

RECONCILING SECURITY AND COMMUNICATION EFFICIENCY IN FEDERATED LEARNING

Anonymous authors

Paper under double-blind review

ABSTRACT

Cross-device Federated Learning is an increasingly popular machine learning setting to train a model by leveraging a large population of client devices with high privacy and security guarantees. However, communication efficiency remains a major bottleneck when scaling federated learning to production environments, particularly due to bandwidth constraints during uplink communication. In this paper, we formalize and address the problem of compressing client-to-server model updates under the Secure Aggregation primitive, a core component of Federated Learning pipelines that allows the server to aggregate the client updates without accessing them individually. In particular, we adapt standard scalar quantization and pruning methods to Secure Aggregation and propose Secure Indexing, a variant of Secure Aggregation that supports quantization for extreme compression. We establish state-of-the-art results on LEAF benchmarks in a secure Federated Learning setup with up to $40\times$ compression in uplink communication with no meaningful loss in utility compared to uncompressed baselines.

1 INTRODUCTION

Federated Learning (FL) is a distributed machine learning (ML) paradigm that trains a model across a number of participating entities holding local data samples. In this work, we focus on *cross-device* FL that harnesses a large number (hundreds of millions) of edge devices with disparate characteristics such as availability, compute, memory, or connectivity resources (Kairouz et al., 2019).

Two challenges to the success of cross-device FL are privacy and scalability. FL was originally motivated for improving privacy since data points remain on client devices. However, as with other forms of ML, information about training data can be extracted via membership inference or reconstruction attacks on a trained model (Carlini et al., 2021a;b; Watson et al., 2022), or leaked through local updates (Melis et al., 2019; Geiping et al., 2020). Consequently, Secure Aggregation (SECAGG) protocols were introduced to prevent the server from directly observing individual client updates, which is a major vector for information leakage (Bonawitz et al., 2019; Huba et al., 2022). Additional mitigations such as Differential Privacy (DP) may be required to offer further protection against attacks (Dwork et al., 2006; Abadi et al., 2016), as discussed in Section 6.

Ensuring scalability to populations of heterogeneous clients is the second challenge for FL. Indeed, wall-clock training times are highly correlated with increasing model and batch sizes (Huba et al., 2022), even with recent efforts such as FedBuff (Nguyen et al., 2022), and communication overhead between the server and clients dominates model convergence time. Consequently, compression techniques were used to reduce the communication bandwidth while maintaining model accuracy. However, a fundamental problem has been largely overlooked in the literature: in their native form, standard compression methods such as scalar quantization and pruning are not compatible with SECAGG. This makes it challenging to ensure both security and communication efficiency.

In this paper, we address this gap by adapting compression techniques to make them compatible with SECAGG. We focus on compressing *uplink* updates from clients to the server for three reasons. First, uplink communication is more sensitive and so is subject to a high security bar, whereas downlink updates broadcast by the server are deemed public. Second, upload bandwidth is generally more restricted than download bandwidth. For instance, according to the most recent FCC¹ report, the

¹US Federal Communications Commission.

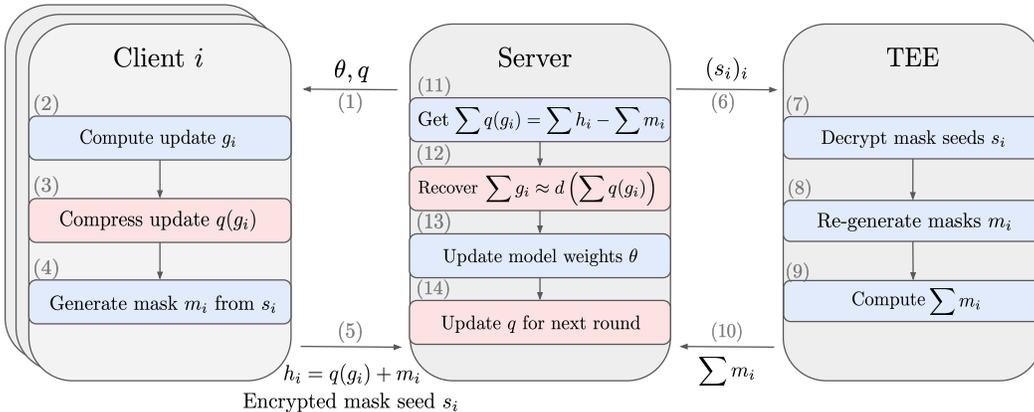


Figure 1: Summary of the proposed approach for one FL round, where we omit the round dependency and Differential Privacy (DP) for clarity. Blue boxes denote standard steps and red boxes denote additional steps for uplink compression. Client i computes local model update g_i , compresses it with the compression operator q , and encrypts it by adding a random mask m_i in the compressed domain, hence reducing the uplink bandwidth (steps 2–4). The server recovers the aggregate in the compressed domain by leveraging any SECAGG protocol (steps 7–13, with a TEE-based SECAGG, see Section 3.1). Since the decompression operator d is linear, the server can convert the aggregate back to the non-compressed domain, up to compression error (step 12). As with the model weights θ , the compression operator q are also periodically updated and broadcast by the server (step 14). In Section 4, we apply the proposed method to scalar quantization and pruning without impacting SECAGG and propose Secure Indexing, a variant of SECAGG for extreme uplink compression with product quantization. See Section 3.1 for details about SECAGG and Section 6 for a discussion on DP.

ratio of download to upload speeds for DSL and cable providers² in the US ranges between $3 \times$ to $20 \times$ (FCC, 2021). Finally, efficient uplink communication brings several benefits beyond speeding up convergence: lowering communication cost reduces selection bias due to under-sampling clients with limited connectivity, improving fairness and inclusiveness. It also shrinks the carbon footprint of FL, the fraction of which attributable to communication can reach 95% (Qiu et al., 2021).

In summary, we present the following contributions in this paper:

- We highlight the fundamental mismatch between two critical components of the FL stack: SECAGG protocols and uplink compression mechanisms.
- We formulate solutions by imposing a linearity constraint on the decompression operator, as illustrated in Figure 1 in the case of TEE-based SECAGG.
- We adapt the popular scalar quantization and (random) pruning compression methods for compatibility with the FL stack that require no changes to the SECAGG protocol.
- For extreme uplink compression without compromising security, we propose Secure Indexing (SECIND), a variant of SECAGG that supports product quantization.

2 RELATED WORK

Communication is identified as a primary efficiency bottleneck in FL, especially in the cross-device FL setting (Kairouz et al., 2019). This has led to significant interest in reducing FL’s communication requirements. In what follows, we might refer to any local model update in a distributed training procedure as a *gradient*, including model updates computed following multiple local training steps.

Efficient Distributed Optimization. There is a large body of literature on reducing the communication cost for distributed training. Seide et al. (2014) proposes quantizing gradients to one bit while carrying the quantization error forward across mini-batches with error feedback. Similarly, Wen et al. (2017) proposes layer-wise ternary gradients and Bernstein et al. (2018) suggests using only the sign

²FL is typically restricted to using unmetered connections, usually over Wi-Fi (Huba et al., 2022).

of the gradients. Gradient sparsity is another related area that is extensively studied (Wangni et al., 2018; Aji & Heafield, 2017; Lin et al., 2018; Renggli et al., 2019; Parcollet et al., 2022). For instance, Chen et al. (2018) and Han et al. (2020) explore adapting the degree of sparsity to the distribution of local client data. Another method, QSGD, tunes the quantization level to trade possibly higher variance gradients for reduced communication bandwidth (Alistarh et al., 2017). Researchers also studied structured and sketched model updates (Konečný et al., 2016). For example, Wang et al. (2018) proposes expressing gradients as a linear combination of basis vectors common to all workers and Wang et al. (2022) propose to cluster the gradients and to implement error correction on the client side. Besides gradient compression, other methods such as Vepakomma et al. (2018); Hu et al. (2019) propose reducing the communication cost by partitioning the model such that each client learns a portion of it, while He et al. (2020) proposes training small models and periodically distilling them to a larger central model. However, as detailed in Section 3 and below, most of the proposed methods are not readily compatible with SECAGG and cannot be used in secure FL.

Bi-directional Compression. In addition to uplink gradient compression, a line of work also focuses on downlink model compression. In a non-distributed setup, Zhou et al. (2016); Courbariaux et al. (2015) demonstrates that it is possible to meaningfully train with low bit-width models and gradients. In FL, Jiang et al. (2019) proposes adapting the model size to the device to reduce both communication and computation overhead. Since the local models are perturbed due to compression, researchers propose adapting the optimization algorithm for better convergence (Liu et al., 2020; Sattler et al., 2020; Tang et al., 2019; Zheng et al., 2019; Amiri et al., 2020; Philippenko & Dieuleveut, 2021). Finally, pre-conditioning models during FL training can allow for quantized on-device inference, as demonstrated for non-distributed training by Gupta et al. (2015); Krishnamoorthi (2018). As stated in Section 1, we do not focus on downlink model compression since uplink bandwidth is the main communication bottleneck and since SECAGG only involves uplink communication.

Aggregation in the Compressed Domain. In the distributed setting, Yu et al. (2018) propose to leverage both gradient compression and parallel aggregation by performing the *ring all-reduce* operation in the compressed domain and decompressing the aggregate. To do so, the authors exploit temporal correlations of the gradients to design a linear compression operator. Another method, PowerSGD (Vogels et al., 2019), leverages a fast low-rank gradient compressor. However, both aforementioned methods are not evaluated in the FL setup and do not mention SECAGG. Indeed, the proposed methods focus on decentralized communication between the workers by leveraging the all-reduce operation. Moreover, Power SGD incorporates (stateful) error feedback on all distributed nodes, which is not readily adaptable to cross-device FL in which clients generally participate in a few (not necessarily consecutive) rounds. Finally, Rothchild et al. (2020) proposes FetchSGD, a compression method relying on a CountSketch, which is compatible with SECAGG.

3 BACKGROUND

In this section, we recall the SECAGG protocol first, then the compression methods that we wish to adapt to SECAGG, namely, scalar quantization, pruning, and product quantization.

3.1 SECURE AGGREGATION

SECAGG refers to a class of protocols that allow the server to aggregate client updates without accessing them individually. While SECAGG alone does not entirely prevent client data leakage, it is a powerful and widely-used component of current at-scale cross-device FL implementations (Kairouz et al., 2019). Two main approaches exist in practice: software-based protocols relying on Multiparty Computation (MPC) (Bonawitz et al., 2019; Bell et al., 2020; Yang et al., 2022), and those that leverage hardware implementations of Trusted Execution Environments (TEEs) (Huba et al., 2022).

SECAGG relies on additive masking, where clients protect their model updates g_i by adding a uniform random mask m_i to it, guaranteeing that each client’s masked update is statistically indistinguishable from any other value. At aggregation time, the protocol ensures that all the masks are canceled out. For instance, in an MPC-based SECAGG, the pairwise masks cancel out within the aggregation itself, since for every pair of users i and j , after they agree on a matched pair of input perturbations, the masks $m_{i,j}$ and $m_{j,i}$ are constructed so that $m_{i,j} = -m_{j,i}$. Similarly and as illustrated in Fig. 1, in a TEE-based SECAGG, the server receives $h_i = g_i + m_i$ from each client as well as the sum of the

masks $\sum_i m_i$ from the TEE and recovers the sum of the updates as

$$\sum_i g_i = \sum_i h_i - \sum_i m_i.$$

We defer the discussion of DP noise addition by SECAGG protocols to Section 6.

Finite Group. SECAGG requires that the plaintexts—client model updates—be elements of a finite group, while the inputs are real-valued vectors represented with floating-point types. This requirement is usually addressed by converting client updates to fixed-point integers and operating in a finite domain (modulo 2^p) where p is typically set in prior literature to 32 bits. The choice of SECAGG bit-width p must balance communication costs with the accuracy loss due to rounding and overflows.

Minimal Complexity. TEE-based protocols offer greater flexibility in how individual client updates can be processed; however, the code executed inside TEE is part of the trusted computing base (TCB) for all clients. In particular, it means that this code must be stable, auditable, defects- and side-channel-free, which severely limits its complexity. Hence, in practice, we prefer compression techniques that are either oblivious to SECAGG’s implementation or require minimal changes to the TCB.

3.2 COMPRESSION METHODS

In this subsection, we consider a matrix $W \in \mathbb{R}^{C_{in} \times C_{out}}$ representing the weights of a linear layer to discuss three major compression methods with distinct compression/accuracy tradeoffs and identify the challenges SECAGG faces to be readily amenable to these popular quantization algorithms.

3.2.1 SCALAR QUANTIZATION

Uniform scalar quantization maps floating-point weight w to 2^b evenly spaced bins, where b is the number of bits. Given a floating-point scale $s > 0$ and an integer shift parameter z called the zero-point, we map any floating-point parameter w to its nearest bin indexed by $\{0, \dots, 2^b - 1\}$:

$$w \mapsto \text{clamp}(\text{round}(w/s) + z, [0, 2^b - 1]).$$

The tuple (s, z) is often referred to as the quantization parameters (`qparams`). With $b = 8$, we recover the popular `int8` quantization scheme (Jacob et al., 2018), while setting $b = 1$ yields the extreme case of binarization (Courbariaux et al., 2015). The quantization parameters s and z are usually calibrated after training a model with floating-point weights using the minimum and maximum values of each layer. The compressed representation of weights W consists of the `qparams` and the integer representation matrix W_q where each entry is stored in b bits. Decompressing any integer entry w_q of W_q back to floating point is performed by applying the (linear) operator $w_q \mapsto s \times (w_q - z)$.

Challenge. The discrete domain of quantized values and the finite group required by SECAGG are not natively compatible because of the overflows that may occur at aggregation time. For instance, consider the extreme case of binary quantization, where each value is replaced by a bit. We can represent these bits in SECAGG with $p = 1$, but the aggregation will inevitably result in overflows.

3.2.2 PRUNING

Pruning is a class of methods that remove parts of a model such as connections or neurons according to some pruning criterion, such as weight magnitude (Le Cun et al. (1989); Hassibi & Stork (1992); see Blalock et al. (2020) for a survey). Konečný et al. (2016) demonstrate client update compression with random sparsity for federated learning. Motivated by previous work and the fact that random masks do not leak information about the data on client devices, we will leverage random pruning of client updates in the remainder of this paper. A standard method to store a sparse matrix is the coordinate list (COO) format³, where only the non-zero entries are stored (in floating point or lower precision), along with their integer coordinates in the matrix. This format is compact, but only for a large enough compression ratio, as we store additional values for each non-zero entry. Decompression is performed by re-instantiating the uncompressed matrix with both sparse and non-sparse entries.

Challenge. Pruning model updates on the client side is an effective compression approach as investigated in previous work. However, the underlying assumption is that clients have different

³See the `torch.sparse` documentation.

masks, either due to their seeds or dependency on client update parameters (*e.g.* weight magnitudes). This is a challenge for SECAGG as aggregation assumes a dense compressed tensor, which is not possible to construct when the coordinates of non-zero entries are not the same for all clients.

3.2.3 PRODUCT QUANTIZATION

Product quantization (PQ) is a compression technique developed for nearest-neighbor search (Jégou et al., 2011) that can be applied for model compression (Stock et al., 2020). Here, we show how we can re-formulate PQ to represent model updates. We focus on linear layers and refer the reader to Stock et al. (2020) for adaptation to convolutions. Let the *block size* be d (say, 8), the number of *codewords* be k (say, 256) and assume that the number of input channels, C_{in} , is a multiple of d . To compress W with PQ, we evenly split its columns into subvectors or blocks of size $d \times 1$ and learn a *codebook* via k -means to select the k codewords used to represent the $C_{\text{in}} \times C_{\text{out}}/d$ blocks of W . PQ with block size $d = 1$ amounts to non-uniform scalar quantization with $\log_2 k$ bits per weight.

The PQ-compressed matrix W is represented with the tuple (C, A) , where C is the codebook of size $k \times d$ and A gives the assignments of size $C_{\text{in}} \times C_{\text{out}}/d$. Assignments are integers in $[0, k - 1]$ and denote which codebook a subvector was assigned to. To decompress the matrix (up to reshaping), we index the codebook with the assignments, written in PyTorch-like notation as

$$\widehat{W} = C[A].$$

Challenge. There are several obstacles to making PQ compatible with SECAGG. First, each client may have a different codebook, and direct access to these codebooks is needed to decode each client’s message. Even if all clients share a (public) codebook, the operation to take assignments to produce an (aggregated) update is not linear, and so cannot be directly wrapped inside SECAGG.

4 METHOD

In this section, we propose solutions to reconcile security (SECAGG) and communication efficiency. Our approach is to modify compression techniques to share some hyperparameters globally across all clients so that aggregation can be done by uniformly combining each client’s response, while still ensuring that there is scope to achieve accurate compressed representations. As detailed below, each of the proposed methods offers the same level of security as standard SECAGG without compression.

4.1 SECURE AGGREGATION AND COMPRESSION

We propose to compress the uplink model updates through a compression operator q , whose parameters are round-dependent but the same for all clients participating in the same round. Then, we will add a random mask m_i to each quantized client update $q(g_i)$ in the compressed domain, thus effectively reducing uplink bandwidth while ensuring that $h_i = q(g_i) + m_i$ is statistically indistinguishable from any other representable value in the finite group (see Section 3.1). In this setting, SECAGG allows the server to recover the aggregate of the client model updates in the compressed domain: $\sum_i q(g_i)$. If the decompression operator d is linear, the server is able to recover the aggregate in the non-compressed domain, up to quantization error, as illustrated in Figure 1:

$$d\left(\sum_i h_i - \sum_i m_i\right) = d\left(\sum_i q(g_i)\right) = \sum_i d(q(g_i)) \approx \sum_i g_i.$$

The server periodically updates the quantization and decompression operator parameters, either from the aggregated model update, which is deemed public, or by emulating a client update on some similarly distributed public data. Once these parameters are updated, the server broadcasts them to the clients for the next round. This adds overhead to the downlink communication payload, however, this is negligible compared to the downlink model size to transmit. For instance, for scalar quantization, q is entirely characterized by one `fp32` scale and one `int32` zero-point per layer, the latter of which is unnecessary in the case of a symmetric quantization scheme. Finally, this approach is compatible with both synchronous FL methods such as FedAvg (McMahan et al., 2017) and asynchronous methods such as FedBuff (Nguyen et al., 2022) as long as SECAGG maintains the mapping between the successive versions of quantization parameters and the corresponding client updates.

4.2 APPLICATION

Next, we show how we adapt scalar quantization and random pruning with no changes required to SECAGG. We illustrate our point with TEE-based SECAGG while these adapted uplink compression mechanisms are agnostic of the SECAGG mechanism. Finally, we show how to obtain extreme uplink compression by proposing a variant of SECAGG, which we call SECIND. This variant supports product quantization and is provably secure.

4.2.1 SCALAR QUANTIZATION AND SECURE AGGREGATION

As detailed in Section 3.2.1, a model update matrix g_i compressed with scalar quantization is given by an integer representation in the range $[0, 2^b - 1]$ and by the quantization parameters *scale* (s) and *zero-point* (z). A sufficient condition for the decompression operator to be linear is to broadcast common quantization parameters per layer for each client. Denote $q(g_i)$ as the integer representation of quantized client model update g_i corresponding to a particular layer for client $1 \leq i \leq N$. Set the scale of the decompression operator to s and its zero-point to z/N . Then, the server is able to decompress as follows (where the decompression operator is defined in Section 3.2.1):

$$d\left(\sum_i q(g_i)\right) = s \sum_i q(g_i) - \frac{z}{N} = \sum_i (s(q(g_i)) - z) \approx \sum_i g_i$$

Recall that all operations are performed in a finite group. Therefore, to avoid overflows at aggregation time, we quantize with a bit-width b but take SECAGG bit-width $p > b$, thus creating a margin for potential overflows (see Section 5.3). This approach is related to the fixed-point aggregation described in (Bonawitz et al., 2019; Huba et al., 2022), but we calibrate the quantization parameters and perform the calibration per layer and periodically, unlike the related approaches.

Privacy, Security and Bandwidth. Scales and zero points are determined from public data on the server. Downlink overhead is negligible: the server broadcasts the per-layer quantization parameters. The upload bandwidth is p bits per weight, where p is the SECAGG finite group size (Section 3.1). Since the masks m_i are chosen in the integer range $[0, 2^p - 1]$, any masked integer representation taken modulo 2^p is statistically indistinguishable from any other vector.

4.2.2 PRUNING AND SECURE AGGREGATION

To enable linear decompression with random pruning, all clients will share a common pruning mask for each round. This can be communicated compactly before each round as a seed for a pseudo-random function. This pruning mask seed is different from the SECAGG mask seed introduced in Section 3.1 and has a distinct role. Each client uses the pruning seed to reconstruct a pruning mask, prunes their model update g_i , and only needs to encrypt and transmit the unpruned parameters. The trade-off here is that some parameters are completely unobserved in a given round, as opposed to traditional pruning. SECAGG operates as usual and the server receives the sum of the tensor of unpruned parameters computed by participating clients in the round, which it can expand using the mask seed. We denote the pruning operator as ϕ applied to the original model update g_i , and the decompression operator as d applied to a compressed tensor $\phi(g_i)$. Decompression is an expansion operation equivalent to multiplication with a sparse permutation matrix P_i whose entries are dependent on the i 'th client's mask seed. Crucially, when all clients share the same mask seed within each round, we have $P_i = P$ for all i and linearity of decompression is maintained:

$$d\left(\sum_i \phi(g_i)\right) = P\left(\sum_i \phi(g_i)\right) = \sum_i P_i \phi(g_i) = \sum_i d(\phi(g_i)) \approx \sum_i g_i.$$

Privacy, Security and Bandwidth. Since the mask is random, no information leaks from the pruning mask. The downlink overhead (the server broadcasts one integer mask seed) is negligible. The upload bandwidth is simply the size of the sparse client model updates. Finally, there is no loss in security since each client uses standard SECAGG mechanism on the non-pruned entries.

4.2.3 PRODUCT QUANTIZATION AND SECURE INDEXING

We next describe the Secure Indexing (SECIND) primitive, and discuss how to instantiate it. Recall that with PQ, each layer has its own codebook C as explained in Section 4. Let us fix one particular

Algorithm 1 Secure Indexing (SECIND)

```

1: procedure SECUREINDEXING( $C$ )                                ▷ This happens inside the TEE
2:   Receive common codebook  $C$  from server                    ▷  $C$  is periodically updated by the server
3:   Initialize histograms  $H_{m,n}$  to 0                        ▷ Each histogram for block  $(m, n)$  has size  $k$ 
4:   for each client  $i$  do
5:     Receive and decrypt assignment matrix  $A^i$ 
6:     for each block index  $(m, n)$  do
7:        $r \leftarrow A_{m,n}^i$                                 ▷ Recover assignment of client  $i$  for block  $(m, n)$ 
8:        $H_{m,n}[r] \leftarrow H_{m,n}[r] + 1$                 ▷ Update global count for codeword index  $r$ 
9:     Send back histograms  $H_{m,n}$  to the server

```

layer compressed with codebook C , containing k codewords. We assume that C is common to all clients participating in the round. Consider the assignment matrix of a given layer $(A^i)_{m,n}$ for client i . From these, we seek to build the *assignment histograms* $H_{m,n} \in \mathbb{R}^k$ that satisfy

$$H_{m,n}[r] = \sum_i \mathbf{1}(A_{m,n}^i = r),$$

where the indicator function $\mathbf{1}(A_{m,n}^i = r) = 1$ if $A_{m,n}^i = r$ and 0 otherwise. A *Secure Indexing* primitive will produce $H_{m,n}$ while ensuring that no other information about client assignments or partial aggregations is revealed. The server receives assignment histograms from SECIND and is able to recover the aggregated update for each block indexed by (m, n) as

$$\sum_r H_{m,n}[r] \cdot C[r].$$

We describe how SECIND can be implemented with a TEE in Algorithm 1. Each client encrypts the assignment matrix, for instance with additive masking as described in Section 3.1, and sends it to the TEE via the server. Hence, the server does not have access to the plaintexts client-specific assignments. TEE decrypts each assignment matrix and for each block indexed by (m, n) produces the assignment histogram, which can then be mapped to an update via the (public) codebook. Compared to SECAGG, where the TEE receives an encrypted seed per client (a few bytes per client) and sends back the sum of the masks m_i (same size as the considered model), SECIND receives the (masked) assignment matrices and sends back the aggregated update for the round. SECIND implementation feasibility is briefly discussed in Appendix A.3.

Privacy, Security and Bandwidth. Codebooks are computed from public data while individual assignments are never revealed to the server. The downlink overhead of sending the codebooks is negligible as demonstrated in Section 5. The upload bandwidth in the TEE implementation is the assignment size, represented in k bits (the number of codewords). For instance, with a block size $d = 8$ and $k = 32$ codewords, assignment storage costs are 5 bits per 8 weights, which converts to 0.625 bits per weight. The tradeoff compared to non-secure PQ is the restriction to a global codebook for all clients (instead of one tailored to each client), and the need to instantiate SECIND instead of SECAGG. Since the assignments are encrypted before being sent to the TEE, there is no loss in security. Here, any encryption mechanism (not necessarily relying on additive masking) would work.

5 EXPERIMENTS

In this section, we numerically evaluate the performance of the proposed approaches when adapted to SECAGG protocols. We study the relationship between uplink compression and model accuracy for the LEAF benchmark tasks. In addition, for scalar and product quantization we also analyze the impact of refresh rate for compression parameters on overall model performance.

5.1 EXPERIMENTAL SETUP

We closely follow the setup of Nguyen et al. (2022) and use the FLSim library for our experiments. All experiments are run on a single V100 GPU 16 GB (except for Sent140 where we use one V100 32 GB) and typically take a few hours to run. More experiment details can be found in Appendix A.1.



Figure 2: We adapt scalar quantization (SQ) and pruning to the SECAGG protocol to enable efficient and secure uplink communications. We also present results for product quantization (PQ) under the proposed novel SECIND protocol. *The x axis is log-scale* and represents the uplink message size. Baseline refers to SECAGG FL run without any uplink compression, displayed as a horizontal line for easier comparison. Model size is indicated in the plot titles. Uncompressed client updates are as large as the models when $p = 32$ (see Section 3.1, represented as stars). We refer to the Appendix A.2.1 for the matching tables where we additionally report the standard deviation of each data point.

Tasks. We run experiments on three datasets from the LEAF benchmark (Caldas et al., 2018): CelebA (Liu et al., 2015), Sent140 (Go et al., 2009) and FEMNIST (LeCun & Cortes, 2010). For CelebA, we train the same convolutional classifier as Nguyen et al. (2022) with BatchNorm layers replaced by GroupNorm layers and 9,343 clients. For Sent140, we train an LSTM classifier for binary sentiment analysis with 59,400 clients. Finally, for FEMNIST, we train a GroupNorm version of the ResNet18 (He et al., 2016) for digit classification with 3,550 clients. For all compression methods, we do not compress biases and norm layers due to their small overhead.

Baselines. We focus here on the (synchronous) FedAvg approach although, as explained in Section 4, the proposed compression methods can be readily adapted to asynchronous FL aggregation protocols. As done in the literature, we keep the number of clients per round to at most 100 (see Appendix A.1), a small fraction of the total considered population size (Chen et al., 2019; Charles et al., 2021). We report the average and standard deviation of accuracy over three independent runs for all tasks at different uplink byte sizes corresponding to various configurations of the compression operator.

Implementation Details. We refer the reader to Appendix A.1. The downlink overhead of sending the per-layer codebooks for product quantization is negligible as shown in Appendix A.2.4. Finally, the convergence time in terms of rounds is similar for PQ runs and the non-compressed baseline, as illustrated in Appendix A.2.5. Note that outside a simulated environment, the wall-clock time convergence for PQ runs would be *lower* due to uplink communication savings.

5.2 RESULTS AND COMPARISON WITH PRIOR WORK

Results for efficient and secure uplink communications are displayed in Figure 2. We observe that PQ yields a consistently better trade-off curve between model update size and accuracy. For instance, on CelebA, PQ achieves $\times 30$ compression with respect to the non-compressed baseline at iso-accuracy. The iso-accuracy compression rate is $\times 32$ on Sent140 and $\times 40$ on FEMNIST (see Appendix A.2.1 for detailed tables). Scalar quantization accuracy degrades significantly for larger compression rates due to the overflows at aggregation as detailed in Appendix A.2.2.

The line of work that develops FL compression techniques mainly includes FetchSGD (Rothchild et al., 2020) as detailed in Section 2, although the authors do not mention SECAGG. Their results are not directly comparable to ours due to non-matching experimental setups (e.g., datasets and architectures). However, Figure 6 in the appendix of Rothchild et al. (2020) mentions upload compression rates at iso-accuracy that are weaker than those obtained here with product quantization.

5.3 ABLATION STUDIES

We investigate the influence of the frequency of updates of the compression operator q for scalar quantization and pruning, and study the influence of the SECAGG bit-width p on the number of overflows for scalar quantization.

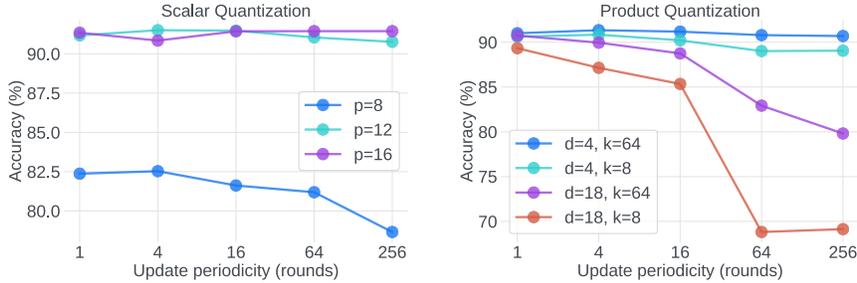


Figure 3: Impact of the refresh rate of the compression operator by the server on the CelebA dataset. **Left:** for scalar quantization (quantization parameters), where we fix the quantization bit-width $b = 8$ (p denotes the SECAGG bit-width). **Right:** for product quantization (codebooks), where k denotes the number of codewords and d the block size.

Update frequency of the compression operators. In Figure 3, we show that for scalar quantization, the update periodicity only plays a role with low SECAGG bit-width values p compared to the quantization bit-width b . For product quantization, the update periodicity plays an important role for aggressive compression setups corresponding to large block sizes d or to a smaller number of codewords k . For pruning, we measure the impact of masks that are refreshed periodically. We observe that if we refresh the compression operator more frequently, staleness is reduced, leading to accuracy improvements. We present our findings in Appendix A.2.6.

Overflows for scalar quantization. As discussed in Section 4.2.1, we choose the SECAGG bit-width p to be greater than quantization bit-width b in order to avoid aggregation overflows. While it suffices to set p to be $\lceil \log_2 n_c \rceil$ more than b , where n_c is the number of clients participating in the round, reducing p is desirable to reduce uplink size. We study the impact of p on the percentage of parameters that suffer overflows and present our findings in Appendix A.2.2.

6 LIMITATIONS AND FUTURE WORK

Compatibility with DP. As mentioned in Section 1, we may want both SECAGG and Differential Privacy (Abadi et al., 2016) to realize the full promise of FL as a privacy-enhancing technology. While our primary focus is on enabling efficient and secure uplink communication, we emphasize that the proposed approaches are compatible with user-level DP. For instance, at the cost of increasing the complexity of the trusted computing base, DP noise can be added natively by the TEE with our modified random pruning or scalar quantization approaches. For PQ and SECIND, we can have the TEE to add noise in the assignment space (outputting a noisy histogram), or to map the histogram to the codeword space and add noise there. Each option offers a different tradeoff between privacy, trust, and accuracy; we leave detailed evaluation to future work.

Efficiency and Privacy. A separate line of work aims to combine communication efficiency and privacy. For instance, Triastcyn et al. (2021) develop a method that unifies compressed communication and DP (where integration with SECAGG is left as an open problem), while Chaudhuri et al. (2022) design a privacy-aware scalar compression mechanism within the *local* differential privacy model.

7 CONCLUSION

In this paper, we reconcile efficiency and security for uplink communication in Federated Learning. We propose to adapt existing compression mechanisms such as scalar quantization and pruning to the secure aggregation protocol by imposing a linearity constraint on the decompression operator. Our experiments demonstrate that we can adapt both quantization and pruning mechanisms to obtain a high degree of uplink compression with minimal degradation in performance and higher security guarantees. For achieving the highest rates of compression, we introduce SECIND, a variant of SECAGG well-suited for TEE-based implementation that supports product quantization while maintaining a high security bar. We plan to extend our work to other federated learning scenarios, such as asynchronous FL, and further investigate the interaction of compression and privacy.

REFERENCES

- Martín Abadi, Andy Chu, Ian Goodfellow, H. Brendan McMahan, Ilya Mironov, Kunal Talwar, and Li Zhang. Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security (CCS)*, pp. 308–318, 2016.
- Alham Fikri Aji and Kenneth Heafield. Sparse communication for distributed gradient descent. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 440–445, 2017.
- Dan Alistarh, Demjan Grubic, Jerry Li, Ryota Tomioka, and Milan Vojnovic. QSGD: Communication-efficient SGD via gradient quantization and encoding. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 30, 2017.
- Mohammad Mohammadi Amiri, Deniz Gunduz, Sanjeev R. Kulkarni, and H. Vincent Poor. Federated learning with quantized global model updates, 2020.
- James Henry Bell, Kallista A Bonawitz, Adrià Gascón, Tancrede Lepoint, and Mariana Raykova. Secure single-server aggregation with (poly) logarithmic overhead. In *Proceedings of the 2020 ACM SIGSAC Conference on Computer and Communications Security (CCS)*, pp. 1253–1269, 2020.
- Jeremy Bernstein, Yu-Xiang Wang, Kamyar Azizzadenesheli, and Animashree Anandkumar. signSGD: Compressed optimisation for non-convex problems. In *Proceedings of the 35th International Conference on Machine Learning (ICML)*, volume 80 of *PMLR*, pp. 560–569, 2018.
- Davis Blalock, Jose Javier Gonzalez Ortiz, Jonathan Frankle, and John Guttag. What is the state of neural network pruning? In *Proceedings of Machine Learning and Systems (MLSys)*, volume 2, pp. 129–146, 2020.
- Kallista A. Bonawitz, Fariborz Salehi, Jakub Konečný, Brendan McMahan, and Marco Gruteser. Federated learning with autotuned communication-efficient secure aggregation. In *53rd Asilomar Conference on Signals, Systems, and Computers (ACSCC)*, pp. 1222–1226. IEEE, 2019.
- Elette Boyle, Niv Gilboa, and Yuval Ishai. Function secret sharing: Improvements and extensions. In *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security (CCS)*, pp. 1292–1303, 2016.
- Sebastian Caldas, Peter Wu, Tian Li, Jakub Konečný, H. Brendan McMahan, Virginia Smith, and Ameet Talwalkar. LEAF: A benchmark for federated settings. *CoRR*, abs/1812.01097, 2018.
- Nicholas Carlini, Steve Chien, Milad Nasr, Shuang Song, Andreas Terzis, and Florian Tramèr. Membership inference attacks from first principles. *CoRR*, abs/2112.03570, 2021a.
- Nicholas Carlini, Florian Tramèr, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom B. Brown, Dawn Song, Úlfar Erlingsson, Alina Oprea, and Colin Raffel. Extracting training data from large language models. In Michael Bailey and Rachel Greenstadt (eds.), *30th USENIX Security Symposium*, pp. 2633–2650, 2021b.
- Zachary Charles, Zachary Garrett, Zhouyuan Huo, Sergei Shmulyian, and Virginia Smith. On large-cohort training for federated learning, 2021.
- Kamalika Chaudhuri, Chuan Guo, and Mike Rabbat. Privacy-aware compression for federated data analysis, 2022.
- Chia-Yu Chen, Jungwook Choi, Daniel Brand, Ankur Agrawal, Wei Zhang, and Kailash Gopalakrishnan. AdaComp: Adaptive residual gradient compression for data-parallel distributed training. In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence (AAAI)*, pp. 2827–2835, 2018.
- Mingqing Chen, Rajiv Mathews, Tom Ouyang, and Françoise Beaufays. Federated learning of out-of-vocabulary words, 2019.

- Matthieu Courbariaux, Yoshua Bengio, and Jean-Pierre David. Binaryconnect: Training deep neural networks with binary weights during propagations, 2015.
- Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. Calibrating noise to sensitivity in private data analysis. In *Proceedings of the Third Conference on Theory of Cryptography*, pp. 265–284, 2006.
- FCC. The eleventh Measuring Broadband America fixed broadband report, 2021. URL <https://www.fcc.gov/reports-research/reports/measuring-broadband-america/measuring-fixed-broadband-eleventh-report>.
- Jonas Geiping, Hartmut Bauermeister, Hannah Dröge, and Michael Moeller. Inverting gradients—How easy is it to break privacy in federated learning? In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 33, pp. 16937–16947, 2020.
- Alec Go, Richa Bhayani, and Lei Huang. Twitter sentiment classification using distant supervision. CS224N Project Report, Stanford, 2009.
- Suyog Gupta, Ankur Agrawal, Kailash Gopalakrishnan, and Pritish Narayanan. Deep learning with limited numerical precision. In *Proceedings of the 32nd International Conference on Machine Learning (ICML)*, volume 37, pp. 1737–1746, 2015.
- Pengchao Han, Shiqiang Wang, and Kin K. Leung. Adaptive gradient sparsification for efficient federated learning: An online learning approach. In *40th IEEE International Conference on Distributed Computing Systems (ICDCS)*, pp. 300–310, 2020.
- Babak Hassibi and David G. Stork. Second order derivatives for network pruning: Optimal Brain Surgeon. In *Advances in Neural Information Processing Systems*, volume 5, pp. 164–171, 1992.
- Chaoyang He, Murali Annavaram, and Salman Avestimehr. Group knowledge transfer: Federated learning of large CNNs at the edge. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 33, pp. 14068–14080, 2020.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, (CVPR)*, pp. 770–778, 2016.
- Chuang Hu, Wei Bao, Dan Wang, and Fengming Liu. Dynamic adaptive DNN surgery for inference acceleration on the edge. *IEEE INFOCOM 2019 - IEEE Conference on Computer Communications*, pp. 1423–1431, 2019.
- Dzmitry Huba, John Nguyen, Kshitiz Malik, Ruiyu Zhu, Mike Rabbat, Ashkan Yousefpour, Carole-Jean Wu, Hongyuan Zhan, Pavel Ustinov, Harish Srinivas, Kaikai Wang, Anthony Shoumikhin, Jesik Min, and Mani Malek. Papaya: Practical, private, and scalable federated learning. In *Proceedings of Conference on Systems and Machine Learning Foundation (MLSys)*, 2022.
- Benoit Jacob, Skirmantas Kligys, Bo Chen, Menglong Zhu, Matthew Tang, Andrew Howard, Hartwig Adam, and Dmitry Kalenichenko. Quantization and training of neural networks for efficient integer-arithmetic-only inference. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 2704–2713, June 2018.
- Hervé Jégou, Matthijs Douze, and Cordelia Schmid. Product quantization for nearest neighbor search. *IEEE Trans. Pattern Anal. Mach. Intell.*, 33(1):117–128, 2011.
- Yuang Jiang, Shiqiang Wang, Bong Jun Ko, Wei-Han Lee, and Leandros Tassioulas. Model pruning enables efficient federated learning on edge devices. *CoRR*, abs/1909.12326, 2019.
- Peter Kairouz, H. Brendan McMahan, Brendan Avent, Aurélien Bellet, Mehdi Bennis, Arjun Nitin Bhagoji, Kallista Bonawitz, Zachary Charles, Graham Cormode, Rachel Cummings, Rafael G. L. D’Oliveira, Hubert Eichner, Salim El Rouayheb, David Evans, Josh Gardner, Zachary Garrett, Adrià Gascón, Badih Ghazi, Phillip B. Gibbons, Marco Gruteser, Zaid Harchaoui, Chaoyang He, Lie He, Zhouyuan Huo, Ben Hutchinson, Justin Hsu, Martin Jaggi, Tara Javidi, Gauri Joshi,

- Mikhail Khodak, Jakub Konečný, Aleksandra Korolova, Farinaz Koushanfar, Sanmi Koyejo, Tancrède Lepoint, Yang Liu, Prateek Mittal, Mehryar Mohri, Richard Nock, Ayfer Özgür, Rasmus Pagh, Mariana Raykova, Hang Qi, Daniel Ramage, Ramesh Raskar, Dawn Song, Weikang Song, Sebastian U. Stich, Ziteng Sun, Ananda Theertha Suresh, Florian Tramèr, Praneeth Vepakomma, Jianyu Wang, Li Xiong, Zheng Xu, Qiang Yang, Felix X. Yu, Han Yu, and Sen Zhao. Advances and open problems in federated learning, 2019.
- Jakub Konečný, H. Brendan McMahan, Felix X. Yu, Peter Richtarik, Ananda Theertha Suresh, and Dave Bacon. Federated learning: Strategies for improving communication efficiency. In *NIPS Workshop on Private Multi-Party Machine Learning*, 2016.
- Raghuraman Krishnamoorthi. Quantizing deep convolutional networks for efficient inference: A whitepaper. *CoRR*, abs/1806.08342, 2018.
- Yann Le Cun, John S. Denker, and Sara A. Solla. Optimal brain damage. In *Proceedings of the 2nd International Conference on Neural Information Processing Systems*, pp. 598–605, 1989.
- Yann LeCun and Corinna Cortes. MNIST handwritten digit database. 2010. URL <http://yann.lecun.com/exdb/mnist/>.
- Yujun Lin, Song Han, Huizi Mao, Yu Wang, and William Dally. Deep gradient compression: Reducing the communication bandwidth for distributed training. In *6th International Conference on Learning Representations (ICLR)*, 2018.
- Xiaorui Liu, Yao Li, Jiliang Tang, and Ming Yan. A double residual compression algorithm for efficient distributed learning. In *The 23rd International Conference on Artificial Intelligence and Statistics (AISTATS)*, volume 108 of *PMLR*, pp. 133–143, 2020.
- Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *IEEE International Conference on Computer Vision (ICCV)*, pp. 3730–3738, 2015.
- Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Agüera y Arcas. Communication-efficient learning of deep networks from decentralized data. In *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics (AISTATS)*, volume 54 of *PMLR*, pp. 1273–1282, 2017.
- Luca Melis, Congzheng Song, Emiliano De Cristofaro, and Vitaly Shmatikov. Exploiting unintended feature leakage in collaborative learning. In *2019 IEEE Symposium on Security and Privacy (S&P)*, pp. 691–706, 2019.
- John Nguyen, Kshitiz Malik, Hongyuan Zhan, Ashkan Yousefpour, Michael Rabbat, Mani Malek, and Dzmityr Huba. Federated learning with buffered asynchronous aggregation. In *Proceedings of The 25th International Conference on Artificial Intelligence and Statistics (AISTATS)*, volume 151 of *PMLR*, pp. 3581–3607, 2022.
- Titouan Parcollet, Javier Fernandez-Marques, Pedro PB Gusmao, Yan Gao, and Nicholas D. Lane. ZeroFL: Efficient on-device training for federated learning with local sparsity. In *International Conference on Learning Representations (ICLR)*, April 2022.
- Constantin Philippenko and Aymeric Dieuleveut. Preserved central model for faster bidirectional compression in distributed settings. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 34, pp. 2387–2399, 2021.
- Xinchi Qiu, Titouan Parcollet, Javier Fernandez-Marques, Pedro Porto Buarque de Gusmao, Daniel J Beutel, Taner Topal, Akhil Mathur, and Nicholas D Lane. A first look into the carbon footprint of federated learning. *arXiv preprint arXiv:2102.07627*, 2021.
- Cedric Renggli, Saleh Ashkboos, Mehdi Aghagolzadeh, Dan Alistarh, and Torsten Hoefler. SparCML: High-performance sparse communication for machine learning. In *SC'19: Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*, 2019.

- Daniel Rothchild, Ashwinee Panda, Enayat Ullah, Nikita Ivkin, Ion Stoica, Vladimir Braverman, Joseph Gonzalez, and Raman Arora. FetchSGD: Communication-efficient federated learning with sketching. In *Proceedings of the 37th International Conference on Machine Learning (ICML)*, volume 119 of *PMLR*, pp. 8253–8265, 2020.
- Felix Sattler, Simon Wiedemann, Klaus-Robert Müller, and Wojciech Samek. Robust and communication-efficient federated learning from non-i.i.d. data. *IEEE Transactions on Neural Networks and Learning Systems*, 31:3400–3413, 2020.
- Frank Seide, Hao Fu, Jasha Droppo, Gang Li, and Dong Yu. 1-bit stochastic gradient descent and its application to data-parallel distributed training of speech DNNs. In *INTERSPEECH 2014, 15th Annual Conference of the International Speech Communication Association*, pp. 1058–1062, 2014.
- Pierre Stock, Armand Joulin, Rémi Gribonval, Benjamin Graham, and Hervé Jégou. And the bit goes down: Revisiting the quantization of neural networks. In *International Conference on Learning Representations (ICLR)*, 2020.
- Hanlin Tang, Chen Yu, Xiangru Lian, Tong Zhang, and Ji Liu. DOUBLESQUEEZE: Parallel stochastic gradient descent with double-pass error-compensated compression. In *Proceedings of the 36th International Conference on Machine Learning (ICML)*, volume 97 of *PMLR*, pp. 6155–6165, 2019.
- Aleksei Triastcyn, Matthias Reisser, and Christos Louizos. DP-REC: Private & communication-efficient federated learning, 2021.
- Praneeth Vepakomma, Otkrist Gupta, Tristan Swedish, and Ramesh Raskar. Split learning for health: Distributed deep learning without sharing raw patient data, 2018.
- Thijs Vogels, Sai Praneeth Karimireddy, and Martin Jaggi. PowerSGD: Practical low-rank gradient compression for distributed optimization. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 32, pp. 14236–14245, 2019.
- Hongyi Wang, Scott Sievert, Shengchao Liu, Zachary Charles, Dimitris Papailiopoulos, and Stephen Wright. ATOMO: Communication-efficient learning via atomic sparsification. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 31, 2018.
- Jianyu Wang, Hang Qi, Ankit Singh Rawat, Sashank Reddi, Sagar Waghmare, Felix X. Yu, and Gauri Joshi. Fedlite: A scalable approach for federated learning on resource-constrained clients, 2022.
- Jianqiao Wangni, Jialei Wang, Ji Liu, and Tong Zhang. Gradient sparsification for communication-efficient distributed optimization. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 31, 2018.
- Lauren Watson, Chuan Guo, Graham Cormode, and Alexandre Sablayrolles. On the importance of difficulty calibration in membership inference attacks. In *International Conference on Learning Representations (ICLR)*, 2022.
- Wei Wen, Cong Xu, Feng Yan, Chunpeng Wu, Yandan Wang, Yiran Chen, and Hai Li. TernGrad: Ternary gradients to reduce communication in distributed deep learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 30, 2017.
- Chien-Sheng Yang, Jinhyun So, Chaoyang He, Songze Li, Qian Yu, and Salman Avestimehr. Light-SecAgg: Rethinking secure aggregation in federated learning. In *Proceedings of Conference on Systems and Machine Learning Foundation (MLSys)*, 2022.
- Mingchao Yu, Zhifeng Lin, Krishna Narra, Songze Li, Youjie Li, Nam Sung Kim, Alexander Schwing, Murali Annavaram, and Salman Avestimehr. GradiVeQ: Vector quantization for bandwidth-efficient gradient aggregation in distributed CNN training. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems (NeurIPS)*, pp. 5129–5139, 2018.
- Shuai Zheng, Ziyue Huang, and James T. Kwok. Communication-efficient distributed blockwise momentum SGD with error-feedback. In *Proceedings of the 33rd International Conference on Neural Information Processing Systems (NeurIPS)*, pp. 11450–11460, 2019.

Shuchang Zhou, Zekun Ni, Xinyu Zhou, He Wen, Yuxin Wu, and Yuheng Zou. DoReFa-Net: Training low bitwidth convolutional neural networks with low bitwidth gradients. *CoRR*, abs/1606.06160, 2016.

A APPENDIX

A.1 EXPERIMENTAL DETAILS

In this section, we provide further details of the experimental setup described in Section 5.1 and the hyper-parameters used for all the runs in Table 1. For all the tasks, we use a mini-batch SGD optimizer for local training at the client and FEDAVG optimizer for global model update on the server. The LEAF benchmark is released under the BSD 2-Clause License.

Baselines. We run hyper-parameter sweeps to tune the client and server learning rates for the uncompressed baselines. Then, we keep the same hyper-parameters in all the runs involving uplink compression. We have observed that tuning the hyper-parameters for each compression factor does not provide significantly different results than using those for the uncompressed baselines, in addition to the high cost of model training involved.

Compression details. For scalar quantization, we use per-tensor quantization with MinMax observers. We use the symmetric quantization scheme over the integer range $[-2^{b-1}, 2^{b-1} - 1]$. For pruning, we compute the random mask separately for each tensor, ensuring all pruned layers have the same target sparsity in their individual updates. For product quantization, we explore various configurations by choosing the number of codewords per layer k in $\{8, 16, 32, 64\}$ and the block size d in $\{4, 9, 18\}$. We automatically adapt the block size for each layer to be the largest allowed one that divides C_{in} (in the fully connected case).

Table 1: Hyper-parameters used for all the experiments including baselines. η is the learning rate.

Dataset	Users per round	Client epochs	Max. server epochs	η_{SGD}	η_{FedAvg}
CelebA	100	1	30	0.90	0.08
Sent140	100	1	10	5.75	0.24
FEMNIST	5	1	5	0.01	0.24

A.2 EXPERIMENTAL RESULTS

We provide various additional experimental results that are referred to in the main paper.

A.2.1 TABLES CORRESPONDING TO FIGURE 2

We provide the detailed results corresponding to Figure 2 along with standard deviation over 3 runs in Tables 4, 5, and 6.

A.2.2 AGGREGATION OVERFLOWS WITH SCALAR QUANTIZATION

We discussed the challenge of aggregation overflows of quantized values with restricted SECAGG finite group size in Section 3.2.1 and noted in Section 4.2.1 that it suffices for SECAGG bit-width p to be greater than quantization bit-width b by at most $\lceil \log_2 N \rceil$, where N is the number of clients participating in a given round. However, the overflow margin increases the client update size by $p - b$ per weight. To optimize this further, we explore the impact of p on aggregation overflows and accuracy, and present the results in Table 2. As expected, we observe a decrease in percentage of weights that overflow during aggregation with the increase in the overflow margin size. However, while there is some benefit to non-zero overflow margin size, there is no strong correlation between the overflow margin size and accuracy, indicating the potential to achieve better utility even in the presence of overflows.

A.2.3 WEIGHTED AGGREGATION AND SCALAR QUANTIZATION

Following the setup of Nguyen et al. (2022), we weight each client update by the number of samples the client trained on. Denoting the weight associated with the client i with ω_i and following the same notations as in Section 4.1, weighted update is obtained as $h_i = (q(g_i) \times \omega_i) + m_i$. Since this is

Table 2: Percentage of aggregation overflows (among all model parameters) for the CelebA dataset over various SQ configurations. b is Quantization bit-width, p is SECAGG bit-width, $p - b$ is overflow margin size in bits.

b	p	$p - b$	Overflows (% of parameters)	Accuracy
4	4	0	3.71±1.53	49.33±2.03
4	5	1	1.43±0.55	50.44±1.77
4	6	2	0.68±0.43	49.67±1.56
4	7	3	0.17±0.12	51.58±0.66
4	8	4	0.06±0.00	87.30±0.36
4	9	5	0.06±0.00	89.19±0.20
4	10	6	0.06±0.00	88.52±0.07
4	11	7	0.05±0.00	87.68±1.24
8	8	0	2.28±0.11	82.11±0.90
8	9	1	1.06±0.06	90.49±0.27
8	10	2	0.39±0.04	90.97±0.50
8	11	3	0.14±0.01	91.08±0.45
8	12	4	0.06±0.00	91.29±0.13
8	13	5	0.04±0.00	90.49±0.93
8	14	6	0.02±0.00	91.31±0.24
8	15	7	0.01±0.00	91.19±0.33

Table 3: Cost of broadcasting codebooks (for downlink communications) is negligible compared to model sizes. Recall that k denotes the number of codebooks and d the block size.

Dataset	Codebook size k	Block size d	Codebooks size (% of model size)
CelebA	8	4	0.6 KB (0.5%)
	8	18	2.5 KB (2.2%)
	64	4	4.2 KB (3.7%)
	64	18	14.6 KB (12.8%)
Sent140	8	4	0.9 KB (0.0%)
	8	18	2.3 KB (0.0%)
	64	4	5.4 KB (0.0%)
	64	18	15.4 KB (0.1%)
FEMNIST	8	4	2.6 KB (0.0%)
	8	18	11.2 KB (0.0%)
	64	4	20.8 KB (0.0%)
	64	18	89.8 KB (0.2%)

a synchronous FL setup, we do not set staleness factor. This weighted aggregation has no impact on pruning and product quantization, but can lead to overflows with scalar quantization. Therefore, we skip the weighting of quantized parameters of client updates and only weight non-quantized parameters (such as bias). For completion, we study with unweighted aggregation of client updates (including bias parameters) for scalar quantization experiments and present the result in Table 7. As expected, these results are similar to the ones with weighted aggregation.

A.2.4 PQ CODEBOOK SIZE IS NEGLIGIBLE

We demonstrate in Table 3 that the overhead of sending codebooks (for all layers) is negligible compared to the model size. When the model is very small (CelebA model is 114 KB), reducing k and d makes the overhead negligible without hurting performance.

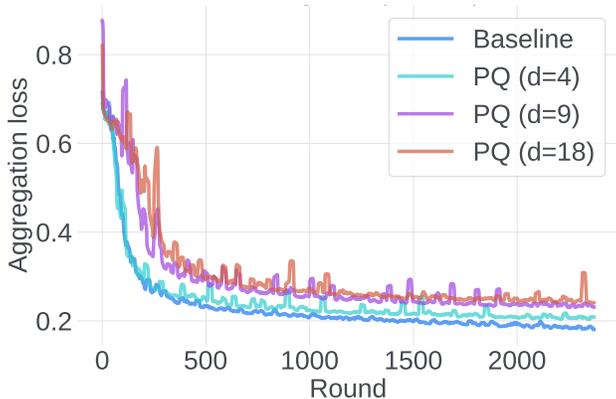


Figure 4: Number of rounds to convergence is similar for PQ-compressed runs compared to the non-compressed baseline (on CelebA). Note that outside a simulated environment, the wall-clock time convergence for PQ runs would be lower than the baseline since uplink communications would be faster.

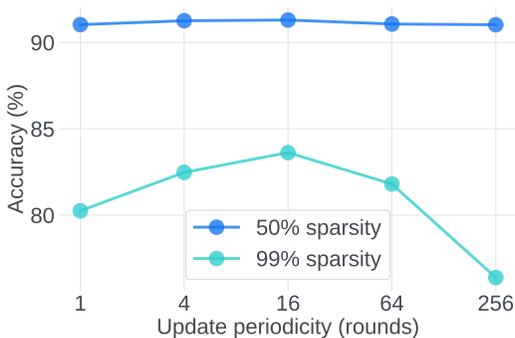


Figure 5: Impact of pruning mask refresh intervals on model performance for the CelebA dataset. Note that the effect of refreshing the pruning masks is more apparent at higher sparsity levels, and generalization performance decreases when masks are stale for longer during training.

A.2.5 CONVERGENCE CURVES

We also provide convergence curves for PQ-compressed and baseline runs to demonstrate similar number of rounds needed to convergence in Figure 4.

A.2.6 PERFORMANCE IMPACT OF SPARSITY MASK REFRESH

In addition to scalar and product quantization as described in Section 5.3, we also conduct experiments with varying the interval for refreshing pruning masks. We consider two levels of sparsity, 50% and 99% and our experiments are on the CelebA dataset. We present our results in Figure 5. Overall we find that the model accuracy is robust to the update periodicity unless at very high sparsities, where accuracy decreases when mask refresh periodicity increases. This is important for future directions such as in asynchronous FL where clients have to maintain the same mask across successive global updates.

A.3 SECIND IMPLEMENTATIONS

SECIND can be extended to other settings, such as multi-party computation (using two or more servers to operate on shares of the input), where each client can send evaluations of *distributed point functions* to encode each assignment (Boyle et al., 2016). These are represented compactly, but

may require longer codewords to overcome the overheads. We leave the study of such software implementations of SECIND to future work.

Table 4: Results of client update compression with SECAGG-compatible scalar quantization on LEAF datasets over three runs. We fix p across runs as this defines the uplink size, but not b . We pick the run with the best accuracy and report the corresponding b .

Dataset	b	p	Uplink size (in KB)	Compression factor	Accuracy
CelebA Baseline: 91.2±0.2	1	1	4.3	26.6	49.3±2.7
	1,2	2	7.9	14.6	51.8±0.4
	2,3	3	11.4	10.0	52.0±0.7
	1,4	4	15.0	7.6	51.8±0.5
	3,4	5	18.5	6.2	53.9±1.7
	1,3,5	6	22.1	5.2	52.6±0.6
	1,2	7	25.6	4.5	52.8±1.3
	6	8	29.2	3.9	89.6±0.2
	6,7	9	32.7	3.5	91.2±0.1
	6,7	10	36.3	3.2	91.2±0.3
	6,7,8	11	39.8	2.9	91.1±0.1
	6,8	12	43.4	2.6	91.4±0.0
	6,7	13	46.9	2.4	91.3±0.2
	7,8	14	50.5	2.3	91.3±0.1
	8	15	54.0	2.1	91.2±0.2
Sent140 Baseline: 70.8±0.4	1	1	399.3	31.5	46.2±0.0
	1	2	792.3	15.9	53.8±0.0
	2,3	3	1185.4	10.6	53.8±0.0
	1,2	4	1578.4	8.0	53.8±0.0
	2,3,5	5	1971.4	6.4	53.8±0.0
	2,6	6	2364.5	5.3	53.8±0.0
	1,7	7	2757.5	4.6	53.8±0.0
	2,5	8	3150.6	4.0	53.8±0.0
	5,6,7	9	3543.6	3.6	53.9±0.0
	4,6	10	3936.6	3.2	53.8±0.1
	6,8	11	4329.7	2.9	53.8±0.1
	5,7	12	4722.7	2.7	51.3±4.4
	6,8	13	5115.7	2.5	50.6±2.9
	7	14	5508.8	2.3	63.1±3.1
	8	15	5901.8	2.1	68.5±1.6
FEMNIST Baseline: 84.8±0.7	1	1	1404.0	31.2	2.2±3.0
	1,2	2	2770.3	15.8	2.4±1.5
	3	3	4136.5	10.6	12.2±3.9
	4	4	5502.8	8.0	65.0±3.3
	5	5	6869.0	6.4	80.7±0.4
	6	6	8235.3	5.3	83.8±0.2
	5,6,7	7	9601.5	4.6	84.3±0.4
	6,7	8	10967.8	4.0	85.1±0.1
	6,8	9	12334.1	3.6	84.9±0.2
	7,8	10	13700.3	3.2	85.0±0.3
	8	11	15066.6	2.9	84.6±0.3

Table 5: Results of client update compression with SECAGG-compatible random mask pruning on LEAF datasets

Dataset	Sparsity	Uplink size (in KB)	Compression factor	Accuracy
CelebA Baseline: 91.2±0.2	0.1	103.0	1.1	91.2±0.2
	0.2	91.6	1.2	91.3±0.0
	0.3	80.3	1.4	91.1±0.2
	0.4	68.9	1.7	91.1±0.2
	0.5	57.5	2.0	91.1±0.1
	0.6	46.1	2.5	91.2±0.1
	0.7	34.8	3.3	91.1±0.1
	0.8	23.4	4.9	90.9±0.2
	0.9	12.0	9.5	90.4±0.2
	0.91	10.9	10.5	90.4±0.1
	0.92	9.8	11.7	90.3±0.2
	0.93	8.6	13.3	89.9±0.2
	0.94	7.5	15.3	90.0±0.3
	0.95	6.3	18.1	89.6±0.1
	0.96	5.2	22.0	89.0±0.4
	0.97	4.1	28.2	88.9±0.1
	0.98	2.9	39.1	87.0±0.5
0.99	1.8	64.1	83.8±0.2	
Sent140 Baseline: 70.8±0.4	0.1	11325.7	1.1	70.5±0.3
	0.2	10068.0	1.2	70.6±0.5
	0.3	8810.3	1.4	70.6±0.3
	0.4	7552.6	1.7	70.6±1.6
	0.5	6294.9	2.0	70.6±0.8
	0.6	5037.1	2.5	70.3±0.9
	0.7	3779.4	3.3	70.5±0.7
	0.8	2521.7	5.0	67.9±0.7
	0.9	1264.0	10.0	68.6±1.8
	0.91	1138.2	11.1	69.4±1.1
	0.92	1012.4	12.4	66.4±3.1
	0.93	886.7	14.2	65.8±2.5
	0.94	760.9	16.5	60.8±0.0
	0.95	635.1	19.8	51.6±5.7
	0.96	509.4	24.7	52.7±3.9
	0.97	383.6	32.8	58.2±7.8
	0.98	257.8	48.8	60.3±5.9
0.99	132.0	95.3	49.3±4.1	
FEMNIST Baseline: 84.8±0.7	0.1	39384.2	1.1	84.6±0.3
	0.2	35010.3	1.2	84.7±0.3
	0.3	30636.4	1.4	84.7±0.2
	0.4	26262.5	1.7	84.6±0.2
	0.5	21888.5	2.0	84.3±0.3
	0.6	17514.7	2.5	84.2±0.2
	0.7	13140.8	3.3	83.6±0.6
	0.8	8766.9	5.0	83.1±0.2
	0.9	4393.0	10.0	81.7±0.4
	0.91	3955.6	11.1	81.1±0.4
	0.92	3518.2	12.4	80.6±0.2
	0.93	3080.8	14.2	80.6±0.2
	0.94	2643.4	16.6	80.2±0.3
	0.95	2206.0	19.8	79.4±0.1
	0.96	1768.6	24.7	78.4±0.3
	0.97	1331.2	32.9	77.0±0.5
	0.98	893.9	49.0	73.3±0.4
0.99	456.5	95.9	65.6±0.1	

Table 6: Results of client update compression with Product quantization and SECIND on LEAF datasets

Dataset	k	d	Uplink size (in KB)	Compression factor	Accuracy
CelebA Baseline: 91.2±0.2	8	4	3.4	33.2	91.0±0.2
	8	9	1.9	58.9	89.9±0.2
	8	18	1.4	83.7	89.1±0.4
	16	4	4.3	26.3	91.2±0.0
	16	9	2.3	49.0	90.7±0.0
	16	18	1.6	73.0	89.7±0.4
	32	4	5.2	21.8	91.4±0.1
	32	9	2.7	42.2	90.9±0.3
	32	18	1.8	65.2	90.1±0.1
	64	4	6.1	18.7	91.1±0.3
	64	9	3.1	37.1	90.9±0.1
	64	18	1.9	58.9	90.5±0.4
Sent140 Baseline: 70.8±0.4	8	4	301.0	41.8	69.5±1.5
	8	9	204.2	61.6	69.0±0.8
	8	18	86.3	145.8	66.4±0.3
	16	4	399.3	31.5	69.8±1.1
	16	9	270.2	46.6	69.3±0.1
	16	18	113.0	111.3	67.7±0.6
	32	4	497.5	25.3	70.6±0.2
	32	9	336.2	37.4	69.4±0.4
	32	18	139.7	90.1	67.7±2.3
	64	4	595.8	21.1	70.7±0.3
	64	9	402.2	31.3	70.1±1.0
	64	18	166.4	75.6	68.7±0.7
FEMNIST Baseline: 84.8±0.7	8	4	1063.3	41.2	84.4±0.4
	8	9	494.6	88.5	82.5±0.2
	8	18	266.4	164.3	81.5±0.4
	16	4	1405.1	31.1	84.7±0.2
	16	9	646.8	67.6	83.3±0.1
	16	18	342.6	127.7	82.2±0.5
	32	4	1747.0	25.0	84.7±0.3
	32	9	799.1	54.8	83.9±0.6
	32	18	418.8	104.5	83.1±0.5
	64	4	2088.9	20.9	84.4±0.2
	64	9	951.4	46.0	83.8±0.8
	64	18	495.1	88.4	83.5±0.7

Table 7: Results of scalar quantization on LEAF datasets with unweighted client update aggregation over three runs. We fix p across runs as this defines the uplink size, but not b . We pick the run with the best accuracy and report the corresponding b .

Dataset	b	p	Uplink size (in KB)	Compression factor	Accuracy
CelebA Baseline: 91.2 ± 0.2	1	1	4.3	26.5	50.3 ± 1.92
	1,2	2	7.9	14.6	50.9 ± 0.54
	1,2,3	3	11.4	10.0	52.0 ± 0.46
	2,3,4	4	15.0	7.6	51.8 ± 0.17
	1,3,4	5	18.5	6.2	52.6 ± 1.23
	3,4	6	22.1	5.2	52.8 ± 1.21
	2,3	7	25.6	4.5	51.9 ± 0.05
	6	8	29.2	3.9	90.2 ± 0.19
	6	9	32.7	3.5	90.8 ± 0.04
	6	10	36.3	3.2	91.2 ± 0.08
	6,7	11	39.8	2.9	91.2 ± 0.20
	6,8	12	43.4	2.6	91.2 ± 0.11
	6	13	46.9	2.4	91.4 ± 0.13
	7,8	14	50.5	2.3	91.3 ± 0.15
	8	15	54.0	2.1	91.3 ± 0.22
Sent140 Baseline: 70.8 ± 0.4	1	1	399.3	31.5	51.3 ± 4.4
	1,2	2	792.3	15.9	51.3 ± 4.4
	1,2	3	1185.4	10.6	53.8 ± 0.0
	1,2	4	1578.4	8.0	53.8 ± 0.0
	2,3	5	1971.4	6.4	53.8 ± 0.0
	1,5	6	2364.5	5.3	53.8 ± 0.0
	1,3,4	7	2757.5	4.6	53.8 ± 0.0
	2,4,7	8	3150.6	4.0	53.8 ± 0.1
	3,5	9	3543.6	3.6	53.8 ± 0.0
	4,5	10	3936.6	3.2	53.8 ± 0.1
	5,7	11	4329.7	2.9	52.0 ± 3.2
	6,7	12	4722.7	2.7	53.8 ± 0.0
	6	13	5115.7	2.5	60.3 ± 2.9
	7	14	5508.8	2.3	65.9 ± 2.9
	8	15	5901.8	2.1	67.7 ± 0.9
FEMNIST Baseline: 84.8 ± 0.7	1	1	1404.0	31.2	2.0 ± 1.3
	1,2	2	2770.3	15.8	4.8 ± 2.8
	1,2,3	3	4136.5	10.6	13.6 ± 2.4
	4	4	5502.8	8.0	66.3 ± 3.5
	5	5	6869.0	6.4	79.7 ± 0.8
	5	6	8235.3	5.3	84.0 ± 0.9
	6,7	7	9601.5	4.6	84.3 ± 0.1
	6,7,8	8	10967.8	4.0	84.8 ± 0.6
	7,8	9	12334.1	3.5	85.1 ± 0.1
	7,8	10	13700.3	3.2	85.0 ± 0.4
	8	11	15066.6	2.9	83.5 ± 1.7