold ε_{th}^r , we will stop increasing $H(i, r)$, and let $H(i, r) = H(i, r - u_i^r - 1)$.

C. REWA Participant Selection Procedure

Given the residual energy aware PS utility function and the corresponding REWA local computing policy, we summarize the REWAFI PS procedure in Algorithm 1. Specifically, at the beginning of each FL training round, each mobile device determines $H(i, r)$ according to Eqn. (3) (Line 8). Given $H(i, r)$ of each mobile device, $|B_i^r|$, $Loss$, $t(i, r)$, $e(i, r)$ and E_i^r can be calculated/estimated (Line 9). Next, the mobile devices send those values to the edge server with negligible communication cost (Line 10). When the edge server receives those values from mobile devices at the r -th round (Line 13), it calculates $Util(i, r)$ according to Eqn. 2 (Line 14). Then, the edge server sorts $Util(i, r)$ in the descending order, and selects the top K devices for participation (Line 15). After that, the edge server broadcasts the PS decision to mobile devices (Line 16). After receiving the decision, if mobile device i is selected to participate in the r -th round training, it will update its residual energy (Line 20, i.e., $E_i^r \leftarrow E_i^{r-1} - e(i, r)$), the number of non-participating rounds (Line 21, $u_i^r \leftarrow 0$), and the number of local iterations (Line 22, $H(i, r) \leftarrow \lceil H(i, r - u_i^r - 1) + \psi(s(i, r))\Delta H \rceil$). By contrast, if mobile device i is not selected for participation in the r -th round, E_i^r and $H(i, r)$ remain unchanged, while u_i^r increases by 1 (Line 24-26).

D. Self-Contained Staleness Solution

Since not every mobile device participates in every round of FL training, PS may lead to parameter staleness and biased training, and eventually degrade FL convergence or model accuracy. Most prior PS designs exploit a separated function to enforce those high-staleness mobile devices to participate and update their model parameters. For example, Oort [12] forcibly increases the utility of long-neglected devices and

Algorithm 1 REWAFI: Residual Energy and Wireless Aware Participant Selection.

- 1: **Input:** Mobile device set \mathcal{S} , participant size K , weighting parameter α and β , threshold energy reserve E_0 , stop threshold ε_{th}^r .
 - 2: **Output:** Participants \mathcal{P} , Local iterations \mathcal{H} .
/* Initialize global variables. */
 - 3: At the edge server, initialize participants $\mathcal{P} \leftarrow \emptyset$, training round $r \leftarrow 0$, and utility value of each mobile device $Util \leftarrow \emptyset$.
 - 4: At the mobile device, initialize the residual energy $E^r \leftarrow E^0$, and the number of non-participating rounds $u^r \leftarrow 0$.
 - 5: For the r -th round FL training, where $r \in \{1, \dots, R\}$,
 - 6: **On Mobile Devices:**
 - 7: **for** $i \in \mathcal{S}$ **do**
/*Calculate the number of local iterations. */
 - 8: Calculate $H(i, r)$ according to Eqn. (3)
/*Calculate/estimate the utility parameters. */
 - 9: Calculate/estimate $|B_i^r|$, $Loss$, $t(i, r)$, $e(i, r)$ and E_i^r based on $H(i, r)$
 - 10: Send $|B_i^r|$, $Loss$, $t(i, r)$, $e(i, r)$ and E_i^r to FL server
 - 11: **end for**
 - 12: **On Edge Server:**
 - 13: Receive the utility parameters from mobile devices
/*Calculate mobile device utility. */
 - 14: Calculate $Util(i, r)$ according to Eqn. (2)
/*Sort mobile device utility and select participants.*/
 - 15: $\mathcal{P} = \text{RankingDevice}(Util(i, r), K)$
 - 16: Broadcast the device selection decision $\mathbb{V}(i, r)$
 - 17: **On Mobile Devices:**
 - 18: **for** $i \in \mathcal{S}$ **do**
 - 19: **if** $\mathbb{V}(i, r) = 1$ **then**
/*Update residual energy*/
 - 20: $E_i^r \leftarrow E_i^{r-1} - e(i, r)$
/*Update number of non-participating rounds*/
 - 21: $u_i^r \leftarrow 0$
/*Update number of local iterations*/
 - 22: $H(i, r) \leftarrow \lceil H(i, r - u_i^r - 1) + \psi(s(i, r))\Delta H \rceil$
 - 23: **else**
 - 24: $E_i^r \leftarrow E_i^{r-1}$ ▷ Residual energy
 - 25: $u_i^r \leftarrow u_i^{r-1} + 1$ ▷ Non-participating rounds
 - 26: $H(i, r) \leftarrow H(i, r - u_i^r - 1)$ ▷ Local iteration
 - 27: **end if**
 - 28: **end for**
-

makes them participate by using a ‘‘temporal uncertainty’’ mechanism, which is decoupled from Oort’s original utility function design.

Different from existing PS designs, the proposed REWAFI addresses the staleness issue in a self-contained manner and buries the solution into the REWAFI’s utility function and its local computing policy. Specifically, mobile devices with larger utility are more likely to be selected for training, where the utility of a mobile device is jointly determined by statistical utility, latency utility, and energy utility. As illustrated in Sec. III-B, given the training dataset, learning task and residual energy, a mobile device’s statistical utility, latency utility and

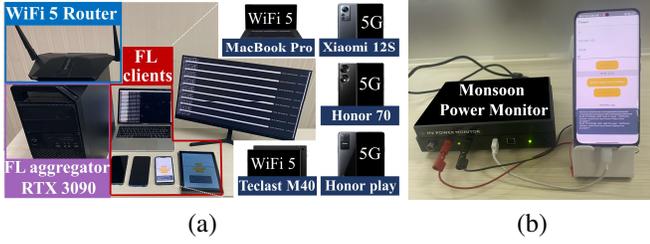


Fig. 3: REWAFL testbed in the lab: (a) FL testbed configuration; (b) Xiaomi 12S’ energy consumption measured by the Monsoon Power Monitor [24].

energy utility are all related to its local computing policy. Following the wireless aware local SGD computing strategy in Eqn. (3), if mobile device i is frequently selected to participate over rounds, its $H(i)$ continues increasing. Generally speaking, with $H(i)$ ’s increase over rounds, device i ’s local training energy consumption increases and residual energy decreases, which leads to a decrease of its energy utility; its training latency increases, which leads to a decrease of its latency utility; its marginal contributions to FL’s global model training are diminishing, which leads to a decrease in its statistical utility. The utility $Util(i)$ of frequently selected device i keeps decreasing until it is less than that $Util(j)$ of a long-neglected device j , whose $H(j)$ has a small value, and then mobile device j will be selected to participate in this training round. In this way, staleness issue can be addressed by REWAFL’s inherent local SGD computing strategy (i.e., the wireless aware adaptive increasing of $H(\cdot)$) in a self-contained manner.

To further illustrate REWA’s self-contained staleness solution, we take two mobile devices i and j with the same data distribution, computing capabilities and initial battery energy but different wireless transmission rates, say $s(i) > s(j)$, as a simple example. At the beginning of FL training, since $s(i) > s(j)$, we have $H(i) < H(j)$ following Eqn. (3), so that $e(i) < e(j)$ and $t(i) < t(j)$. Given the fact that the local model updates of two mobile devices have similar contributions to learning at the beginning of FL training, we have $Util(i) > Util(j)$, and device i will be selected for participation. As FL training proceeds, $H(i)$ continuously increases over rounds following Eqn. (3), while $H(j)$ remains unchanged. Until the r -th round, when $H(i, r)$ is larger than a certain value that triggers $Util(i, r) < Util(j, r)$, the staled mobile device j will be selected to participate in this round of FL training. Similar analysis can be applied to mobile devices with data or/and device heterogeneity, where $H(\cdot)$ inherently serves as a tuning knob to tackle staleness in REWAFL.

IV. IMPLEMENTATION & EXPERIMENTAL SETUP

A. REWAFL Implementation

We present the setup of testbed in Fig. 3. Our REWAFL testbed consists of one model aggregation server and different types of mobile devices. On FL aggregator side, we use a NVIDIA RTX 3090. On FL client side, we consider five types of mobile devices: (1) Xiaomi 12S smartphone with

Qualcomm Snapdragon 8+ Gen 1 CPU, Adreno 730 GPU, 8GB RAM, and a battery capacity of 4500mAh; (2) Honor 70 smartphone with Qualcomm Snapdragon 778G Plus CPU, Qualcomm Adreno 642L GPU, 12GB RAM, and a battery capacity of 5000mAh; (3) Honor Play 6T smartphone with MediaTek Dimensity 700 CPU, Mali-G57 MC2 GPU, 8GB RAM, and a battery capacity of 5000mAh; (4) Teclast M40 tablet with Unisoc Tiger T618 CPU, 6GB RAM, and a battery capacity of 7000mAh; (5) MacBook Pro 2018 laptop with Intel Core i5-8259U CPU, Intel Iris Plus Graphics 655 GPU, 8GB RAM, and a battery capacity of 58Wh. For each mobile device, its initial battery energy follows a normal distribution within the range of its battery’s maximum capacity. We set up a REWAFL system consisting of 100 mobile devices, with 20 devices for each type of mobile devices mentioned above. We select 20 mobile devices for participation in every training round.

For communication between FL clients and FL aggregator, we employ a complex wireless transmission scenario, which consists of a mixture of Wi-Fi 5 and 5G networks. Specifically, the Teclast M40 tablet and MacBook Pro 2018 laptop are connected to FL aggregator via Wi-Fi 5, which uses the WebSocket communication protocol [25]. The Xiaomi 12S smartphone, Honor 70 smartphone, and Honor Play 6T smartphone are connected to FL aggregator via 5G networks, which follows the 5G NR standard. We further consider two types of transmission rates (high and low), and each device can be configured to one of these types. As for a Wi-Fi transmission environment, the command line tools nmcli and wondershaper are used for reporting network status and controlling different wireless transmission rates. As for 5G transmission environments, we consider two transmission environments, i.e., indoor and outdoor, which represent different 5G transmission rates.

REWAFL is implemented by building on top of FLOWER [26]. In total, we added a module to control the number of local iterations on the client side and a PS module on the server side, based on the original FLOWER code. In order to realize REWAFL, we need to know the relevant parameters, such as transmission rate, training latency and energy consumption. To estimate on-device transmission rates, we use the Network Monitor toolbox in the Android kernel on the smartphone. To measure the latency, the time of model training and transmission are recorded on mobile devices. To obtain the energy consumption, the Monsoon Power Monitor [24] is employed to measure the power consumption of smartphones and tablet. As for laptop, HWiNFO is used to monitor the real-time power consumption.

B. Models, Datasets and Parameters

We evaluate REWAFL’ performance for three different FL tasks: image classification, next word prediction, and human activity recognition. We consider two DNN models: 2-layer CNN [27] and LSTM [28].

As for the image classification task, we use two datasets. One is a MNIST dataset [29] consisting of 10 categories ranging from digit “0” to “9”, which includes 60,000 training images and 10,000 validation images. The other is CIFAR10 dataset [30] consisting of 10 categories, which includes 50000

training images and 10,000 validation images. As for data distribution among mobile devices, we denote λ as the non-i.i.d. levels, where $\lambda = 0$ indicates that the data among mobile devices is i.i.d., $\lambda = 0.8$ indicates that 80% of the data belong to one label and the remaining 20% data belong to other labels, and $\lambda = 1$ indicates that each mobile device owns a disjoint subset of data with one label. We train the image classification data samples with 2-layer CNN.

As for next word prediction task, we use Shakespeare dataset [31] that consists of 1,129 roles and each of which is viewed as a device. Since the number of lines and speaking habits of each role vary greatly, the dataset is non-i.i.d. and unbalanced. We train the Shakespeare data samples with LSTM.

As for task of human activity recognition, we focus on a publicly accessible dataset generated by having volunteers wear Samsung Galaxy S2 smartphones equipped with accelerometer and gyroscope sensors [32]. The dataset includes six categories of activities, i.e., walking, walking upstairs, walking downstairs, sitting, standing and laying. The dataset is non-i.i.d. due to the different behavioral habits of the volunteers. It contains 10,299 data samples. We employ a 2-layer CNN to train it.

The default parameter settings are given as follows unless specified otherwise. For the MNIST and CIFAR10 datasets, the non-i.i.d. level of data distribution among mobile devices is set as $\lambda = 0.8$. Besides, the scaling coefficients in the proposed REA PS utility function are set as $\alpha = 1$ and $\beta = 1$.

C. Baselines for Comparison

We compared REWAFI with the following peer FL designs under different learning tasks, DNN models and datasets.

Random [33]: The FL server selects participants randomly, and the selected mobile devices utilize a fixed local computing policy.

Oort [12]: The FL server selects mobile devices based on global FL training utility function as defined in Eqn. (1). The selected mobile devices utilize a fixed local computing policy.

AutoFL [20]: The FL server selects mobile devices based on mobile devices' energy consumption, and the selected mobile devices utilize a fixed local computing policy.

REAFI: The FL server selects mobile devices based on REA PS utility function as defined in Eqn. (2), and the selected mobile devices utilize a fixed local computing policy.

REAFI+LUPA: The FL server selects mobile devices based on REA PS utility function as defined in Eqn. (2), and the selected mobile devices utilize the AdaH strategy during FL training [23], i.e., $H(i, r) = \left| H(0) + \sum_{l=1}^r \psi \Delta H \right|$.

V. EVALUATION RESULTS AND ANALYSIS

A. Advantages of REA PS Utility Function

1) Awareness of devices' heterogeneous residual energy:

As shown in Table II, given a targeted testing accuracy, the dropout ratio of REAFI is much lower than that of *Random*, *Oort*, or *AutoFL* for all learning tasks. The reason is simple. REAFI employs the proposed residual energy aware PS utility

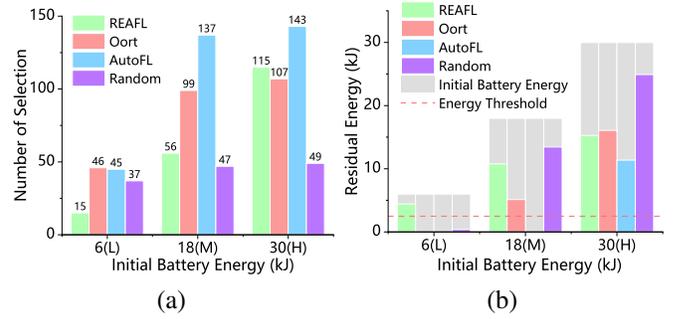


Fig. 4: Performance comparison of different PS utility designs (CNN@MNIST, Xiaomi 12S smartphone, 5G-79.60Mbps). (a) the number of selections, and (b) energy consumption and residual energy.

function and avoids frequently selecting mobile devices with low residual energy, while other SOTA PS designs are not aware of participating mobile devices' residual energy levels at all. That may result in the depletion of energy for frequently selected high-end mobile devices (i.e., mobile devices with good learning contribution, short latency, and/or small energy consumption), especially for those mobile devices with low initial residual energy.

To further verify the above analysis, we have recorded the number of selections and energy consumption for the high-end mobile devices (Xiaomi 12S smartphone with the averaged 5G uplink transmission rate of 79.60Mbps) with heterogeneous initial residual energy levels (i.e., 6kJ-low, 18kJ-medium, and 30kJ-high) under different PS utility function designs, and shown the results in Fig. 4. From Fig. 4(a) and Fig. 4(b), we find that the proposed REA PS utility function represented by REAFI is aware of heterogeneous residual energy among mobile devices, so it doesn't select the high-end mobile device with low residual energy too often, and keeps the necessary energy for participating devices' mandatory operations. By contrast, other PS utility function designs (i.e., *Random*, *Oort*, or *AutoFL*) are not aware of mobile devices' residual energy, and often select high-end mobile devices to participate in training even when their residual energy is low, thus may drain out the battery of such devices.

2) *REAFI energy and latency efficiency analysis:* Besides cutting off dropout, REAFI also helps to reduce the overall energy consumption and latency of FL training to achieve a target testing accuracy. As shown in Table II, in terms of overall latency, REAFI reduces around 51.2%, 45.0%, 37.9% and 46.4% for learning tasks CNN@MNIST, CNN@CIFAR10, LSTM@Shakespeare, and CNN@HAR, compared with *Oort*, respectively, around 75.3%, 74.9%, 65.4% and 67.0% compared with *AutoFL*, and around 59.2%, 67.9%, 64.0% and 53.8% compared with *Random*; in terms of overall energy consumption, REAFI saves around 54.3%, 23.9%, 18.2% and 54.9% for learning tasks CNN@MNIST, CNN@CIFAR10, LSTM@Shakespeare, and CNN@HAR, compared with *Oort*, around 47.4%, 7.6%, 11.7% and 53.0% compared with *AutoFL*, and around 50.5%, 16.9%, 15.2% and 54.2% compared

TABLE II
Performance Comparison: Dropout Ratio (DR), Overall Latency (OL) and Overall Energy Consumption (OEC) of FL Training with Non-i.i.d. Data to Reach the Target Testing Accuracy.

Task Type	CV						NLP			HAR		
Local Model	CNN@MNIST			CNN@CIFAR10			LSTM@Shakespeare			CNN@HAR		
Target Accuracy	91.0%			72.2%			50.3%			89.3%		
Methods	DR (%)	OL (h)	OEC (kJ)	DR (%)	OL (h)	OEC (kJ)	DR (%)	OL (h)	OEC (kJ)	DR (%)	OL (h)	OEC (kJ)
Random	7.0	4.9	1137.5	18.0	22.4	4349.5	38.0	32.8	5031.2	28.0	6.5	1538.4
Oort	46.0	4.1	1230.3	46.0	13.1	4750.5	66.0	19.0	5217.1	85.0	5.6	1561.9
AutoFL	37.0	8.1	1069.5	25.0	28.7	3911.9	53.0	34.1	4835.4	55.0	9.1	1498.6
REAFI	0.0	2.0	562.8	0.0	7.2	3613.6	2.0	11.8	4267.7	0.0	3.0	704.7

TABLE III
Summary of REWAFI Performance Improvement over REAFI and REAFI + LUPA

Task Type	CV				NLP		HAR	
Local Model	CNN@MNIST		CNN@CIFAR10		LSTM@Shakespeare		CNN@HAR	
Target Accuracy	91.0%		72.2%		50.3%		89.3%	
Methods	OL (h)	OEC (kJ)	OL (h)	OEC (kJ)	OL (h)	OEC (kJ)	OL (h)	OEC (kJ)
REAFI	2.0	562.8	7.2	3613.6	11.8	4267.7	3.0	704.7
REAFI+LUPA	1.7	473.6	7.0	3503.7	11.4	4004.8	2.3	603.5
REWAFI	1.3	357.6	6.8	3203.2	10.9	3716.1	2.1	569.7

with *Random*. The rationale behind it is that SOTA PS utility function designs (i.e., *Random*, *Oort*, or *AutoFL*) are not aware of mobile devices' residual energy, and thus frequently select the mobile devices with good learning contributions and small latency/energy consumption at early training stages, which may drain up their energy. Thus, at the late stages of FL training, those depleted mobile devices are not able to participate. Their absence will slow down FL convergence, in particular for non-i.i.d. cases, and thus increase overall latency/energy consumption.

B. Further Performance Improvement by the REWA Local Computing Policy

Under the proposed *REAFI* PS utility function framework, we further evaluate the potential performance improvement brought by REWA local computing policy in terms of overall latency and overall energy consumption for FL training.

As shown in Table III, compared with *REAFI* (i.e., the best PS design in last subsection), when FL training reaches the target accuracy, *REWAFI* reduces around 35.0%, 5.6%, 7.6% and 30.0% latency for learning tasks CNN@MNIST, CNN@CIFAR10, LSTM@Shakespeare and CNN@HAR, and 36.5%, 11.4%, 12.9% and 19.2% energy consumption. Different from *REAFI*'s fixed local computing policy, the performance improvement comes from *REWAFI*'s REWA local computing policy, which increases global statistically utility, latency utility, and local device's energy utility, and helps to speed up convergence and reduce overall latency/energy consumption.

Compared with *REAFI+LUPA*, when FL model reaches target accuracy, *REWAFI* reduces around 23.5%, 2.9%, 4.4% and 8.7% latency for learning tasks CNN@MNIST, CNN@CIFAR10, LSTM@Shakespeare and CNN@HAR, and

24.5%, 8.6%, 7.2% and 5.6% energy consumption. Although *REAFI+LUPA* employs the AdaH policy [23], its performance is not as good as *REWAFI* because *REAFI+LUPA*'s local computing policy fails to consider the impacts of devices' heterogeneous wireless transmissions, and ignores the trade-off between learning performance benefits/statistical utility and the energy cost/energy utility of increasing H . Purely reducing the number of communication rounds while keeping increasing H without stopping criteria will inevitably increase the overall latency/energy consumption. By contrast, *REWAFI* has awareness of wireless transmission heterogeneity and develops an energy utility aware stopping criterion for increasing H , and exploit them to improve the latency/energy efficiency of FL training over mobile devices.

C. REWAFI's Staleness Analysis

To analyze *REWAFI*'s self-contained staleness solution, we employ two types of mobile devices: Xiaomi 12S and Honor 70, representing high-end and low-end mobile devices, respectively. Figure 5 presents the changes of H with high-end and low-end mobile devices, which shows the growth frequency, increment of H and the final stopping value. According to Eqn. (3), H increases only if the device is selected, and the increment of H is determined by the wireless aware local computing policy. Besides, the saturated value of H is determined by Eqn. (4).

Figure 5(a) shows the changes of H for mobile devices with different initial battery energy levels (i.e., high, medium and low). Here, we configure the same transmission rate for the same type of mobile devices, i.e., Xiaomi 12S with the averaged 5G uplink transmission rate of 79.6 Mbps, and Honor 70 smartphones with the averaged 5G uplink transmission rate of 45.0 Mbps. As shown in the figure, for the same type of

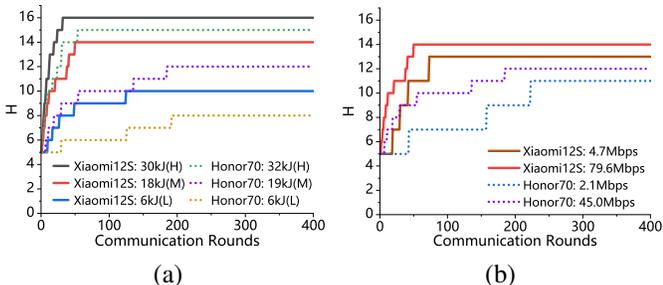


Fig. 5: REWAFI’s staleness solution. (a) different initial energy; (b) different tx. rates (CNN@MNIST).

mobile devices, the H ’s saturated values of the mobile devices with high initial energy are larger than those of devices with low initial energy. That is consistent with Eqn. (4) and the energy utility aware stopping criterion analysis in Sec. III-B2. Moreover, as for different types of mobile devices with the same/similar level of initial battery energy, H ’ growth frequency is high for high-end mobile devices (i.e., Xiaomi 12S) at the beginning of the training process. It is because high-end mobile devices have low energy consumption/latency, leading to frequent selections of these mobile devices for participation. However, in the later stages of training, the H growth frequency of low-end mobile devices (i.e., Honor 70) is higher than that of high-end mobile devices. Similar to the analysis in Sec. III-D, as the high-end mobile device is consecutively selected, its H increases and PS utility keeps decreasing until it is less than that of the low-end mobile device. Then, the low-end mobile device will be selected to participate in the FL training. The results demonstrate that the proposed REWAFI addresses the staleness issue in a self-contained manner.

Figure 5(b) presents the changes of H for mobile devices with different wireless transmission rates. First, we have similar observations as those in Fig. 5(a) for high-end and low-end devices. Besides, for the same type of devices, high transmission rate leads to a small H increment in each communication round, and a high growth frequency of H in the early stages of training. That is because REWAFI employs a wireless aware local SGD computing strategy, which makes mobile devices with higher transmission rate have a smaller increase in H , so that its computing latency and energy consumption can be reduced. Thus, the mobile devices with higher transmission rate has a better chance to be selected, which increases their growth frequency of H at the early stages of training. However, the utility of devices with higher transmission rate decreases as H increases. Once it falls below the utility of mobile devices with lower transmission rate, the mobile device with low transmission rate is then selected to participate. Therefore, in the later stages of training, H of the mobile devices with lower transmission rate grows due to being selected. This is consistent with our analysis of REWAFI’s wireless aware local computing policy in Sec. III-B and self-contained staleness solution in Sec. III-D.

Then, we demonstrate the advantages of REWAFI’s stal-

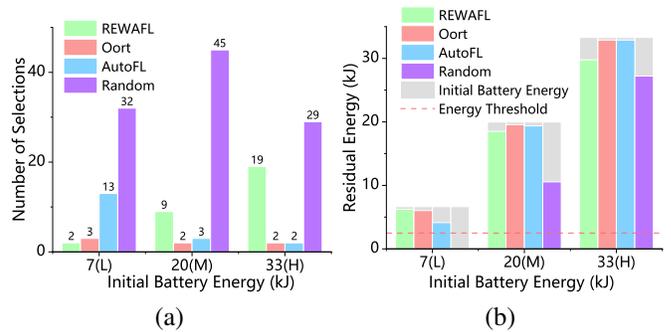


Fig. 6: Performance comparison of staleness solutions in different PS approaches (CNN@MNIST, Honor Play 6T, 5G-0.64Mbps). (a) the number of selections, and (b) energy consumption and residual energy.

ness solution by comparing the number of selections and residual energy for the low-end mobile devices (i.e., Honor Play 6T smartphones) with low transmission rates (i.e., averaged 5G uplink transmission rate of 0.64Mbps) under different PS designs. The results are shown in Fig. 6. Compared with *Random*, *Oort* and *AutoFL*, the low-end mobile devices have a better chance to be selected for participation in training without depleting their energy in REWAFI. It indicates that REWAFI can address the staleness issue in a self-contained manner while outperforming SOTA PS designs and/or their isolated staleness solutions.

D. Sensitivity Analysis

To analyze REWAFI’s sensitivity to different parameter settings, we further evaluate the performance of REWAFI with different α and β values. Figure 7 shows that REWAFI outperforms its counterparts across different α and β values. As shown in Fig. 7(a), a larger α makes the system focus more on latency utility, which allows the global model to reach the same learning performance in a smaller latency. Correspondingly, a larger β makes the system focus more on improving the energy utility, and allows the global model to reach the same test accuracy with less energy consumption as shown in Fig. 7(b). Besides, as shown in Fig. 7(c), a larger β leaves more residual energy for the frequently selected high-end devices (i.e., Xiaomi 12S) and less residual energy for the infrequently selected low-end devices (i.e., Honor Play 6T), which balances the energy consumption among participating mobile devices. Figure 7 indicates that α and β can serve as scaling coefficients to associate/balance different utilities in terms of accuracy, latency and energy efficiency.

E. Impacts of Data Heterogeneity

We further evaluate the performance of REWAFI with different non-i.i.d levels of training data samples in CNN@MNIST case. The results are also applicable to other models and datasets. As shown in Table IV, REWAFI has the best performance in terms of dropout ratio. For non-i.i.d. cases (i.e., $\lambda = 0.8$ and $\lambda = 1$), REWAFI has a better overall latency and overall energy consumption to reach the

TABLE IV
Performance Evaluation under Data Heterogeneity.

Local Model	CNN@MNIST								
Data Heterog.	$\lambda = 0$			$\lambda = 0.8$			$\lambda = 1$		
Target Accuracy	97.0%			91.0%			89.9%		
Methods	OL (h)	OEC (kJ)	DR (%)	OL (h)	OEC (kJ)	DR (%)	OL (h)	OEC (kJ)	DR (%)
Random	3.1	803.8	9.0	4.9	1137.5	7.0	5.0	1215.6	25.0
Oort	1.3	650.1	24.0	4.1	1230.3	46.0	5.0	1451.8	48.0
AutoFL	1.6	411.2	5.0	8.1	1069.5	37.0	8.7	1118.2	38.0
REWAFL	2.2	624.9	0.0	1.3	357.6	0.0	2.6	650.1	0.0

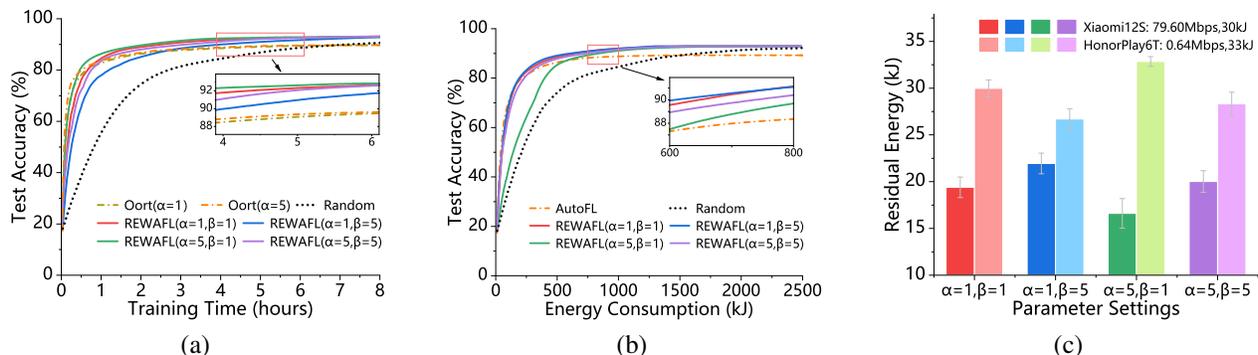


Fig. 7: REWAFL's learning performance and efficiency with different α and β values (CNN@HAR, $\lambda = 0.8$). (a) Test accuracy v.s. training time, (b) Test accuracy v.s. energy consumption, and (c) Residual energy.

target accuracy. For i.i.d. data distribution case (i.e., $\lambda = 0$), we find that *REWAFL* has similar performance to *Oort* in terms of overall latency to reach the target accuracy, and similar performance to *AutoFL* in terms of overall energy consumption. The potential reason is that the impact of device dropouts on convergence rate is not significant due to the similar contribution of local data from different devices to the global model. Unsurprisingly, even for i.i.d. case, *REWAFL* can still maintain the lowest dropout ratio, which reserves enough residual energy at mobile devices for their mandatory operations while enjoying the global learning results. Besides, we find that the FL system with non-i.i.d. data has a higher latency, energy consumption and dropout ratio than that with i.i.d. data to achieve the same accuracy goal.

VI. CONCLUSION

In this paper, we have developed a residual energy and wireless aware PS design for efficient FL training over mobile devices (*REWAFL*). We have observed that most existing participation selection approaches are (i) unaware of heterogeneous residual battery energy, and (ii) unaware of heterogeneous wireless transmission rates among candidate/participating mobile devices. To address these challenges, *REWAFL* introduces a novel PS utility function considering statistical utility, global latency utility and local energy utility. Under the proposed PS utility function framework, *REWAFL* further provides a residual energy and wireless aware local computing policy to improve FL efficiency. Moreover, *REWAFL* buries the staleness solution into its utility function and local computing policy development without incurring any

additional mechanisms. Extensive experimental results have demonstrated the effectiveness of the proposed *REWAFL*, and its efficiency superiority over peer PS designs under various learning tasks, DNN models, datasets for FL training over mobile devices.

REFERENCES

- [1] L. Li, D. Shi, R. Hou, H. Li, M. Pan, and Z. Han, "To talk or to work: Flexible communication compression for energy efficient federated learning over heterogeneous mobile edge devices," in *IEEE INFOCOM 2021 - IEEE Conference on Computer Communications*, 2021, pp. 1–10.
- [2] G. Zhu, D. Liu, Y. Du, C. You, J. Zhang, and K. Huang, "Toward an intelligent edge: Wireless communication meets machine learning," *IEEE Communications Magazine*, vol. 58, no. 1, pp. 19–25, 2020.
- [3] W. Guo, R. Li, C. Huang, X. Qin, K. Shen, and W. Zhang, "Joint device selection and power control for wireless federated learning," *IEEE Journal on Selected Areas in Communications*, vol. 40, no. 8, pp. 2395–2410, 2022.
- [4] Y. J. Cho, J. Wang, and G. Joshi, "Client selection in federated learning: Convergence analysis and power-of-choice selection strategies," *ArXiv*, vol. abs/2010.01243, 2020.
- [5] A. D. P. Team, "Learning with privacy at scale differential," in *Apple Machine Learning Journal*, 2017.
- [6] R. C. Geyer, T. J. Klein, and M. Nabi, "Differentially private federated learning: A client level perspective," 2019. [Online]. Available: <https://openreview.net/forum?id=SkVRTj0cYQ>
- [7] F. Hartmann, S. Suh, A. Komarzewski, T. D. Smith, and I. Segall, "Federated Learning for Ranking Browser History Suggestions," *arXiv e-prints*, Nov. 2019.
- [8] F. Lai, J. You, X. Zhu, H. V. Madhyastha, and M. Chowdhury, "Sol: Fast distributed computation over slow networks," in *17th USENIX Symposium on Networked Systems Design and Implementation (NSDI 20)*. Santa Clara, CA: USENIX Association, Feb. 2020, pp. 273–288.
- [9] T. Li, A. K. Sahu, A. Talwalkar, and V. Smith, "Federated learning: Challenges, methods, and future directions," *IEEE Signal Processing Magazine*, vol. 37, no. 3, pp. 50–60, 2020.

- [10] D. C. Nguyen, M. Ding, P. N. Pathirana, A. Seneviratne, J. Li, and H. V. Poor, "Federated learning for internet of things: A comprehensive survey," *arXiv preprint arXiv:2104.07914*, 2021.
- [11] D. Ye, R. Yu, M. Pan, and Z. Han, "Federated learning in vehicular edge computing: A selective model aggregation approach," *IEEE Access*, vol. 8, pp. 23 920–23 935, 2020.
- [12] F. Lai, X. Zhu, H. V. Madhyastha, and M. Chowdhury, "Oort: Efficient federated learning via guided participant selection," in *15th USENIX Symposium on Operating Systems Design and Implementation (OSDI 21)*. USENIX Association, Jul. 2021, pp. 19–35.
- [13] J. Ren, Y. He, D. Wen, G. Yu, K. Huang, and D. Guo, "Scheduling for cellular federated edge learning with importance and channel awareness," *IEEE Transactions on Wireless Communications*, vol. 19, no. 11, pp. 7690–7703, 2020.
- [14] T. Chen, G. B. Giannakis, T. Sun, and W. Yin, "Lag: Lazily aggregated gradient for communication-efficient distributed learning," in *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, ser. NIPS'18. Red Hook, NY, USA: Curran Associates Inc., 2018, p. 5055–5065.
- [15] A. Katharopoulos and F. Fleuret, "Not all samples are created equal: Deep learning with importance sampling," in *International Conference on Machine Learning*, 2018.
- [16] F. Salehi, P. Thiran, and L. E. Celis, "Coordinate descent with bandit sampling," in *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, ser. NIPS'18. Red Hook, NY, USA: Curran Associates Inc., 2018, p. 9267–9277.
- [17] H. Wu and P. Wang, "Node selection toward faster convergence for federated learning on non-iid data," *IEEE Transactions on Network Science and Engineering*, vol. 9, no. 5, pp. 3099–3111, 2022.
- [18] G. Zhu, Y. Wang, and K. Huang, "Broadband analog aggregation for low-latency federated edge learning," *IEEE Transactions on Wireless Communications*, vol. 19, no. 1, pp. 491–506, 2020.
- [19] M. M. Amiri, D. Gündüz, S. R. Kulkarni, and H. V. Poor, "Convergence of update aware device scheduling for federated learning at the wireless edge," *IEEE Transactions on Wireless Communications*, vol. 20, no. 6, pp. 3643–3658, 2021.
- [20] Y. G. Kim and C.-J. Wu, "Autofl: Enabling heterogeneity-aware energy efficient federated learning," in *MICRO-54: 54th Annual IEEE/ACM International Symposium on Microarchitecture*, 2021, pp. 183–198.
- [21] R. Chen, Q. Wan, X. Zhang, X. Qin, Y. Hou, D. Wang, X. Fu, and M. Pan, "EEFL: High-speed wireless communications inspired energy efficient federated learning over mobile devices," in *Proceedings of the 21st Annual International Conference on Mobile Systems, Applications and Services*, ser. MobiSys '23. New York, NY, USA: Association for Computing Machinery, 2023, p. 544–556.
- [22] Y. Li, X. Qin, H. Chen, K. Han, and P. Zhang, "Energy-aware edge association for cluster-based personalized federated learning," *IEEE Transactions on Vehicular Technology*, vol. 71, no. 6, pp. 6756–6761, 2022.
- [23] F. Haddadpour, M. M. Kamani, M. Mahdavi, and V. R. Cadambe, *Local SGD with Periodic Averaging: Tighter Analysis and Adaptive Synchronization*. Red Hook, NY, USA: Curran Associates Inc., 2019.
- [24] Monsoon-solutions, "High voltage power monitor." <http://www.monsoon.com/LabEquipment/PowerMonitor/>, December 2019.
- [25] I. Fette and A. Melnikov, "The websocket protocol," <https://www.hjp.at/doc/rfc/rfc6455.html>, RFC 6455, IETF. Accessed April 4, 2021.
- [26] D. J. Beutel, T. Topal, A. Mathur, X. Qiu, T. Parcollet, P. P. B. de Gusmao, and N. D. Lane, "Flower: A friendly federated learning framework," *arXiv preprint arXiv:2007.14390*, 2021.
- [27] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y. Arcas, "for federated learning in wireless networks," in *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics*, ser. Proceedings of Machine Learning Research, vol. 54. PMLR, 20–22 Apr 2017, pp. 1273–1282.
- [28] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [29] L. Deng, "The mnist database of handwritten digit images for machine learning research," *IEEE Signal Processing Magazine*, vol. 29, no. 6, pp. 141–142, 2012.
- [30] A. Krizhevsky, G. Hinton *et al.*, "Learning multiple layers of features from tiny images," 2009.
- [31] S. Caldas, S. Meher Karthik Duddu, P. Wu, T. Li, J. Konečný, H. B. McMahan, V. Smith, and A. Talwalkar, "LEAF: A Benchmark for Federated Settings," *arXiv e-prints*, Dec. 2018.
- [32] A. Stisen, H. Blunck, S. Bhattacharya, T. S. Prentow, M. B. Kjærgaard, A. K. Dey, T. Sonne, and M. M. Jensen, "Smart devices are different: Assessing and mitigating mobile sensing heterogeneities for activity recognition," *Proceedings of the 13th ACM Conference on Embedded Networked Sensor Systems*, 2015.
- [33] Z. Yang, M. Chen, W. Saad, C. S. Hong, and M. Shikh-Bahaei, "Energy efficient federated learning over wireless communication networks," *IEEE Transactions on Wireless Communications*, vol. 20, no. 3, pp. 1935–1949, 2021.