# Communication-Efficient Satellite-Ground Federated Learning Through Progressive Weight Quantization

Chen Yang<sup>®</sup>, Jinliang Yuan<sup>®</sup>, Yaozong Wu<sup>®</sup>, Qibo Sun<sup>®</sup>, Ao Zhou<sup>®</sup>, Shangguang Wang<sup>®</sup>, *Senior Member, IEEE*, and Mengwei Xu<sup>®</sup>

Abstract-Large constellations of Low Earth Orbit (LEO) satellites have been launched for Earth observation and satellite-ground communication, which collect massive imagery and sensor data. These data can enhance the AI capabilities of satellites to address global challenges such as real-time disaster navigation and mitigation. Prior studies proposed leveraging federated learning (FL) across satellite-ground to collaboratively train a share machine learning (ML) model in a privacy-preserving mechanism. However, they mostly focus on single unique challenges such as limited ground-to-satellite bandwidth, short connection window, and long connection cycle, while ignoring the completeness of these challenges in deploying efficient FL frameworks in space. In this paper, we propose an efficient satellite-ground FL framework, SatelliteFL, to address these three challenges collectively. Its key idea is to ensure that each satellite must complete per-round training within each connection window. Moreover, we design a progressive blockwise quantization algorithm that determines a unique bitwidth for each block of the ML model to maximize the model utility while not exceeding the connection window. We evaluate SatelliteFL by plugging an implemented FL platform into real-world satellite networks and satellite images. The results show that SatelliteFL highly accelerates the convergence by up to  $2.8 \times$  and improves the bandwidth utilization ratio by up to  $9.3 \times$  compared to the state-of-the-art methods.

Index Terms—In-orbit computing, satellite network, federated learning.

# I. INTRODUCTION

**O** VER the years, Earth observation satellites have consistently provided rich informational support in critical areas such as food security, disaster navigation, climate change, and disease spread [1], [2], [3]. In recent years, significant advancements in satellite technology have substantially reduced satellite deployment costs, leading to the emergence of low Earth orbit (LEO) satellite constellations as a mainstream trend. For example, companies like Planet [4] are now able to collect

The authors are with the State Key Laboratory of Networking and Switching Technology, Beijing University of Posts and Telecommunications, Beijing 100876, China (e-mail: yangc@bupt.edu.cn; yuanjinliang@bupt.edu. cn; yzwu@bupt.edu.cn; qbsun@bupt.edu.cn; aozhou@bupt.edu.cn; sgwang @bupt.edu.cn; mwx@bupt.edu.cn).

Digital Object Identifier 10.1109/TMC.2024.3358804



Fig. 1. (a) LEO satellite constellations comprise many low-Earth orbit satellites, which orbit 500 to 2,000 kilometers from Earth and offer communication services collaborated with the ground stations; (b) Satellite-ground FL framework in space without downloading satellite data to the ground.

over 350 million square kilometers of images daily. Massive imagery and sensor data (10 s–100 s of TB/day) collected by these satellites can effectively enhance machine learning (ML) capabilities to address numerous global challenges. For example, the Sunflower Satellites, with the aid of AI model, can spot wildfires in less than 1 minute, which could have been used to avoid the 2019-2020 Australian 'Black Summer' bushfire with over \$70 billion in property and economic losses [5].

Traditional satellite-based machine learning tasks are mostly completed on the ground. Those collected data by LEO constellation, in Fig. 1(a), is transmitted to the ground station across satellite-to-ground link, and used to train ML models with a powerful computation cluster. However, this paradigm is becoming increasingly infeasible for the following reasons: 1) downloading the raw data collected from satellites would raise substantial overhead to the satellite-ground link. Nevertheless, the current satellite-ground link bandwidth does not even support the download of all the collected data. 2) sharing high-resolution Earth observation images, such as rare natural resources and significant economic activities may not always be feasible due to regulatory restrictions imposed by different countries [6].

There have been several efforts that dedicate to deploying in-orbit computing to address these challenges, which is motivated by the improvement of satellite computation capabilities. OEC [7] is a proposed orbital edge computing platform to support various computing tasks on satellites so that those collected data can be processed locally. SmallSats [8] aims to bring

See https://www.ieee.org/publications/rights/index.html for more information.

Manuscript received 28 February 2023; revised 11 December 2023; accepted 15 January 2024. Date of publication 26 January 2024; date of current version 6 August 2024. This work was supported in part by the National Key R&D Program of China under Grant 2021ZD0113001, in part by Shenzhen Science and Technology Innovation Commission Free Exploring Basic Research Project under Grant 2021Szvup011, and in part by the National Natural Science Foundation of China under Grants U21B2016 and 62032003. Recommended for acceptance by L. Mottola. (*Corresponding author: Qibo Sun.*)

<sup>1536-1233 © 2024</sup> IEEE. Personal use is permitted, but republication/redistribution requires IEEE permission.

efficient inference and training of DNNs in space by introducing lightweight hardware accelerators and designing compact ML algorithms. However, it also faces two key challenges with the deployment of on-board training. First, centralized ML on a single satellite with limited training data is hard to obtain an advanced model, due to privacy concerns when exchanging sensitive imagery data across LEO satellites [2]. Second, considering the memory and power cost of inter-satellite communication, cross-satellite distributed ML in space is currently infeasible in LEO satellite constellations [9].

Our goal in this work is to address these challenges by designing an efficient FL framework that enables satellites to collaboratively train a shared ML model with the coordination of ground stations. Our system is designed with several steps as shown in Fig. 1(b): (1) the global model is dispatched to those connected satellites through ground-to-satellite link; (2) connected satellites run on-board training to update local models without sharing their raw data; (3) these models are transmitted to the ground station through the satellite-to-ground link for global model aggregation. These steps form per-round FL training and iterate with another batch of connected satellites in the next round until the global model converges. Some prior literature [2], [10], [11], [12] also introduces FL framework across satellites and ground stations to enable efficient AI applications in space. Most of them focus on tackling the straggler and staleness problems of synchronous or asynchronous FL algorithms, due to the system heterogeneity of satellites [2], [13]. Another important problem in this satellite-ground FL framework is the slow model convergence caused by limited ground-to-satellite bandwidth [10], [12]. However, it is still fundamentally challenging to apply existing solutions to this environment when considering all the unique challenges introduced by orbital dynamics of LEO satellites:

- Limited ground-to-satellite bandwidth: Dispatching the global model to connected satellites through the ground-to-satellite link is a key step in this environment to guarantee the success of FL training. However, the ground-satellite link is asymmetric. Specifically, the ground-to-satellite link bandwidth is extremely limited, often only about 200 Kbps [14]. Thus, even if the satellite-to-ground download bandwidth is as high as a few Gbps, the entire FL training speed will be severely constrained by the ground-to-satellite link.
- Short connection window: Considering the orbital dynamics of LEO satellites and limited deployment of ground stations on Earth, each satellite tends to stay a short time (≤ 10 minutes [14]) in the coverage of a particular ground station before flying away. This causes the per-round FL training to be interrupted because it is often impossible to complete the above steps in such a short connection window. Especially, the first step of per-round FL training across the limited ground-to-satellite link makes this interruption more frequent.
- *Long connection cycle:* Furthermore, each satellite would only get several connection periods with ground stations in one day, and it often takes over 6 hours [14] for any two connection periods. Therefore, once a satellite fails to

transmit its model update to the ground in the assigned connection window, it takes excessive time to wait till the next connection window and the update is likely to be expired.

In this paper, we present an efficient satellite-ground FL framework, namely SatelliteFL, that leverages progressive weight quantization to compress the communication data size across ground-to-satellite links. (Quantization involves representing model parameters with fewer bitwidth to reduce model size while maintaining acceptable neural network model performance degradation.) We first highlight the importance of ensuring that each satellite must complete the per-round FL training within each connection window to tackle the above three challenges. Then, we formulate this per-round FL training as a delay-constrained optimization problem to maximize the model utility while not exceeding the connection window. However, solving this optimization problem is also challenging due to the complex trade-off between model utility and constrained connection window caused by the limited decision space of quantization bitwidth. Finally, we propose a progressive block-wise quantization (ProBQ) algorithm that fine-grained quantizes each block of model with a unique quantization bitwidth, which can solve the problem and achieve near-optimal performance. The rationale is that we cherry-pick a proper quantization bitwidth for each block of ML model, instead of the whole model; it makes full use of the ground-to-satellite link bandwidth therefore retains more weight information to improve the model accuracy.

We evaluate SatelliteFL on our implemented FL platform plugged in real-world satellite networks. Extensive experiments are performed on classical ML models and open-source satellite imagery data, and we will also open-source our platform once published. The results show that, SatelliteFL can accelerate the convergence by  $1.8 \times ($ on average, up to  $2.8 \times )$  and improve the bandwidth utilization ratio by  $4.5 \times ($ on average, up to  $9.3 \times )$ compared to the state-of-the-art satellite-ground FL methods.

We make the following contributions:

- We formulate the collaborative satellite-ground training problem, with ground-to-satellite bandwidth, short connection window and long connection cycle, as a delayconstrained optimization problem. To tackle these challenges, we introduce an efficient satellite-ground FL framework that each satellite should complete per-round FL training within each connection window.
- We propose a progressive block-wise quantization algorithm to maximize the model utility while not exceeding the connection window.
- We validate the effectiveness of our solution with realworld satellite networks and satellite imagery dataset, and show that it significantly accelerates FL training and improve the bandwidth utilization ratio over the state-of-theart methods.

#### II. BACKGROUND AND MOTIVATIONS

In this section, we first introduce the communication model of satellite-ground dynamic links and the system model of FL at LEO satellites. Then, we formulate the optimization model of



Fig. 2. LEO satellite operates in orbits with two angles and its connection window is related to the signal coverage of ground station. (1) The location of the satellite with respect to a ground station on Earth is defined by the two angles azimuth and elevation; (2) The connection window of the satellite is determined by its orbit and the signal coverage of ground station on Earth, which also indicates the connection window of this satellite at a time.

deploying FL framework across satellites and ground stations based on the unique communication model. Finally, we highlight the key obstacles and the limitations of existing solutions to enable communication-efficient FL training in this satellite-ground cooperative environment.

#### A. Communication Model of Satellite-Ground Link

In an Earth observation scenario, we first illustrate the background of LEO satellites related to the dynamic communication link. Then, we introduce two important definitions of *connected satellite* and *connection window*.

In Fig. 2, there are two observation angles: azimuth and elevation that decide the location of a satellite with respect to a ground station on Earth. The former denotes the angle between the North pole, measured clockwise around the ground station's horizon, and the satellite, while the latter denotes the vertical angular distance between the satellite and the ground station's local horizon. Specifically, when the satellite passes directly over the ground station, the elevation angle reaches its maximum value of 90 degrees. Due to the satellite's low orbit, it often takes only a few minutes for a satellite with a low elevation ( $\beta_1$ in Fig. 2) from entering the signal coverage of a ground station to leaving it. But, with a higher elevation ( $\beta_2 > \beta_1$ ), this satellite gets closer to Earth's surface, thus delivering a longer connection window as shown in Fig. 2. These observations reveal that: (1) satellites aren't always connected to ground station, nor are they intended to be; (2) and even when they do, each connection window is short.

Consider a collection LEO satellites S and ground stations G for Earth observation. We introduce a continuous wall clock time t and a discrete time index  $i \in \{0, 1, 2, ...\}$  with each adjacent time indexes having  $\tau$  wall clock time interval. We denote the wall clock time interval from  $i\tau$  to  $(i + 1)\tau$  as [i]. If a satellite  $s \in S$  is captured by a ground station  $g \in G$  at any time  $t \in [i]$ , their communication link would be established. We define the satellite s as a *connected satellite* in time interval [i]. As

s moves out of signal coverage of g, the link breaks, and s is no longer a *connected satellite*. We define the duration of this link as *connection window*, which is denoted as  $\tau_d$ , and we have  $\tau_d \leq \tau$ .

In this work, we assume that all the ground stations work as a powerful computing cluster without considering the transmission delay across those geographically-distributed ground stations. This assumption is theoretically sound because the fabric connectivity among them is always well-provisioned and much faster than the satellite-ground links [15]. Therefore, a satellite connected to anyone ground station can be considered to be connected to a cluster with all the ground stations  $\mathcal{G}$ .

## B. Optimization Model of Satellite-Ground FL

LEO satellites in this Earth observation scenario have collected a lot of imagery and sensor data to train a global model. Due to data privacy and bandwidth limitation concerns, a satellite-ground cooperative FL framework is proposed to collaboratively learn this model with on-board training and on-ground coordination in Fig. 1(a). We follow the classical synchronized stochastic gradient descent (SGD) method to formulate this cooperative FL process.

Consider K satellites in S and a cluster with ground stations  $\mathcal{G}$ . For each satellite  $k \in \{1, 2, ..., K\}$ , it collects and stores an imagery dataset  $\mathcal{D}_k$ . These satellites aim to collaboratively learn a global model  $\omega$  by minimizing a global objective function  $F(\omega)$  as follows:

$$\min_{\omega} \left\{ F(\omega) = \sum_{k=1}^{K} \frac{n_k}{n} f_k(\omega) \right\},\tag{1}$$

where  $n_k = |\mathcal{D}_k|$  denotes the samples of training data at satellite k,  $f_k$  denotes the local objective function of satellite k, and we have  $n = \sum_{k=1}^{K} n_k$ . The global model  $\omega$  is maintained and advanced on the cluster of ground stations. We introduce the most popular FL protocol, FedAvg [16], that performs local SGD algorithm on satellites and aggregates those models of satellites on the ground to advance global model. In particular, we only randomly select  $K_s$  satellites from  $K_c$  connected satellites in each round, where  $K_c \leq K$ . The local update on each selected satellite  $k \in K_s$  can be formulated as

$$\omega_{r+1}^k = \omega_r - \eta \nabla f_k(\omega_r), \qquad (2)$$

where  $\eta$  denotes the hyper-parameter learning rate and r denotes the global training round across satellites and ground stations. The global update is to aggregate the received models from  $K_s$ selected satellites, which can be formulated as

$$\omega_{r+1} = \sum_{k \in K_s} \frac{n_k}{n} \omega_{r+1}^k. \tag{3}$$

Consider the wall clock time during per-round of FL training, which includes three key steps: 1) dispatching global model to each satellite; 2) on-board training to update local model; 3) sending local model to ground station. We denote the time of above three steps as  $t_{g,s}$ ,  $t_c$  and  $t_{s,g}$ , respectively. In this satellite-ground FL framework, it is highly likely that some satellites' per-round of FL training will be forced to stop due



Fig. 3. LEO satellite's three distinctive characteristics on satellite-ground connection [17], [18], [19].



to the disconnection of satellite-ground links. We describe the procedures of this framework in Algorithm 1.

# C. Challenges of Deploying Efficient FL in Space

As illustrated above, those participating satellites in each round are selected from those *connected satellites* with determined *connection window*, which makes the satellite-ground communication crucial in the FL optimization. However, satellite-ground communication substantially differs from traditional communication on the ground in three aspects.

Limited Ground-to-Satellite Bandwidth. In Earth observation scenario, the communication links from the ground station to satellite are typically designed to transmit control signaling, thus resulting in very low bandwidth. In addition, expanding this ground-to-satellite bandwidth is often highly expensive due to the huge monetary cost and serious heat dissipation when replacing Earth-dial links with satellite communication links [20], [21], [22]. Therefore, the ground-to-satellite bandwidth is often hundreds of Kbps as shown in Fig. 3(a), although the ground station today supports Gbps satellite-to-ground bandwidth [1], [14], [23]. Our preliminary experiments show that transmitting a ResNet18 [24] model to ground station only costs 0.4 seconds, yet it takes nearly one hour for the aggregated model on ground station to be dispatched to connected satellites. The extreme ground-to-satellite bandwidth breaks down the widely-accepted assumption of symmetric upload/download bandwidth in traditional FL.

Short Connection Window. Traditional FL on ground runs in a well-connected distributed system with a number of highperformance network infrastructures to guarantee their superior connection conditions. However, due to the orbital dynamics of LEO satellites and limited deployment of ground stations on Earth, the established communication link usually only lasts for a short time. For example, it is often several minutes at a time, and in most cases no more than ten minutes as shown in Fig. 3(b). Our preliminary experiments on ResNet18 model show that 80% of the satellites fail to complete the per-round FL training under these stringent connection windows.

Long Connection Cycle. Given the frequent failure in one connection window mentioned above, many satellites would attempt to complete the previously failed steps in the next connection window. However, apart from the short connection window, the orbital dynamics of satellites and limited deployment of ground stations also lead to a very long connection cycle of satellite-ground communication links. Fig. 3(c) shows that each satellite would only get several connection periods with ground stations in one day, and it often takes many hours for any two connection periods. Therefore, this makes it infeasible to pick up the failed steps in the next connection window, due to the expired weights with large staleness.

# D. Limitations of Existing Solutions

The aforementioned unique challenges have not been thoroughly studied, leaving us with ample room for improvement. This drives us to delve deeper into how to enable satellites in space with efficient intelligent processing capabilities. Based on existing efforts made from various perspectives, we group related work into three categories:

AI on Satellites. Several categories of work relate to enabling AI capability of LEO satellites. OEC [7] and Tiansuan constellation [25] proposed to build the orbital edge computing platform to support various computing tasks on satellites so that those collected data can be processed in space. Model inference on satellite [8], [26], [27], [28], [29] is the key to improving the image process capacity, which can help to efficiently analyze the massive collected Earth observation image. To tackle the challenges of satellite networks, some early research [30], [31] proposed to execute distributed machine learning in space. Their goal is to enable ML model to work better under a limited computing power and communication bandwidth of satellites. However, considering the memory and power cost of intersatellite communication, cross-satellite distributed ML in space is currently infeasible in LEO satellite constellations.

Federated Learning in Space. Federated learning [32], [33], a collaborative machine learning paradigm, has been introduced to enable satellites to collaboratively train an ML model with the coordination of ground stations [2], [10], [11], [13], [34], [35]. Most of them focus on tackling the straggler and staleness problems of synchronous or asynchronous FL algorithms, due to the system heterogeneity of satellites [2], [13]. Other research aim to accelerate the slow model convergence caused by limited ground-satellite bandwidth [10], [12], [34], [35]. However, it is still fundamentally challenging to apply existing solutions to this environment when considering all the unique challenges introduced by orbital dynamics of LEO satellites: limited ground-satellite bandwidth, short connection window and long connection cycle. This motivates us to tackle those challenges together to enable communication efficient FL in space.

*Model Quantization for Communication Efficient FL*. Quantization for FL [36], [37] is an approach that allows several devices to update models using low bitwidth gradients, maintaining accuracy while reducing their communication cost. Most existing neural quantization approaches focus on how to prune the redundant gradient information in the training processing, e.g., replacing the default numerical FP32 gradients with INT8 and even INT4 [38], [39], [40], [41], [42], [43], [44], [45]. Instead of contributing a novelty accuracy-first FL quantization training algorithm, our goal is to design a generic system to efficiently support in-orbit satellite training under dynamic ground-satellite connection in reality.

In this paper, our goal is to enable satellites with effective intelligent processing capabilities under aforementioned unique challenges. Notably, some new Earth observation satellites have been launched in recent years with ground-satellite bandwidth up to tens of Mbps [46], but the majority of orbiting satellites, particularly those launched earlier, maintain only hundreds of Kbps bandwidth [47], [48], [49], [50]. Besides, as the scale of a satellite network increases, the actual bandwidth allocation per satellite diminishes. Thus, deploying communication-efficient FL framework across satellites and ground stations to enable satellites with effective intelligent processing capabilities is crucial and meaningful.

# III. DESIGN OF SATELLITEFL FRAMEWORK

This section proposes a satellite-ground FL framework (SatelliteFL) with a progressive block-wise quantization algorithm



Fig. 4. Overview of our SatelliteFL framework with K satellites and a cluster of ground stations. Each satellite has c connection windows to the ground station in one day.

(ProBQ) to improve communication efficiency during FL training across satellites and ground stations. Its key idea is to ensure that each satellite completes per-round FL training within each connection window by progressive weight quantization. Section III-A first illustrates the overview of SatelliteFL. Section III-B1 formulates the objective of SatelliteFL as a delayconstrained optimization problem. Section III-B2 proposes the ProBQ algorithm, which can approximately solve this problem.

# A. SatelliteFL Overview

We show the detailed illustration of our SatelliteFL framework in Fig. 4. Overall, SatelliteFL also adopts a C/S architecture, where a central server served by ground stations maintains and keeps advancing a global model, and the client refers to a collection of LEO satellites. These satellites dynamically establish communication links with ground stations in the process of orbiting Earth.

Each satellite in SatelliteFL is responsible for three tasks: (1) It dequantizes the received integer (INT) model into 32-bit float (FP32) one; (2) The dequantized FP32 model, along with the local data, are used to obtain an advanced model by onboard training; (3) It maintains the connection information, such as link bandwidth and connection window. The updated FP32 model and connection information are transmitted to the ground stations for global model updating. Considering the relatively fixed satellite orbit and ground station deployment, the daily connection of each satellite is predictable, but the connection's bandwidth and window can only be determined once the link is established. Section III-B specifically analyzes the impact of these dynamic connections on the satellite-ground FL training.

Ground station serves as two key roles: aggregator and quantizer. The aggregator is responsible for aggregating the received FP32 models from satellites, and obtaining an advanced global model. The quantizer is responsible for quantizing the global model into multiple INT models, which are adaptive to the connection information in the satellite profiles. Its goal is to improve the model utility under the communication constraints, which is formulated as a delay-constrained optimization problem (details in Section III-B1). Then, it adopts a progressive block-wise quantization algorithm that fine-grained quantizes each block of model with a unique quantization bitwidth to solve this problem (details in Section III-B2). Finally, these quantized INT models are dispatched to corresponding satellites through the ground-to-satellite link.

#### B. Delay-Constrained Optimization Problem

1) Problem Formulation: Consider a connected satellite participating in this round of FL training, the ground station obtains its connection window  $\tau_d$ , ground-to-satellite bandwidth  $b_{g,s}$ and satellite-to-ground bandwidth  $b_{s,g}$ . The ground station also maintains a global model  $\omega$  and launches a round of FL training with three steps. First, it would quantize this model into INT format to accelerate the communication from ground station to satellite, which can be formulated as

$$\omega(\gamma) = quantize(\omega, \gamma), \gamma \in \mathcal{S},\tag{4}$$

where  $\gamma$  is the quantization bitwidth to represent the INT format model  $\omega(\gamma)$  and  $\mathcal{N}$  is a set of available bitwidth specifications for hardware devices on satellite. Then, the quantized model  $\omega(\gamma)$  with  $\gamma$ -bit INT format is sent to this satellite through the ground-to-satellite link. After receiving this quantized model, the satellite would dequantize this INT model to obtain dequantized model  $\omega'(\gamma)$  with FP32 format, and run on-board training with  $\omega'(\gamma)$ . Finally, the locally updated model with FP32 format is transmitted to the ground station through the satellite-to-ground link. The original model  $\omega$  and dequantized model are represented in FP32 format. We denote the data size of model as  $|\cdot|$ , thus

$$\frac{\gamma}{32} = \frac{|\omega(\gamma)|}{|\omega|}, \text{ and } |\omega'(\gamma)| = |\omega|.$$
(5)

The implication of (5) is that the compression rate of the  $\omega(\gamma)$  to the original  $\omega$  is  $\frac{\gamma}{32}$ , and the data size of dequantized model  $\omega'(\gamma)$  is the same as global model  $\omega$ . So, a smaller  $\gamma$  attributes to fewer communication data, thus accelerating the model transmission from the ground station to the satellites.

Let t denote the per-round time for anyone connected satellite participating in the FL training to complete the above three steps. We do not consider the computation time of model quantization and dequantization, because it only needs to conduct a few scalar multiplications that are much less complex than local model updates. Therefore, we have the per-round time t as follows.

$$t = t_{g,s} + t_c + t_{s,g},$$
 (6)

where  $t_{g,s}$  is the communication time for  $\omega(\gamma)$  to be transmitted from the ground station to the satellite,  $t_c$  is the computation time for the model training on satellite, and  $t_{s,g}$  is the communication time for updated model based on  $\omega'(\gamma)$  to be transmitted from the satellite to the ground station. We calculate the two communication time based on the transmission data size and link bandwidth as follows:

$$t_{g,s} = \frac{|\omega(\gamma)|}{b_{g,s}} = \frac{|\omega| \cdot \gamma}{32 \cdot b_{g,s}},\tag{7}$$

and

$$t_{s,g} = \frac{|\omega'(\gamma)|}{b_{s,g}} = \frac{|\omega|}{b_{s,g}}.$$
(8)

The computation time  $t_c$  is spent updating  $\omega'(\gamma)$  by SGD algorithm, which is related to the hardware devices on satellite and  $\omega$ .

Recall the global objective function  $F(\omega)$  in (1). Our goal is to minimize each  $F(\omega)$  in our SatelliteFL with two key considerations: 1) ensuring that per-round FL training does not exceed the connection window  $\tau_d$  of this satellite in this round; 2) the numerical error between the dequantized  $\omega'(\gamma)$  and the original  $\omega$  without quantization should not exceed the threshold  $\epsilon$  to ensure the global model accuracy. We introduce *weight divergence* between  $\omega'(\gamma)$  and  $\omega$  to quantify their numerical error:  $\frac{|||\omega|| - ||\omega'(\gamma)|||}{||\omega||}$ . So, we formulate this problem

$$\mathbf{P1} : \min_{\gamma} \left\{ F(\omega) = \sum_{k \in \mathcal{S}} f_k(\omega'(\gamma)) \right\}$$
  
s.t. 
$$\frac{|\omega| \cdot \gamma}{32 \cdot b_{g,s}} + t_c + \frac{|\omega|}{b_{s,g}} \le \tau_d,$$
$$\frac{||\omega|| - ||\omega'(\gamma)||}{||\omega||} \le \epsilon,$$
$$\gamma \in \mathcal{N}.$$

Solving P1 is challenging for the following reasons. First, a smaller  $\gamma$  attributes to fewer communication data, thus accelerating the transmission across ground-to-satellite. However, a larger  $\gamma$  is required to represent more information of original  $\omega$ , thus reducing the numerical error introduced by quantization. Second, it is generally impossible to obtain an exact analytical relationship to connect local objective function f with the dequantized  $\omega'(\gamma)$ .

2) Analysis of **P1**: Low bitwidth quantization, while transmitting-friendly, severely limits the expressiveness of the updated information. Updating the local parameters on satellite with such limited information thus may not improve the quality of local models, even making them worse. We first analyze the deficiency of naive uniform quantization to solve **P1**. Then, we dive into the sensitivity of the bitwidth of quantization to model structure, and translate the above problem to **P2**.

Naive uniform quantization is to decide a unique  $\gamma$  to quantize all the weights of global FP32 model, which aims to ensure that each satellite must complete per-round training within each connection window. Based on the estimation of computation time and communication time, we choose as large bitwidth as possible to quantize model, while satisfying the constraint of connection window. However, this scheme leads to a serious deficiency in the model accuracy and bandwidth utilization as shown in Fig. 5. The reason behind is that limited decision space on  $\gamma$  makes it difficult to achieve a complex balance between the constraints of connection window and quantization error.

Therefore, we dive into the impact of model architecture on accuracy performance when deciding different quantized bitwidth  $\gamma$ . Since the classical model architectures tend to have several blocks, conceptualized as analogous to the ventral visual blocks [51]. We then conduct a comprehensive measurement on each block with different low bitwidth, separately. We obtain two key observations from these experimental results that motivate



Fig. 5. Deficiency of naive uniform quantization on model accuracy and bandwidth utilization.



Fig. 6. Comparison of traditional fully FP32 blocks (FP32) and only one quantized block with INT4 bitwidth representation (B1-B6: Block1-Block6).



Fig. 7. Hybrid bitwidth quantization across FL training round with DenseNet121 on two datasets. The deep V regions marked in red represent the lower accuracy when using fully INT4 quantization in this round.

our further design. The settings of measurements are consistent with Section IV-A.

• Observation 1: The impact on the accuracy varies greatly when quantizing different blocks with low bitwidth, blocks in the middle layer are less suitable for low bitwidth quantization. Fig. 6 shows the end-to-end accuracy performance when using 4-bit quantization in each single block, while other blocks in this model use FP32 representation. Most of the blocks suffer accuracy degradation when using low bitwidth quantization. For example, block2, block3, and block4 of DenseNet121 cause up to a 35%–67% accuracy loss, while there is less than 6% accuracy loss when quantizing block1 and block6. Therefore, it is worth that we can improve the efficiency of model accuracy and bandwidth utilization (in Fig. 5) by cherry-picking suitable blocks for low bitwidth quantization.

• Observation 2: High bitwidth quantization on blocks can make up for the accuracy loss caused by low bitwidth quantization quickly. Fig. 7 shows the accuracy performance when using hybrid bitwidth quantization across FL round on all blocks. Note that the four marked regions denoting the lower obtained accuracy with full INT4 quantization. The lower accuracy immediately reverts to a higher level in the next few rounds when turning to INT8 (or INT16, INT32) quantization. The behind reason is that a higher bitwidth representation retains more correct gradient information than a lower one, which can revise the incorrect updating direction and mitigate the accuracy degradation. This motivates us to introduce hybrid bitwidth quantization among different blocks for efficiency improvements during the FL training.

• Implications. In summary, low bitwidth quantization as a strategy to reduce model parameters, make an obvious accuracy degradation. However, some specific blocks will not suffer this serious performance loss, and the high bitwidth quantization can mitigate the accuracy degradation caused by the low bitwidth quantization. To enable a practical scheme with tolerable accuracy degradation and full link utilization, the quantization paradigm needs to be re-architected.

Therefore, we introduce the block-wise decision of  $\gamma$  to **P1**. For anyone model  $\omega$ , we have m blocks

$$\omega = \{\omega_1, \omega_2, \dots, \omega_m\},\tag{9}$$

where  $\omega_i$  denotes a block of  $\omega$  and  $i \in [1, 2, ..., m]$ . For anyone block  $\omega_i \in \omega$ , we decide a unique bitwidth  $\gamma_i$  to quantize this block with FP32 format into INT format with  $\gamma_i$  bit representation. Then, we get a new decision vector  $\Gamma$ 

$$\Gamma = \{(\omega_1, \gamma_1), (\omega_2, \gamma_2), \dots, (\omega_m, \gamma_m)\}.$$
(10)

Consider the new quantization decision on  $\omega$ , we re-calculate the communication time  $t_{q,s}$  from ground station to satellite

$$t_{g,s} = \frac{\sum_{i=1}^{m} |\omega_i(\gamma_i)|}{b_{g,s}} = \frac{\sum_{i=1}^{m} |\omega_i| \cdot \gamma_i}{32 \cdot b_{g,s}},$$
 (11)

where  $\omega_i(\gamma_i)$  denotes the quantized block  $\omega_i$  with bitwidth decision  $\gamma_i$ . We also dequantize block  $\omega_i(\gamma_i)$   $(i \in [1, 2, ..., m])$  one by one, and organize them to obtain the dequantized model  $\omega'(\Gamma)$  with FP32 format.

Finally, we translate **P1** into **P2** as:

$$\begin{aligned} \mathbf{P2} : \min_{\Gamma} \left\{ F(\omega) &= \sum_{k \in \mathcal{S}} f_k(\omega'(\Gamma)) \right\} \\ \text{s.t.} \quad \frac{\sum_{i=1}^m |\omega_i| \cdot \gamma_i}{32 \cdot b_{g,s}} \leq \tau_d - t_c - t_{s,g} \\ \frac{|||\omega|| - ||\omega'(\Gamma)|||}{||\omega||} \leq \epsilon, \\ \frac{\forall \gamma_i \in \Gamma, \gamma_i \in \mathcal{N}. \end{aligned}$$

In problem **P2**, the objective is to minimize the each  $F(\omega)$ in our SatelliteFL with two key considerations. Note that this  $F(\omega)$  is controlled by the decision vector  $\Gamma$  in each satellite. One of the key considerations is how to ensure that per-round FL training does not exceed the transmission time  $\tau_d - t_c - t_{s,g}$  of this satellite in current round, and the per-round FL training time is determined by the current decision vector  $\Gamma$  in each satellite

| Algorithm 2: Satellite-Ground FL With ProBQ Design.              |  |  |  |  |  |  |  |
|--|--|--|--|--|--|--|--|
| Ι  | nput   | : initialized global model $\omega_0(\gamma)$ and sorted<br>blocks { $\omega_1, \omega_2, \cdots, \omega_m$ } based on the impact<br>degree, decision of bitwidth<br>$\gamma \in \mathcal{N} = \{b_1, b_2, \cdots, b_n\}$ , ground-to-satellite<br>bandwidth $b_{g,s}$ , set of total satellites $S$ and |  |  |  |  |  |
|  | total training rounds <i>R</i> .   |  |  |  |  |  |  |
| (  | <b>Output:</b> the global model $\omega_r$ .   |  |  |  |  |  |  |
| 1 (  | Ground station executes: // global coordination  |  |  |  |  |  |  |
| 2  | for $r \leftarrow 0$ to R do   |  |  |  |  |  |  |
| 3  |  | $\mathcal{S}_c, \{\tau_d\} \leftarrow \text{connected satellites and their duration}$  |  |  |  |  |  |
| 4  | Sort $S_c$ by its connection duration $\tau_d$ in reverse  |  |  |  |  |  |  |
| 5  | $\overline{\mathcal{S}_s} \leftarrow \text{select top-}K_s \text{ satellites from } \mathcal{S}_c$ |  |  |  |  |  |  |
| 6  |  | for each satellite $k \in S_s$ with duration $\tau_d$ do   |  |  |  |  |  |
| 7  |  | $t_c, t_{s,q} \leftarrow$ estimated training and download time   |  |  |  |  |  |
| 8  |  | $\overline{\Gamma \leftarrow \mathbf{ProBO}(\tau_d, t_c, t_{s,a}, b_{a,s})}$   |  |  |  |  |  |
| 9  |  | Quantize $\omega_r(\gamma)$ to INT model $\omega_r^k(int)$ as $\Gamma$   |  |  |  |  |  |
| 10   |  | Dispatch the $\omega_{\pi}^{k}(int)$ to satellite k  |  |  |  |  |  |
| 10   |  | $\omega_{n+1}^{k} \leftarrow \text{SatelliteUpdate}(\omega_{n}^{k}(int))$  |  |  |  |  |  |
| 11   |  | Update global model $\omega_{r+1}$ as Eq. (3)  |  |  |  |  |  |
| 12 $ProBQ(\tau_d, t_c, t_{s,g}, b_{g,s})$ : // search a feasible |  |  |  |  |  |  |  |
| 13   | De   | efault decision  |  |  |  |  |  |
|  | $\Gamma = \{(\omega_1, b_1), (\omega_2, b_1), \cdots, (\omega_m, b_1)\}$                           |  |  |  |  |  |  |
| 14   | for $i \leftarrow 1$ to $m$ do   |  |  |  |  |  |  |
| 15   |  | if $t_{q,s} = \frac{\sum_{i=1}^{m}  \omega_i  \cdot \gamma_i}{22.b} > \tau_d - t_c - t_{s,q}$ then   |  |  |  |  |  |
| 16   |  | <b>for</b> $i \leftarrow b_1$ <b>to</b> $b_n$ <b>do</b>  |  |  |  |  |  |
| 17   |  | $b_{j} = b_{j+1}$  |  |  |  |  |  |
| 18   |  | Set $\gamma_i = b_j$ for $\omega_i$ and update $\Gamma$  |  |  |  |  |  |
| 19   |  | else   |  |  |  |  |  |
| 20   |  | return Γ   |  |  |  |  |  |
|  | L L  |  |  |  |  |  |  |
| 21 S   | Satell   | <i>iteUpdate(<math>\omega_r^k(int)</math>):</i> // local training  |  |  |  |  |  |
| 22   | De   | equantize $\omega_r^k(int)$ to FP32 model $\omega_r^k$   |  |  |  |  |  |
| 23   | $\overline{\omega_r^k}$  | $_{\pm 1} \leftarrow$ on-board training to update $\omega_r^k$ as Eq. (2)  |  |  |  |  |  |
| 24   | ret  | turn the updated FP32 model $\omega_{r+1}^k$   |  |  |  |  |  |

model. The other key consideration is how to ensure that the numerical error between the dequantized model and original model should not exceed the threshold  $\epsilon$ . Therefore, problem **P2** is a nonlinear integer dynamic optimization problem [52], [53]. To solve **P2**, we need to find out how the value of  $\Gamma$  affects the loss function *F* and the *weight divergence*  $\frac{|||\omega|| - ||\omega'(\gamma)||}{||\omega||}$ . However, it has been proved that problem **P2** is a classical NP-hard problem [54] over  $\Gamma$ , and in reality, algorithm with reduced computation complexity is required to solve **P2**. Meanwhile, considering the relatively limited integer bitwidth supported by satellite hardware, the greedy strategy can be used to solve it well. Thus, we adopt the greedy idea to obtain a suboptimal solution in designing Algorithm **2**.

Intuitively, as long as the bitwidth of quantization is larger, the practical training loss would be closer to the original FP32-based loss and the quantization error would be smaller. Therefore, the key is how to choose the maximum quantization bitwidth. As long as connection window is guaranteed by this decided bit width, the model utility can be improved to the maximum. And the available bitwidth specifications for hardware devices on

satellite are very limited (like INT2, INT4, INT8, etc.), which means the decision space of  $\gamma$  is also small. Therefore, we can adopt a greedy algorithm to solve this problem in a short clock time.

#### C. Progressive Block-Wise Quantization Algorithm

In this section, we propose a progressive block-wise quantization algorithm (ProBQ) that quantizes each block of a model with a suitable bitwidth based on the above observations. This novel design compensates for the quantization error as much as possible on the premise of guaranteeing the model update in each transmission process. Its core technical designs are two-fold: (1) adaptive quantization on fine-grained blocks to achieve a better trade-off between model utility and delay constraints; (2) FP32-guided dequantization on low bitwidth model to enable high precision on-board training. Algorithm 2 describes the workflow of our SatelliteFL with proposed ProBQ algorithm.

We first introduce two key variables in our decision vector to help describe the algorithm:  $\mathcal{N}$  and  $\omega$ . Consider a finite space with *n* positive integers  $\mathcal{N} = \{b_1, b_2, \ldots, b_n\}$  that represents a set of available bitwidth specifications for hardware devices on satellite. In most cases, it supports only several types of bitwidth. (like INT2, INT4, INT8, etc.) We sort this set by the number of bitwidth in reverse, which means that the first element of  $\mathcal{N}$ is the maximum bitwidth to quantize each block. Consider a general ML model  $\omega$  with *m* blocks. Based on the prior profile information that the impact degree on the model accuracy when quantizing different blocks with a low bitwidth as shown in Figs. 6 and 7. We sort model's all blocks as this impact degree from smallest to largest to obtain a sorted list of blocks:  $\omega = {\omega_1, \omega_2, \ldots, \omega_m}$ .

Then, under the satellite-ground cooperative mechanism in this framework, we focus on the introduction of the ProBQ algorithm. To reduce the failure of satellites' FL training in one round, we select  $K_s$  satellites with the longest connection time from the *connected satellites* in each round (Line 4-5). During each satellite's decision phase, we first obtain the transmission time for the quantized INT model in this connected duration by the estimation (Line 7), then we adopt a greedy method to search the feasible solution with the constrained available time  $au_d - t_c - t_{s,g}$  (Line 13-21). Although this method's complexity is very high, the practical execution speed is not slow due to the small number of available satellites, bitwidth supported by the satellites, and the blocks of the model. With the obtained decision vector  $\Gamma$ , global model is quantified into INT format and dispatched to the corresponding satellite. Finally, in order to ensure the training accuracy, each satellite dequantizes the received INT model into FP32 model for on-board training (Line 23-25).

# IV. PERFORMANCE EVALUATION

# A. Experimental Settings

1) Datasets and Models: We adopt a real-world satellite imagery dataset and an open-source image dataset as shown in Table I: Functional Map of the World (fMoW) [55] and

TABLE I DATASETS USED IN EXPERIMENTS FOR ONE TASK: IMAGE CLASSIFICATION. IC (IMAGE CLASSIFICATION), OD (OBJECT DETECTION), IS (INSTANCE SEGMENTATION)

| Dataset   | Task       | Size           | Samples    | Labels |
|-----------|------------|----------------|------------|--------|
| fMoW[55]  | IC, OD, IS | resize 224*224 | 1 million  | 63     |
| L-SUN[56] | IC         | 256*256        | 59 million | 30     |

L-SUN [56] on 3 classic CNN models: DenseNet121 [57], VGG16 [58] and ResNet18 [24]. These three models are popular for image classification, which contains 6, 6, and 5 blocks respectively.

- fMoW is a classic world functional map dataset, which contains more than 1 million high-resolution images with a total of 63 categories. We generate a sub-dataset by random sampling the data for each class, setting 0.2 as the fraction of data to be sampled. In sub-dataset, each high-resolution image is resize to 224\*224.
- L-SUN is a classic image classification dataset, which contains approximately 59 million images with 10 scene categories and 20 object categories, and the fraction of data to be sampled is set as 0.01.

2) Simulated Platform: We have implemented SatelliteFL on a simulated platform atop FedML [59], a popular FL framework with real-world communication and computation simulation. We also followed the prior work [2] to plug all the configurations as an example constellation, named Planet Lab with 12 ground stations and 191 satellites. We follow the prior work [60], [61], [62] to divide the two datasets into non-iid sub-dataset, and assigned them to 191 satellites. Here, we use the cote simulator to obtain the connection information same as FedSpace [2]. Specifically, to establish a clear correspondence between training rounds and the clock time, we set per training round period with  $\tau = 15$  minutes, and there are 96 training rounds in one day. The connected satellites with their connection information in each round are also obtained from this simulation platform. Moreover, to simulate the dynamic variations in the satelliteto-ground link, we attempted to model the ground-to-satellite bandwidth using a Gaussian distribution with a mean of  $b_{q,s}$  and a standard deviation of  $b_{q,s}/3$ . All the experiments are conducted on a Ubuntu 18.04 Linux server with 8 NVIDIA A40 GPUs.

3) Metrics: Apart from the convergence accuracy of the testing data, we also report the following two metrics that closely relate to the satellites. (i) Clock time is the end-to-end training time perceived by the satellites, including multiple rounds of per-round FL training until model convergence. (ii) Bandwidth utilization of ground-to-satellite link refers to the ratio of the valid data size to the total data size that can be transmitted by the ground station during each connection duration, denoted as  $u_ratio$ . The total data size can be calculated as  $b_{g,s} \cdot \tau_d$ , which means if there is enough data required to be transmitted to the satellite throughout the connection duration  $\tau_d$  with the ground-to-satellite bandwidth  $b_{g,s}$ . The valid data size here refers to the size of actual transmitted model parameters and it must be a complete one with all parameters, because part of parameters cannot be used for the satellite's on-board training. So we have

$$u_{\text{ratio}} = \frac{\alpha(|\omega|)}{b_{g,s} \cdot \tau_d}, \text{ and } \alpha(|\omega|) = \begin{cases} |\omega|, & \text{success} \\ 0, & \text{fail.} \end{cases}$$
 (12)

We fix the sum of computation time  $t_c$  and communication time  $t_{s.g}$  as 3 minutes based on the profiled measurement on Jetson TX2 and 1 Gbps satellite-to-ground bandwidth. Note that the average connection window is 6 minutes as shown in Fig. 3(b), and the maximum  $u_ratio$  is about 50% on average.

4) Baselines: We use three baselines in experiments to demonstrate SatelliteFL's benefits are: (i) FedAvg [16]: the traditional FL protocol using fully FP32 format on the ground. (ii) Vanilla SatelliteFL (vanilla): using naive uniform bitwidth quantization in satellite-ground FL framework as described in Analysis of **P1** of Section III-B. (iii) FedSpace [2]: a state-of-the-art satellite-ground FL framework using an adaptive buffer to balance the synchronous and asynchronous training phases.

#### **B.** Experimental Results

1) End-to-End Performance: We show the overall results of SatelliteFL on the two metrics compared to baselines when using  $b_{g,s}$ =400 as the default setting. Fig. 8 shows the time-to-accuracy performance of SatelliteFL compared to three baselines, which demonstrates that our SatelliteFL framework with ProBQ algorithm greatly improves both the convergence speed and the model accuracy. Fig. 9 shows the performance of bandwidth utilization ratio, which reveals that it also highly improves the bandwidth utilization of the ground-to-satellite link.

As shown in Fig. 8, SatelliteFL greatly improves both the convergence speed and model accuracy compared with baselines, Compared to FedAvg with DenseNet121 on fMoW and L-SUN, SatelliteFL obtains 31.6% and 36.9% higher convergence accuracy, and it takes 42.1% and 44.2% fewer days to converge. Moreover, for ResNet18 and VGG16, FedAvg's accuracy is even less than 10% without any useful learned information. And FedSpace achieves the same result when using VGG16 model. This is because they are not designed for transmission delay-constrained scenarios, but only transmit current update information as possible, thus the update information is likely to be interrupted. Compared to FedSpace with DenseNet121 on fMoW and L-SUN, SatelliteFL obtains 10.7% and 1.8% higher convergence accuracy, and it takes 45.8% and 47.5% fewer days to converge. But for ResNet18 on these two datasets, FedSpace only achieves 13.3% and 34.7%, which is 33.8% and 23.4% lower than SatelliteFL, respectively. Vanilla SatelliteFL with uniform quantization method achieves much better performance than the other two baselines, yet, compared with our SatelliteFL with ProBQ, its time-to-accuracy performance still has obvious shortcomings. For example, using ResNet18 on fMoW and L-SUN, ours accuracy is 3.4% and 12.1% higher than Vanilla, respectively. Even on the other two models, it can achieve up to 1.7%-7.4% accuracy improvement. The reason is that ProBQ can retain more effective gradient information when the ground-to-satellite bandwidth fluctuates and even decreases. Such improvements are mostly attributed to the progressive quantization design as described in Section III-B2.

Authorized licensed use limited to: BEIJING UNIVERSITY OF POST AND TELECOM. Downloaded on July 07,2025 at 10:19:01 UTC from IEEE Xplore. Restrictions apply.



Fig. 8. Time-to-accuracy performance of our SatelliteFL with ProBQ (ours) compared with three baselines: FedAvg, Vanilla SatelliteFL (vanilla) and FedSpace.



Fig. 9. Utilization ratio (u\_ratio) of ground-to-satellite link bandwidth when comparing our SatelliteFL with ProBQ (ours) and three baseline algorithms using three models on two datasets. Considering the on-board training time and data transmission time across satellite-to-ground link, the maximum u\_ratio is only about 50% on average (details in Section IV-A).

Apart from the time-to-accuracy performance, we also evaluate the utilization ratio of ground-to-satellite link bandwidth in Fig. 9. The results show that the u ratios of FedAvg and FedSpace are very low, almost close to 0, due to the frequent failure of ensuring the per-round FL training with one connection window. Among these two algorithms, FedSpace respectively achieves 30.2% and 31.4% with DenseNet121 on fMoW and L-SUN, but they are still 17.2% and 16% lower than our SatelliteFL. Our SatelliteFL with ProBQ can also improve the u\_ratio (3.9%-8.3%) of Vanilla SatelliteFL using naive uniform quantization method. For example, using ResNet18 on fMoW, our SatelliteFL with ProBQ achieves 50.5% u\_ratio, which is 8.3% higher than Vanilla SatelliteFL. The reason behind is that ProBQ designed in our SatelliteFL has progressively expanded the quantization bitwidth on the most of blocks, which can improve the system efficiency under delay constraints.

2) Link Analysis: We further study the performance of SatelliteFL under different ground-to-satellite bandwidth compared to baselines. We keep the hyper-parameter settings the same as the above end-to-end settings. Fig. 10 shows the accuracy and u\_ratio of different methods under varying ground-to-satellite bandwidth. We set  $b_{g,s}$ = 160, 240, 320, 400, 480, and 560, and train DenseNet121 on fMoW in this experiment.

As observed, ProBQ only reduces accuracy performance by  $9.1\%{-}15.6\%$  under poor bandwidth conditions, whereas the



Fig. 10. Performance of our SatelliteFL with ProBQ (ours) compared with three baselines under varying bandwidth.

other three methods decrease accuracy by 19.8%-30.3%, 6.8%-44.1%, and 14.9%-36.5%, respectively. The rationale is that, compared to FedAvg and FedSpace, ProBO strives to preserve effective update information as much as possible under limited bandwidth conditions. In contrast to Vanilla SatelliteFL's fixed quantization, ProBQ's progressive block quantization allows sensitive blocks with lower bitwidth to retain more useful update information. More specifically, with lower bandwidth conditions such as  $b_{q,s} = 160$  and 240, ProBQ reduces only 18.6% and 1.2% u\_ratio, whereas FedAvg and FedSpace struggle to keep the system operational with an effective u ratio. This is because ProBQ permits the retention of more fragmented update information and transmits it at a block-wise granularity. In addition, with  $b_{a,s} >$ 240, ProBQ can maintain a stable u\_ratio close to 50% with consistent accuracy performance. Vanilla SatelliteFL also can keep a stable u\_ratio close to 40%, while the other two baselines experience a sharp decline in u\_ratio as bandwidth becomes worse. The rationale is that, within a relatively ample bandwidth range, each block can retain sufficient useful information for updates. In contrast, in more constrained bandwidth conditions, ensuring the availability of enough update information for the entire model becomes challenging.

3) Sensitivity Analysis: We then dive into the performance sensitivity of SatelliteFL under varying environments. Specifically, to further explore SatelliteFL's performance under poor bandwidth conditions, we set  $b_{g,s}=240$  according to the observation drawn from the above link analysis discussed, while keeping



Fig. 11. Impacts of SatelliteFL settings. V (Vanilla). D (Default): order of original front-to-back. S (Size): order of parameter size. R (Random): order of random shuffling. O (Ours).

the remaining hyper-parameters consistent with the end-to-end configuration.

Impacts of Block Ordering Methods. Fig. 11(a) reports the impact of block ordering methods in our SatelliteFL on the performance of the above two metrics. We introduce another three kinds of ordering methods into Vanilla SatelliteFL: (i) The default method (Default) orders the blocks based on the original front-to-back order; (ii) The size method (Size) orders the blocks based on the order of parameter size; (iii) The random method (Random) orders the blocks based on a random shuffle of all the blocks. We set  $b_{g,s}$ =240 and train DenseNet121 on fMoW in this experiment.

The results in Fig. 11(a) show the compared performance to ProBQ ("Ours"), Vanilla, Default, Size, and Random algorithms. For the accuracy performance, ProBQ still outperforms other algorithms 8.4%–14.6%, this is because ProBQ can achieve a near-optimal choice on block priority under poor bandwidth, to give priority guarantee on the sensitive blocks remaining information with a higher bitwidth. Note that, whatever the block priority we choose in SatelliteFL, the last four methods with block-wise quantization can obtain 6.7% (on average) higher accuracy and 19.1% (on average) higher u\_ratio. The reason for such performance improvement is that our block-wise strategy can make more efficient use of redundant fragmented ground-to-satellite links under poor bandwidth.

Impacts of Bandwidth Variation. Fig. 11(b) shows the accuracy and u\_ratio on different settings of bandwidth variation. We set  $b_{g,s} = 240$  and change the standard deviation = 40, 80, 120, 160, and 200 in Gaussian distribution to simulate the dynamic fluctuations of ground-to-satellite bandwidth, and then train DenseNet121 on fMoW in this experiment.

As observed, Fig. 11(b) illustrates the performance of SatelliteFL under varying degrees of bandwidth fluctuations. With an increase in the degree of bandwidth variation, the u\_ratio decreased by 1.1%–2.4%, while the accuracy remained nearly constant, fluctuating within approximately 1%. The underlying reason for this phenomenon is that ProBQ can effectively adjust the different bitwidth for parameter quantization under varying degrees of dynamic bandwidth changing, thereby maintaining model accuracy as stably as possible.

Impacts of Connection Window. Fig. 11(c) shows the accuracy and u\_ratio on different settings of connection window. We set  $b_{g,s} = 240$  and train DenseNet121 on fMoW with connection window = 4, 5, 6, 7, 8, and 9 minutes in this experiment.

The results in Fig. 11(c) show the benefit of increasing the connection window time. It shows that SatelliteFL can improve accuracy by 39.5%–59.6% and  $2.1\times$ – $3.7\times$  u\_ratio, which is due to ProBQ being able to retain more update information and benefit from a richer connection window time. Note that, when the connection window time > 6, the u\_ratio exceeds 50\%, this is because ProBQ has more transmission time to transmit gradients across the ground-to-satellite link and maintain a stable accuracy performance.

Impacts of Connected Satellite Numbers. Fig. 11(d) shows the accuracy and u\_ratio on different settings of connected satellites numbers. We set  $b_{g,s} = 240$  and train DenseNet121 on fMoW with connected satellites number = 3, 4, 5, 6, and 7 in this experiment.

Fig. 11(d) illustrates the performance on SatelliteFL for different numbers of connected satellites. As the number of connected satellites increases, the u\_ratio abnormally decreases by 0.6%–8.9%, the accuracy first increases by 5.1%, and then also abnormally decreases by 0.9%–1.9%. The reason behind the abnormal phenomenon is that when connected satellite numbers increase, there is a greater probability of selecting satellites with worse link conditions, and these satellites have a poor u\_ratio, which also affects the final accuracy.

#### V. DISCUSSION AND LIMITATION

In our research, we build upon the foundational work presented in prior work [2], which assumes the existence of a steady (with some dynamic changes) communication link during each connection window. However, it is important to acknowledge that in real-world satellite-ground communication scenarios, achieving such steadiness can be challenging due to varying distance between the satellite and the ground station or atmospheric conditions like rainfall-induced signal attenuation. The presence of an uncertain communication link can pose difficulties in meeting our objective of completing the per-round FL training within each connection window.

In our future work endeavors, we aim to enhance the success rate of parameter transmission by developing predictive models for link uncertainty. Additionally, we intend to explore redundant training strategies and missing-parameter recovery mechanisms, minimizing the need for re-transmissions following parameter transmission failures based on the concept of bounded-loss tolerance commonly applied in ML tasks [63]. Furthermore, the images collected by Earth observation satellites lack labels, posing a challenge in efficiently utilizing this raw data. We intend to explore effective methods for providing soft labels to these raw data, focusing on unsupervised training [64] and leveraging the few-shot capabilities exhibited by current large models [65].

# VI. CONCLUSION

SatelliteFL is a communication efficient FL framework across satellites and ground stations that focuses on three unique challenges introduced by the dynamic orbits of LEO satellites: limited ground-to-satellite bandwidth, short connection window and long connection cycle. It formulates the goal as a delayconstrained optimization problem to maximize the model utility while guaranteeing not to exceed the connection window. To solve this problem, it leverages a progressive weight quantization method that fine-grained quantizes each block of model with a unique quantization bitwidth. Experiments show that SatelliteFL can accelerate the convergence by  $1.8 \times (\text{on average}, up to 2.8 \times)$  and improve the bandwidth utilization ratio by  $4.5 \times (\text{on average}, up to 9.3 \times)$  with acceptable accuracy loss.

#### REFERENCES

- Z. Lai, Q. Wu, H. Li, M. Lv, and J. Wu, "OrbitCast: Exploiting megaconstellations for low-latency earth observation," in *Proc. Int. Conf. Netw. Protoc.*, 2021, pp. 1–12.
- [2] J. So, K. Hsieh, B. Arzani, S. Noghabi, S. Avestimehr, and R. Chandra, "FedSpace: An efficient federated learning framework at satellites and ground stations," 2022, arXiv:2202.01267.
- [3] S. Yu, X. Gong, Q. Shi, X. Wang, and X. Chen, "EC-SAGINs: Edgecomputing-enhanced space-air-ground-integrated networks for Internet of Vehicles," *IEEE Internet Things J.*, vol. 9, no. 8, pp. 5742–5754, Apr. 2022.
- [4] Remote sensing from space: What norms govern?. 2023. [Online]. Available: https://www.planet.com/our-constellations/
- Black summer. 2021. [Online]. Available: https://www.abc.net.au/news/ science/2021--03-14/bushfires-detecting-them-from-space-fireballsatellite-launch/13203470
- [6] Remote sensing from space: What norms govern?. 2023. [Online]. Available: https://www.justsecurity.org/86114/remote-sensing-fromspace-what-norms-govern/
- [7] B. Denby and B. Lucia, "Orbital edge computing: Nanosatellite constellations as a new class of computer system," in *Proc. Int. Conf. Architectural Support Program. Lang. Operating Syst.*, 2020, pp. 939–954.
- [8] A. Chen, Y. Xie, Y. Wang, and L. Li, "Knowledge graph-based image recognition transfer learning method for on-orbit service manipulation," *Space: Sci. Technol.*, vol. 2021, 2021, Art. no. 9807452.
- [9] M. Handley, "Using ground relays for low-latency wide-area routing in megaconstellations," in *Proc. Workshop Hot Topics Netw.*, 2019, pp. 125–132.
- [10] N. Razmi, B. Matthiesen, A. Dekorsy, and P. Popovski, "Ground-assisted federated learning in LEO satellite constellations," *IEEE Wireless Commun. Lett.*, vol. 11, no. 4, pp. 717–721, Apr. 2022.
- [11] N. Razmi, B. Matthiesen, A. Dekorsy, and P. Popovski, "On-board federated learning for dense LEO constellations," in *Proc. Int. Conf. Commun.*, 2022, pp. 4715–4720.
- [12] F. Tang, C. Wen, X. Chen, and N. Kato, "Federated learning for intelligent transmission with space-air-ground integrated network sagin toward 6G," *IEEE Netw.*, vol. 37, no. 2, pp. 198–204, Mar./Apr. 2023.
- [13] J. Nguyen et al., "Federated learning with buffered asynchronous aggregation," in Proc. Int. Conf. Artif. Intell. Statist., 2022, pp. 3581–3607.
- [14] D. Vasisht and R. Chandra, "L2D2: Low latency distributed downlink for LEO satellites," in *Proc. Special Int. Group Data Commun.*, 2021, pp. 151–164.
- [15] D. Bhattacherjee, S. Kassing, M. Licciardello, and A. Singla, "In-orbit computing: An outlandish thought experiment?," in *Proc. Workshop Hot Topics Netw.*, 2020, pp. 197–204.

- [16] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas, "Communication-efficient learning of deep networks from decentralized data," in *Proc. Int. Conf. Artif. Intell. Statist.*, 2017, pp. 1273–1282.
- [17] C. Foster et al., "Constellation phasing with differential drag on planet labs satellites," J. Spacecraft Rockets, vol. 55, pp. 473–483, 2018.
- [18] V. Ignatenko, P. Laurila, A. Radius, O. Antropov, and D. Muff, "Iceye microsatellite SAR constellation status update: Evaluation of first commercial imaging modes," in *Proc. Int. Geosci. Remote Sens. Symp.*, 2020, pp. 3581–3584.
- [19] M. Safyan, "Planet's dove satellite constellation," in *Handbook of Small Satellites: Technology, Design, Manufacture, Applications, Economics and Regulation*, Berlin, Germany: Springer, 2020, pp. 1–17.
- [20] I. Del Portillo, B. G. Cameron, and E. F. Crawley, "A technical comparison of three low earth orbit satellite constellation systems to provide global broadband," *Acta Astronautica*, vol. 159, pp. 123–135, 2019.
- [21] N. U. Hassan, C. Huang, C. Yuen, A. Ahmad, and Y. Zhang, "Dense small satellite networks for modern terrestrial communication systems: Benefits, infrastructure, and technologies," *IEEE Wireless Commun.*, vol. 27, no. 5, pp. 96–103, Oct. 2020.
- [22] Y. Su, Y. Liu, Y. Zhou, J. Yuan, H. Cao, and J. Shi, "Broadband LEO satellite communications: Architectures and key technologies," *IEEE Wireless Commun.*, vol. 26, no. 2, pp. 55–61, Apr. 2019.
- [23] K. Devaraj et al., "Planet high speed radio: Crossing gbps from a 3U cubesat," 2019.
- [24] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. Conf. Comput. Vis. Pattern Recognit.*, 2016, pp. 770–778.
- [25] S. Wang, Q. Li, M. Xu, X. Ma, A. Zhou, and Q. Sun, "Tiansuan constellation: An open research platform," in *Proc. Int. Conf. Edge Comput.*, 2021, pp. 94–101.
- [26] B. Denby and B. Lucia, "Orbital edge computing: Machine inference in space," *IEEE Comput. Archit. Lett.*, vol. 18, no. 1, pp. 59–62, Jan.–Jun. 2019.
- [27] A. Van Etten, "You only look twice: Rapid multi-scale object detection in satellite imagery," 2018, *arXiv:1805.09512*.
- [28] Y. Li, A. Padmanabhan, P. Zhao, Y. Wang, G. H. Xu, and R. Netravali, "Reducto: On-camera filtering for resource-efficient real-time video analytics," in *Proc. Special Int. Group Data Commun.*, 2020, pp. 359–376.
- [29] Z. Xiong, F. Zhang, Y. Wang, Y. Shi, and X. X. Zhu, "EarthNets: Empowering AI in Earth observation," 2022, arXiv:2210.04936.
- [30] H. Guo et al., "SpaceDML: Enabling distributed machine learning in space information networks," *IEEE Netw.*, vol. 35, no. 4, pp. 82–87, Jul./Aug. 2021.
- [31] H. Li, C. Chen, C. Li, L. Liu, and G. Gui, "Aerial computing offloading by distributed deep learning in collaborative satellite-terrestrial networks," in *Proc. Int. Conf. Wireless Commun. Signal Process.*, 2021, pp. 1–6.
- [32] J. Konečný, H. B. McMahan, F. X. Yu, P. Richtárik, A. T. Suresh, and D. Bacon, "Federated learning: Strategies for improving communication efficiency," 2016, arXiv:1610.05492.
- [33] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas, "Communication-efficient learning of deep networks from decentralized data," in *Proc. Int. Conf. Artif. Intell. Statist.*, 2017, pp. 1273–1282.
- [34] B. Matthiesen, N. Razmi, I. Leyva-Mayorga, A. Dekorsy, and P. Popovski, "Federated learning in satellite constellations," 2022, arXiv:2206.00307.
- [35] H. Chen, M. Xiao, and Z. Pang, "Satellite-based computing networks with federated learning," *IEEE Wireless Commun.*, vol. 29, no. 1, pp. 78–84, Feb. 2022.
- [36] Z. Yan, D. Li, X. Yu, and Z. Zhang, "Latency-efficient wireless federated learning with quantization and scheduling," *IEEE Commun. Lett.*, vol. 26, no. 11, pp. 2621–2625, Nov. 2022.
- [37] P. Prakash et al., "IoT device friendly and communication-efficient federated learning via joint model pruning and quantization," *IEEE Internet Things J.*, vol. 9, no. 15, pp. 13638–13650, Aug. 2022.
- [38] D. Xu et al., "Mandheling: Mixed-precision on-device DNN training with DSP offloading," in *Proc. Int. Conf. Mobile Comput. Netw.*, 2022, pp. 214–227.
- [39] J. Yoon, G. Park, W. Jeong, and S. J. Hwang, "Bitwidth heterogeneous federated learning with progressive weight dequantization," 2022, arXiv:2202.11453.
- [40] Q. Zhou et al., "Octo:INT8 training with loss-aware compensation and backward quantization for tiny on-device learning," in *Proc. Annu. Tech. Conf.*, 2021, pp. 177–191.
- [41] F. Zhu et al., "Towards unified INT8 training for convolutional neural network," in *Proc. Conf. Comput. Vis. Pattern Recognit.*, 2020, pp. 1969–1979.

- [42] R. Banner, I. Hubara, E. Hoffer, and D. Soudry, "Scalable methods for 8-bit training of neural networks," in *Proc. Adv. Neural Inf. Process. Syst.*, 2018, pp. 5151–5159.
- [43] B. Jacob et al., "Quantization and training of neural networks for efficient integer-arithmetic-only inference," in *Proc. Conf. Comput. Vis. Pattern Recognit.*, 2018, pp. 2704–2713.
- [44] Y. Zhou, S.-M. Moosavi-Dezfooli, N.-M. Cheung, and P. Frossard, "Adaptive quantization for deep neural network," in *Proc. Assoc. Advance. Artif. Intell.*, 2018, pp. 4596–4604.
- [45] S. Khoram and J. Li, "Adaptive quantization of neural networks," in *Proc. Int. Conf. Learn. Representations*, 2018.
- [46] A. Farrag, S. Othman, T. Mahmoud, and A. Y. ELRaffiei, "Satellite swarm survey and new conceptual design for earth observation applications," *Egyptian J. Remote Sens. Space Sci.*, vol. 24, pp. 47–54, 2021.
- [47] Earth observation satellites in 2022. 2022. [Online]. Available: https:// www.pixalytics.com/earth-observation-satellites-2022
- [48] D. Giggenbach, J. Horwath, and M. Knapek, "Optical data downlinks from earth observation platforms," in *Proc. Free-Space Laser Commun. Technol. XXI*, 2009, pp. 17–30.
- [49] K. Devaraj et al., "Dove high speed downlink system," 2017.
- [50] Cubesat communications system table. 2019. [Online]. Available: https: //www.klofas.com
- [51] C. Li et al., "Block-wisely supervised neural architecture search with knowledge distillation," in *Proc. Conf. Comput. Vis. Pattern Recognit.*, 2020, pp. 1986–1995.
- [52] T. Ibaraki, "Integer programming formulation of combinatorial optimization problems," *Discrete Math.*, vol. 16, no. 1, pp. 39–52, 1976.
- [53] D. Li et al., Nonlinear Integer Programming, vol. 84, Berlin, Germany: Springer, 2006.
- [54] X. Chen, L. Jiao, W. Li, and X. Fu, "Efficient multi-user computation offloading for mobile-edge cloud computing," *IEEE/ACM Trans. Netw.*, vol. 24, no. 5, pp. 2795–2808, Oct. 2016.
- [55] G. Christie, N. Fendley, J. Wilson, and R. Mukherjee, "Functional map of the world," in *Proc. Conf. Comput. Vis. Pattern Recognit.*, 2018, pp. 6172–6180.
- [56] F. Yu, A. Seff, Y. Zhang, S. Song, T. Funkhouser, and J. Xiao, "LSUN: Construction of a large-scale image dataset using deep learning with humans in the loop," 2015, arXiv:1506.03365.
- [57] G. Huang, Z. Liu, L. van der Maaten, and K. Q. Weinberger, "Densely connected convolutional networks," in *Proc. Conf. Comput. Vis. Pattern Recognit.*, 2017, pp. 2261–2269.
- [58] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," 2014, arXiv:1409.1556.
- [59] C. He et al., "FedML: A research library and benchmark for federated machine learning," 2020, arXiv:2007.13518.
- [60] K. Hsieh, A. Phanishayee, O. Mutlu, and P. Gibbons, "The non-iid data quagmire of decentralized machine learning," in *Proc. Int. Conf. Mach. Learn.*, 2020, pp. 4387–4398.
- [61] Y. Zhao, M. Li, L. Lai, N. Suda, D. Civin, and V. Chandra, "Federated learning with non-iid data," 2018, arXiv:1806.00582.
- [62] T.-M. H. Hsu, H. Qi, and M. Brown, "Federated visual classification with real-world data distribution," in *Proc. Eur. Conf. Comput. Vis.*, 2020, pp. 76–92.
- [63] J. Xia et al., "Rethinking transport layer design for distributed machine learning," in Proc. Asia-Pacific Workshop Netw., 2019, pp. 22–28.
- [64] H. Jang, H. Lee, and J. Shin, "Unsupervised meta-learning via few-shot pseudo-supervised contrastive learning," 2023, arXiv:2303.00996.
- [65] A. Kirillov et al., "Segment anything," 2023, arXiv:2304.02643.



**Jinliang Yuan** is currently working toward the PhD degree with the School of Computer Science, Beijing University of Posts and Telecommunication, Beijing, China. His research interests are distributed machine learning, federated learning, and edge computing.



**Yaozong Wu** is currently working toward the master's degree with the School of Computer Science, Beijing University of Posts and Telecommunication, Beijing, China. His research interests are federated learning and resource-efficient AI systems.



**Qibo Sun** received the PhD degree in communication and electronic system from the Beijing University of Posts and Telecommunication, in 2002. He is an associate professor with the School of Computer Science and Engineering, Beijing University of Posts and Telecommunications, China. His research interests include services computing, satellite computing and space-ground integration network. He has published more than 100 papers. He is a member of the China Computer Federation and Chinese Association for Artificial Intelligence.

Ao Zhou received the PhD degree from the Beijing University of Posts and Telecommunications, Beijing, China, in 2015. She is currently an associate professor with the State Key Laboratory of Networking and Switching Technology, BUPT. She has published more than 50 research papers. She played a key role with many international conferences. Her research interests include cloud computing and mobile edge computing.



Shangguang Wang (Senior Member, IEEE) is a professor with the School of Computer Science, Beijing University of Posts and Telecommunications, China. He is the founder and chief scientist of the Tiansuan Constellation. His research interests include service computing, mobile edge computing, and satellite computing. He is currently serving as chair of IEEE Technical Committee on Services Computing, and vice chair of IEEE Technical Committee on Cloud Computing. He also served as general chairs or program chairs of more than 10 IEEE conferences. He is a fellow of the IET.



**ChenYang** is currently working toward the PhD degree with the School of Computer Science, Beijing University of Posts and Telecommunication, Beijing, China. His research interests are federated learning and edge computing.



**Mengwei Xu** is an assistant professor with the Computer Science Department, Beijing University of Posts and Telecommunications. His research interests cover the broad areas of mobile computing, edge computing, artificial intelligence, and system software.