

---

# Bi-Directional Communication-Efficient Stochastic FL via Remote Source Generation

---

Maximilian Egger\*, Rawad Bitar, Antonia Wachter-Zeh

Technical University of Munich

Munich, Germany

{maximilian.egger, rawad.bitar, antonia.wachter-zeh}@tum.de

Nir Weinberger

Israel Institute of Technology

Haifa, Israel

nirwein@technion.ac.il

Deniz Gündüz

Imperial College London

London, United Kingdom

d.gunduz@imperial.ac.uk

## Abstract

Federated Learning (FL) incurs high communication costs in both uplink and downlink. The literature largely focuses on lossy compression of model updates in deterministic FL. In contrast, stochastic (Bayesian) FL considers distributions over parameters, enabling uncertainty quantification, better generalization, and, crucially, inherent communication-regularized training through a mirror-descent structure. In this paper, we consider both uplink and downlink communication in stochastic FL, and propose a communication framework based on remote source generation. Employing Minimal Random Coding (MRC) for remote generation, we allow the server and the clients to sample from local and global posteriors (sources), respectively, rather than transmitting locally sampled updates. The framework encompasses communication-regularized local optimization and principled compression of model updates, leveraging gradually updated prior distributions as side information. Through extensive simulations, we show that our method achieves  $5 - 32\times$  reduction in total communication cost while preserving accuracy. We further analyze the communication cost, refining existing MRC bounds and enabling precise quantification of uplink and downlink trade-offs. We also extend our method to conventional FL via stochastic quantization and prove a contraction property for the biased MRC compressor to facilitate convergence analysis.

## 1 Introduction

Federated learning (FL) enables collaborative machine learning (ML) across multiple clients orchestrated by a central federator [McMahan et al., 2017]. Communication efficiency, privacy, security, and data heterogeneity are well-established challenges in FL [Zhang et al., 2021, Wen et al., 2023]. As a bi-directional process, FL requires substantial *uplink* and *downlink* communication, posing increasing pressure on communication networks as ML models grow larger. To address this, lossy compression techniques have been widely adopted to reduce uplink gradient transmissions and downlink model broadcasts [Seide et al., 2014, Alistarh et al., 2017, Philippenko and Dieuleveut, 2020, Gruntkowska et al., 2023]. However, these methods almost exclusively focus on conventional (non-stochastic) FL, where clients train deterministic models and transmit fixed updates.

Alternatively, stochastic (Bayesian) FL offers improved generalization, robustness, and inherent uncertainty estimation [Zhang et al., 2022, Milasheuski et al., 2025]. Rather than training deterministic

---

\*Corresponding author

models, clients learn local posterior distributions, aggregated by the federator to obtain a global posterior. Recently, [Isik et al., 2024] empirically demonstrated state-of-the-art performance under limited uplink bandwidth using stochastic compression methods, outperforming classical approaches. This framework can be applied to a variety of Bayesian FL solutions such as QSGD Alistarh et al. [2017], QLSD Vono et al. [2022], dithered quantization Abdi and Fekri [2019] and FedPM Isik et al. [2023], as well as to conventional FL settings augmented with stochastic compression.

A key technique enabling stochastic FL is *remote source generation*, which allows the federator to sample from the clients’ local posterior, rather than obtaining samples locally generated by the clients. This avoids redundant transmission and enables tight, stochastic control over communication. Such remote generation requires common randomness shared between the transmitter and receiver in the form of a common prior Li [2024], which we also refer to as *side information*. If the downlink is unlimited, this allows the server to broadcast the global posterior to all the clients, and this posterior serves as a natural common prior. However, when both uplink and downlink are limited, the possibility for remote source generation is restricted, which challenges the application of efficient stochastic FL. Thus, in this paper, we explore and analyze communication-efficient stochastic FL. The rigorous treatment of stochastic FL in this scenario is further reinforced by its two fundamental advantages: (i) *Communication-regularized training*: we show that stochastic FL, as opposed to conventional solutions, inherently integrates communication constraints into the training process, effectively treating communication as an integral part of the optimization objective; (ii) *Priors as side information*: the probabilistic structure allows principled integration of common priors as side information, reducing communication costs to the update from prior to posterior.

Concretely, we address the following question: *Can joint uplink and downlink compression significantly reduce communication costs in stochastic FL without compromising accuracy?* We answer this question affirmatively, tightening the communication–accuracy trade-off in ways that deterministic methods cannot, and achieving **communication reductions of up to  $32\times$**  without performance loss. Below, we summarize our key contributions.

- We propose two novel bi-directional compression algorithms for stochastic FL with Minimal Random Coding (MRC): one leveraging globally shared randomness, and one requiring private shared randomness between each client and the federator. Both enable efficient sampling-based communication by exploiting carefully selected side information.
- We demonstrate substantial communication savings, reducing total cost by factors of 5 – 32 while maintaining competitive accuracy across standard benchmarks. Our ablation studies analyze the effects of shared randomness and the choice of side information.
- We extend our approach to conventional FL with stochastic quantization, proving a contraction property of the resulting (biased) compression operator. This enables convergence guarantees in both uni- and bi-directional settings.
- We develop a theoretical framework for communication analysis in stochastic FL, quantifying uplink and downlink costs under MRC. Our results refine and generalize bounds from Chatterjee and Diaconis [2018], including tight analyses for Bernoulli distributions and tools applicable to broader distribution classes.

## 2 Preliminaries: Stochastic FL with Bi-Directional Compression

In this section, we shortly review the concepts of stochastic FL and compression based on MRC, which are employed in our proposed stochastic bi-directional algorithm.

**Stochastic FL.** A set of  $n$  clients collaboratively and iteratively train a model, e.g., a neural network, under the orchestration of a federator. We primarily consider cross-silo federated learning with a moderate number of reliable clients [Kairouz and McMahan, 2021]. Client  $i \in [n] := \{1, \dots, n\}$  possesses a dataset  $\mathcal{D}_i$ . We differentiate between homogeneous data, where  $\mathcal{D}_i$  is drawn independently from the same distribution for all clients (i.i.d.), and heterogeneous data, where each  $\mathcal{D}_i$  may come from a different distribution (non i.i.d.). At each training iteration  $t$ , the federator holds a model  $\theta_t$  described by a probability distribution. After downlink transmission, each client  $i$  has an estimate  $\hat{\theta}_{i,t}$  of  $\theta_t$ , and locally optimizes  $\hat{\theta}_{i,t}$  to obtain a local probabilistic model, called *the posterior*  $q_i^t$ . Compressed clients’ posteriors  $q_i^t$  are transmitted back to the federator on the uplink to obtain an estimate  $\hat{q}_i^t$ . The federator aggregates the received posteriors using an aggregation rule  $R(\cdot)$  to obtain

a refined global distribution  $\theta_{t+1} = R(\{\hat{q}_i^t\}_{i \in [n]})$ . A simple aggregation rule  $R(\cdot)$  is the average over all clients' posteriors. This process is repeated until a certain convergence criterion is met. In many stochastic FL settings, the sent client updates  $\hat{q}_i^t$  are just samples from the posterior  $q_i^t$ .

In fact, conventional FL with stochastic quantization can also be described by the procedure above, though with the following differences: (i) the federator holds a model  $\theta_t$  with deterministic parameters; (ii) each client  $i$  locally optimizes  $\hat{\theta}_{i,t}$  to obtain a local gradient  $g_i^t$ . A stochastic compression  $Q_s(\cdot)$  is applied on the client's gradient to obtain a posterior distribution  $q_i^t$  from  $Q_s(g_i^t)$ ; (iii) samples of  $q_i^t$  are transmitted to the federator on the uplink to obtain an estimate of the gradient, which we still denote by  $\hat{q}_i^t$ ; and (iv) the federator updates the global model as  $\theta_{t+1} = \theta_t - \eta R(\{\hat{q}_i^t\}_{i \in [n]})$ , with learning rate  $\eta$ . In this paper, we will investigate both settings, with a prominent focus on the former.

**Stochastic Compression by MRC.** To efficiently transmit samples from the posterior  $q_i^t$ , we employ MRC [Havasi et al., 2019], which allows to leverage shared randomness and side information common to the federator and the clients. MRC is a stochastic compressor  $\mathcal{C}_{\text{mrc}}(\cdot)$ , whose input is a posterior distribution  $Q$  and a prior distribution  $P$ , and its output is a sample from a distribution  $\hat{Q}$  close to  $Q$ . It operates as follows: The encoder and decoder generate  $n_{\text{IS}}$  samples  $\{X_i\}_{i \in [n_{\text{IS}}]}$  from  $P$ . The encoder computes a categorical distribution  $W$ , with  $W(i) \propto Q(X_i)/P(X_i)$ , and transmits an index  $i \sim W$  with  $\log_2(n_{\text{IS}})$  bits. To obtain high accuracy, it is required that  $n_{\text{IS}} = \Theta(\exp(\text{D}_{\text{KL}}(Q||P)))$ , where  $\text{D}_{\text{KL}}(Q||P)$  is the KL-divergence between  $Q$  and  $P$  [Chatterjee and Diaconis, 2018]. For brevity, in what follows, for two Bernoulli distributions with parameters  $q$  and  $p$ , we will use the shorthands  $\text{d}_{\text{KL}}(q||p)$  and  $\mathcal{C}_{\text{mrc}}(q, p)$ .

The choice of MRC for remote generation stems from a practicality aspect and ease of exposition. Ordered random coding (ORC) [Theis and Ahmed, 2022] is a natural advancement, reducing the entropy of the indices selected for transmission using the Gumbel-Max trick. The extension to ORC can be incorporated with our methods using a minor adaption that we omit here for clarity. Theis and Ahmed [2022] further expose an interesting connection to the Poisson functional representation, allowing exact remote generation. Such lossless sampling schemes, however, usually incur substantially larger communication costs and sampling complexities, often rendering those methods impractical. Thus, we herein focus on lossy remote generation schemes for efficiency.

### 3 BiCOMPFL

We next introduce our method BiCOMPFL, a bi-directional stochastic compression strategy, which uses MRC to reduce both uplink and downlink communication costs. The scheme assumes that shared randomness between each of the clients and the federator exists, which can be implemented using pseudo-random sequences generated from a common seed. We distinguish two types of shared randomness: private shared randomness (between individual clients and the federator) and global shared common randomness (among all parties), with the latter being more challenging to achieve in practice. We assume all clients and the federator share the same global model  $\theta_0$  at initialization. This does not incur any communication when global shared randomness is available, but necessitates an initial model transmission from the federator to clients when only private shared randomness exists.

**BiCOMPFL: The General Algorithm.** Our method serves as a general framework for stochastic optimization procedures. We explain BiCOMPFL for Bayesian FL and show in the sequel how it can be used for conventional FL with stochastic quantization. Consider probabilistic mask training (similar to FedPM, [Isik et al., 2023]) as an example of Bayesian FL. The models  $\theta_t \in [0, 1]^d$  of dimension  $d$  are parameters of Bernoulli distributions. Those parameters determine for each weight of a randomly initialized network with fixed weights  $w$  whether it is activated or not. During inference, the weights  $w$  are masked with samples  $x^t \in \{0, 1\}^d \sim \theta_t$ , i.e., the network weights are  $w \odot x^t$ . We start with a general description, which is valid for the cases of global and private shared randomness.

At iteration  $t = 0$ , each client  $i \in [n]$  shares with the federator the same global model, i.e.,  $\hat{\theta}_{i,0} = \theta_0$ , for all  $i \in [n]$ . At iteration  $t$ , each client  $i$  locally trains model  $\hat{\theta}_{i,t}$  in  $L$  local iterations. In our previous example, when training Bernoulli distributions to mask a random network, the parameters are mapped to scores in a dual space, which are then trained for  $L$  local iterations  $m \in [L]$  using stochastic gradient descent. Mapping the trained scores back to the primal space, each client  $i$  obtains a model update in terms of a posterior  $q_i^t$ . We refer to Appendix F for details. This optimization principle is a special instance of mirror descent, which, in the special case of optimizing over Bernoulli

---

**Algorithm 1** BICOMPFL-GR with Global Randomness

---

**Require:** Both clients and federator initialize the same global model  $\theta_0$  using a shared seed

**Ensure:** Clients set prior  $p^t = \hat{\theta}_{i,0} = \theta_0, \forall i \in [n]$

```
1: repeat
2:   for Client  $i \in [n]$  do
3:      $q_i^t \leftarrow$  Local training of  $\hat{\theta}_{i,t}$ 
4:     Employ  $\mathcal{C}_{\text{mrc}}(q_i^t, p^t)$  to sample indices  $I_{i,\ell}^b, \ell \in [n_{\text{UL}}], b \in [B]$  with prior  $p^t$ , transmitted to
       federator to reconstruct  $\hat{q}_i^t$ 
5:   end for
6:   Federator updates global model  $\theta_{t+1} = \frac{1}{n} \sum_{i=1}^n \hat{q}_i^t$ 
7:   Federator relays to client  $j$  the other clients' indices  $\{I_{i,\ell}^b\}_{\ell \in [n_{\text{UL}}], b \in [B], i \in [n] \setminus \{j\}}$ 
8:   for Clients  $i \in [n]$  do
9:     Reconstruct  $\hat{\theta}_{i,t+1} = \frac{1}{n} \sum_{i=1}^n \hat{q}_i^t$  from  $\{I_{i,\ell}^b\}$ 
10:  end for
11:  Clients and federator set prior  $p^t = \hat{\theta}_{t+1}$ 
12:   $t \leftarrow t + 1$ 
13: until Convergence
```

---

distributions, leads to a point-wise minimization with respect to a KL-proximity term (as opposed to the Euclidean distance in standard SGD, cf. Appendix C for details). The KL-divergence between the updated local model and the global model directly determines the communication cost. Hence, we *regularize* the minimization of the loss function by the communication cost, thereby treating communication as an inherent optimization objective.

To convey the model update  $q_i^t$  to the federator, each client employs  $\mathcal{C}_{\text{mrc}}(\cdot)$  in  $B$  blocks of size  $d/B$  each (assuming for simplicity that  $B$  divides  $d$ ) with a prior distribution  $p_{i,u}^t$ , which is set to  $p_{i,u}^0 = \hat{\theta}_{i,0}$  at iteration  $t = 0$ . The choice of  $p_{i,u}^t$  for  $t > 0$  will be clarified later. For each block  $b \in [d/B]$ , client  $i$  conveys  $n_{\text{UL}}$  samples  $\{y_{i,\ell}^t\}_{\ell \in [n_{\text{UL}}]}$  of  $q_i^t$  to the federator by transmitting for each block  $b$  an index  $I_{i,\ell}^b$  with  $\log_2(n_{\text{IS}})$  bits, where  $n_{\text{IS}}$  is the number of samples per block, generated from the prior distribution  $p_{i,u}^t$  at both the client and the federator using the available shared randomness. The samples of all blocks are concatenated for each  $\ell$ . Hence, the federator obtains an estimate of client  $i$ 's posterior distribution using the empirical average  $\hat{q}_i^t = \frac{1}{n_{\text{UL}}} \sum_{\ell=1}^{n_{\text{UL}}} y_{i,\ell}^t$ .

By averaging the estimates  $\hat{q}_i^t$  for all the clients' models, the federator updates the global model as  $\theta_{t+1} = \frac{1}{n} \sum_{i=1}^n \hat{q}_i^t$ . To transmit the new model to each client  $i$ , we assume the existence of a common prior  $p_{i,d}^t$  shared by the federator and the clients. With  $p_{i,d}^t$ , the federator performs MRC in  $B$  blocks of size  $d/B$  to make client  $i$  sample from, and thereby estimate, the latest global model  $\theta_{t+1}$ . The client samples  $n_{\text{DL}}$  masks  $\{x_{i,\ell}^t\}_{\ell \in [n_{\text{DL}}]}$ , each incurring a communication cost of  $B \log_2(n_{\text{IS}})$  bits. An estimate of the updated global model is obtained by concatenating the reconstructed samples for all the blocks  $b \in [B]$ , and averaging over all masks  $\hat{\theta}_{i,t+1} = \frac{1}{n_{\text{DL}}} \sum_{\ell=1}^{n_{\text{DL}}} x_{i,\ell}^t$ .

Since the number of clients is typically large,  $n_{\text{UL}} = 1$  often suffices. The clients' contributions are averaged at the federator, effectively reducing the noise due to the MRC step. This allowed Isik et al. [2024] to theoretically analyze the uplink communication for importance sampling-based stochastic communication of model updates. We will follow a similar approach for downlink communication; however, since downlink communication cannot benefit from the averaging effect of multiple clients, we reduce the variance of the model estimate in the downlink by setting  $n_{\text{DL}} = n \cdot n_{\text{UL}}$ .

The choice of the priors  $p_{i,u}^t$  and  $p_{i,d}^t$  for MRC in the uplink and downlink channels, respectively, crucially affects the performance and the communication cost of the algorithm. As a first-order characterization, the communication cost of MRC is determined by  $D_{\text{KL}}(q_i^t \| p_{i,u}^t)$  in the uplink and by  $D_{\text{KL}}(\theta_{t+1} \| p_{i,d}^t)$  in the downlink. We continue the description with the easier setting in which global shared randomness is available, before turning to the more challenging setting of private randomness.

**Global Randomness.** When global shared randomness is available, all clients can maintain the same priors at each iteration  $t$ , and, thereby, obtain the same global model estimates  $\hat{\theta}_{i,t}$ . The global model is known to the clients and the federator from initialization, and synchronization among all clients

---

**Algorithm 2** BiCOMPFL-PR with Private Randomness

---

**Require:** Both clients and federator initialize the same global model  $\theta_0$  using a shared seed

**Ensure:** Clients set prior  $p_{i,u}^t = p_{i,d}^t = \hat{\theta}_{i,0} = \theta_0, \forall i \in [n]$

```
1: repeat
2:   for Client  $i \in [n]$  do
3:      $q_i^t \leftarrow$  Local training of  $\hat{\theta}_{i,t}$ 
4:     Federator employs  $\mathcal{C}_{\text{mrc}}(q_i^t, p_{i,u}^t)$  to draw  $n_{\text{UL}}$  samples  $y_{i,\ell}^t \sim q_i^t$  using prior  $p_{i,u}^t$ 
5:     Federator estimates client's posterior  $\hat{q}_i^t = \frac{1}{n_{\text{UL}}} \sum_{\ell=1}^{n_{\text{UL}}} y_{i,\ell}^t$ 
6:   end for
7:   Federator updates global model  $\theta_{t+1} = \frac{1}{n} \sum_{i=1}^n \hat{q}_i^t$ 
8:   for Clients  $i \in [n]$  do
9:     Client employs  $\mathcal{C}_{\text{mrc}}(\theta_{t+1}, p_{i,d}^t)$  to draw  $n_{\text{DL}}$  samples  $x_{i,\ell}^t \sim \theta_{t+1}$  using prior  $p_{i,d}^t$ 
10:    Client est. global model:  $\hat{\theta}_{i,t+1} = \frac{1}{n_{\text{DL}}} \sum_{\ell=1}^{n_{\text{DL}}} x_{i,\ell}^t$ 
11:    Clients set prior  $p_{i,u}^t = p_{i,d}^t = \hat{\theta}_{i,t+1}$ 
12:   end for
13:    $t \leftarrow t + 1$ 
14: until Convergence
```

---

is ensured by choosing as prior  $p_{i,u}^t = p_{i,d}^t$  the latest estimate of the global model  $\hat{\theta}_{i,t}$ . The clients utilize the globally shared randomness to sample the exact same samples from the same prior for uplink transmission at all iterations. Selected indices of such samples are transmitted to the federator to convey an estimate  $\hat{q}_i^t$  of the posterior  $q_i^t$ , who reconstructs the global model  $\theta_{t+1}$ . Using the same prior in the downlink, i.e., the global model from the previous iteration, the updated model can be transmitted to the clients through MRC. Leveraging the shared randomness, all clients  $i \in [n]$  sample from the same prior, and thus obtain the exact same estimate of the global model  $\hat{\theta}_{i,t+1} = \hat{\theta}_{t+1}$ , for all  $i \in [n]$ . Hence, we have that  $p_{i,u}^t = p_{i,d}^t = \hat{\theta}_t$  for all  $i \in [n]$ .

In this version, the federator reconstructs the global model from estimates of the client posteriors  $\hat{q}_i^t$ . However, in the uplink, all clients sample from the same prior, which enables further improvements. Naively, the federator will reconstruct the global model using the indices  $I_{i,\ell}^b$  for  $b \in [B], \ell \in [n_{\text{UL}}]$  received by the clients  $i \in [n]$  through MRC, followed by an additional MRC round for downlink transmission. Instead, and more efficiently, the federator can simply relay the indices to the respective other clients (i.e., client  $j$  receives  $I_{i,\ell}^b$  for  $b \in [B], i \in [n] \setminus \{j\}, \ell \in [n_{\text{UL}}]$ ), which reconstruct the updated global model individually. This avoids additional noise by a second compression round and allows better convergence without additional communication facilitated by global randomness. We term this approach BiCOMPFL-GR summarized in Algorithm 1.

**Private Randomness.** Without global randomness, maintaining the same prior among all clients is impossible without additional communication. Instead, an additional round of MRC is needed for the downlink transmission, and each client obtains a different estimate of the global model  $\hat{\theta}_{i,t}$  at each iteration. Hence, the clients' local trainings start from different estimates of the global model. In a non-stochastic setting, this has only been considered by Philippenko and Dieuleveut [2021], Grunkowska et al. [2024]. Understanding the additional cost incurred due to lack of shared randomness in terms of both the convergence speed, communication load, and the choice of the priors  $p_{i,u}^t$  and  $p_{i,d}^t$ , is then of interest.

For the uplink transmission of client  $i$ , any convex combination of  $\hat{\theta}_{i,t}$  and  $\hat{q}_i^t$  can be used as prior, i.e.,  $p_{i,u}^t = \lambda \hat{\theta}_{i,t} + (1 - \lambda) \hat{q}_i^t$ , for some  $0 \leq \lambda \leq 1$  (cf. Appendix I.2 for details). This is due to the availability of both quantities at the federator and client  $i$ . However, small  $\lambda$  values are not expected to reduce the cost of communication reflected by  $d_{\text{KL}}(q_i^t || p_{i,u}^t)$  since the previous global model estimate is likely to be similarly different from the posterior (in terms of the KL-divergence) than the previous posterior estimate of the federator. Indeed, our numerical experiments have shown that the savings from choosing  $\lambda \neq 1$ , i.e., priors other than  $\hat{\theta}_{i,t}$ , are not significant. For simplicity, we thus propose to use  $p_{i,u}^t = p_{i,d}^t = \hat{\theta}_{i,t}$ . We term this approach BiCOMPFL-PR and summarize the procedure in Algorithm 2. Choosing different priors is possible and only affects line 11 in Algorithm 2. We mention in passing that BiCOMPFL-PR allows partial client participation, which is incompatible

with shared randomness and the method BICOMPFL-GR. We further note that our methods are readily compatible with sparsification and pruning techniques (e.g., Wangni et al. [2018], Shi et al. [2019]) by excluding pruned parameters from the block selection.

**Block Allocation.** We consider three different block allocation strategies: 1) fixed block size (referred to as “Fixed” in the experiments), where each block  $b \in [B]$  is of the same size and constant across all  $t$ ; 2) adaptive block allocation (Adaptive) as proposed by Isik et al. [2024], where each block size is separately optimized each iteration  $t$ ; and 3) adaptive average allocation (Adaptive-Avg), where the block sizes are equal but optimized at each iteration  $t$  according to the average KL-divergence per block. We refer the reader to Appendix D for a detailed discussion on this.

**Partial Client Participation.** We note that BICOMPFL-PR supports partial client participation without additional communication overhead, as it does not rely on full prior synchronization among the clients. BICOMPFL-GR on the other hand, requires that all clients’ priors are synchronized at each iteration. To allow partial client participation, the federator is required to transmit the previous global model to the clients that were absent in the previous iteration. This ensures full synchronization of the global model estimate used for efficient uplink communication. Hence, BICOMPFL-GR’s compatibility with partial client participation comes with a potentially increased downlink communication cost for previously passive clients.

## 4 Experiments

We next evaluate the performance of our proposed BICOMPFL-GR and BICOMPFL-PR schemes in experiments, and compare against baseline FL strategies without compression (FedAvg or PSGD) [McMahan et al., 2017] and several non-stochastic bi-directional compression schemes that employ different combinations of compression, error-feedback, and momentum. In particular, we compare against DOUBLESQUEEZE [Tang et al., 2019], MEM-SGD [Stich et al., 2018], NEOLITHIC [Huang et al., 2022], CSER [Xie et al., 2020], and the recently proposed LIEC [Cheng et al., 2024]. SignSGD [Seide et al., 2014] serves to compress the transmitted gradients for all the schemes. We further compare with M3 [Gruntkowska et al., 2024], which partitions the model into disjoint parts for downlink transmission and transmits to each client a different part of the model. While M3 is focused on RandK compression for the uplink (i.e., transmitting random  $K$  entries of the gradient), we use TopK [Wangni et al., 2018, Shi et al., 2019], which achieved much more stable results. As mentioned above, the mirror descent approach outlined in Section 3 inherently minimizes the communication cost as a by-product. Hence, it is a strong candidate for communication-efficient stochastic FL. Nonetheless, we show how our method can also be used to improve the communication efficiency in conventional FL, by using the uplink and downlink compression  $\mathcal{C}_{\text{mrc}}(\cdot)$  combined with stochastic quantizers, e.g., [Alistarh et al., 2017]. In Section 5, we pave the way to convergence guarantees by proving a contraction property of  $\mathcal{C}_{\text{mrc}}(\cdot)$  composed with a stochastic quantization  $Q_s(\cdot)$  of gradients  $g_i^t$ . To compare our method to the baselines that use SignSGD as compressor, we evaluate BICOMPFL-GR in a conventional federated learning (CFL) task with a stochastic variant of SignSGD. We replace mirror descent over Bernoulli masks by a standard learning procedure over a deterministic model, which takes as input the global model estimate  $\hat{\theta}_{i,t}$ , computes a gradient  $g_i^t$  (over  $L$  local epochs, using SGD and cross-entropy losses), and outputs a distribution  $Q_s(g_i^t)$ . In stochastic SignSGD,  $Q_s(\cdot)$  transforms each gradient entry  $g_{i,e}^t$  to a Bernoulli random variable with parameter  $q_{i,e}^t = 1/(1 + \exp(-g_{i,e}^t/K))$  for some  $K > 0$ , where the random variable takes value  $+1$  with probability  $q_{i,e}^t$ , and  $-1$  otherwise. We then employ  $\mathcal{C}_{\text{mrc}}(q_{i,e}^t, p_{i,u}^t)$  to obtain samples  $y_{i,\ell}^t$ , where the compression is performed element-wise. We apply this method to BICOMPFL-GR where Step 6 is replaced by  $\theta_{t+1} = \theta_t - \eta_s \frac{1}{n} \sum_{i=1}^n \hat{q}_i^t$ , where  $\hat{q}_i^t = \frac{1}{n_{\text{UL}}} \sum_{\ell=1}^{n_{\text{UL}}} y_{i,\ell}^t$  and  $\eta_s$  is the federator’s learning rate. Step 9 is modified accordingly. The priors  $p^t$  are chosen as Bernoulli random variables with parameter 0.5. We re-

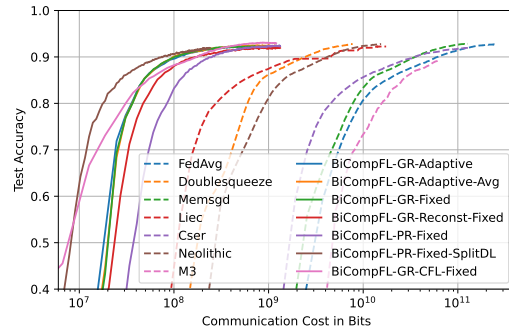


Figure 1: Test accuracy for BICOMPFL and baselines on Fashion MNIST 4CNN on i.i.d. data.

mark that while MRC samples are biased towards  $p^t$  (as we discuss in Section 5), this particular prior choice avoids imbalance in stochastic SignSGD, and rather acts as a regularizer, pulling the clients' posteriors closer to maximum entropy distributions. Consequently, convergence is achieved under bi-directional compression even without error feedback. For general prior choices, error feedback may be needed, see Algorithm 3 and Appendix B. We will refer to this method as BiCOMPFL-GR-CFL.

We study  $n = 10$  clients (see Appendix H for additional experiments with more clients) collaboratively training a convolutional neural network (CNN)-based classifier for the datasets MNIST, Fashion-MNIST and CIFAR-10 under the orchestration of a federator. For MNIST, we use two different models, LeNet-5 [Lecun et al., 1998] and a 4-layer convolutional neural network (4CNN) proposed by Ramanujan et al. [2020]. The latter is also used to train on Fashion MNIST. For CIFAR-10, we use a larger neural network with 6 convolutional layers (6CNN). We train MNIST and Fashion-MNIST for 200 global iterations and CIFAR-10 for 400 global iterations. Through all experiments and datasets, we carry  $L = 3$  local iterations per client. The learning rates are carefully selected to ensure convergence and comparability across all methods. Particularly, we tune the hyperparameters so that all algorithms achieve similar accuracies, allowing a fair comparison of their communication costs (see Appendix I.6 for details). Our main claims are the communication reduction of the bitrates per parameter per epoch, which are orthogonal to the choice of the learning rates of the algorithms. The code to reproduce our experiments is included in the supplementary material.

We evaluate the schemes in two settings: with a uniform data allocation (i.i.d.), to model homogeneous systems, and with a non-i.i.d. allocation, to model heterogeneous systems, where data allocation for each client is drawn from a Dirichlet distribution with parameter  $\alpha = 0.1$ . This regime is challenging due to extreme class imbalance. Each result shows the average across three simulation runs with different seeds. Further details on the simulation setup and the network architectures are deferred to Appendix E. Consistently throughout all experiments, our proposed methods provide **order-wise improvements** in the communication cost, while achieving state-of-the-art accuracies.

We plot in Fig. 1 the test accuracies for all the schemes as a function of the total communication cost in bits per parameter and per global iteration. While all the schemes achieve approximately the same maximum test accuracy, BiCOMPFL-GR and BiCOMPFL-PR require substantially less communication. Hence, when the bandwidths of uplink and downlink transmissions are limited, both variations of the proposed method achieve better test accuracies. Turning our focus to the different variations of our scheme, it can be observed that, without partitioning the model for downlink compression, BiCOMPFL-PR converges significantly slower than BiCOMPFL-GR for any block allocation method. This highlights the intuition above that the additional MRC step in downlink incurs further noise, which reduces the convergence speed. However, when we partition the model in the downlink and only send disjoint parts to each client through MRC (BiCOMPFL-PR-Fixed-SplitDL), the downlink communication cost reduces by a factor of  $n$ . In the regime of Fashion MNIST with a uniform data allocation, this comes without performance degradation, and is hence the method of choice in this regime. We additionally simulated BiCOMPFL-GR with the suboptimal implementation (BiCOMPFL-GR-Reconst-Fixed), in which the federator first reconstructs the global model, and then performs an additional MRC step for downlink transmission. This naturally reduces the convergence speed per iteration without gains in the communication cost. Hence, justifying the choice of BiCOMPFL-GR. We show that, in conventional FL, BiCOMPFL-GR-CFL substantially reduces the communication cost without loss in performance. In some cases, especially for non-i.i.d. data, we even observe improved performance, which we attribute to implicit regularization. Note that BiCOMPFL-GR-CFL provides improvements even without error-feedback or momentum. However, our method is fully compatible with such techniques and can be used as a plug-in approach to further minimize the communication cost in many existing schemes. We study the convergence in Section 5.

We plot in Fig. 2(a) the schemes' average bitrates over the maximum test accuracy for MNIST and 4CNN. The average bitrate is reduced by more than a factor of 1000 compared to FedAvg, and more than a **factor of 32** compared to DOUBLESQUEEZE, NEOLITHIC and LIEC, which perform best among the conventional bi-directional compression methods. We repeat the study for non-i.i.d. data allocation according to a Dirichlet distribution with parameter  $\alpha = 0.1$ , and show maximum test accuracies over average bitrates in Fig. 2(b). Partitioning the model in BiCOMPFL-PR worsens the final accuracy of the model. While the model converges faster, it does not achieve the same accuracies as BiCOMPFL-GR and BiCOMPFL-PR without partitioning. This hints towards hybrid schemes for BiCOMPFL-PR, where the training begins with partitioning on the downlink, which is later switched

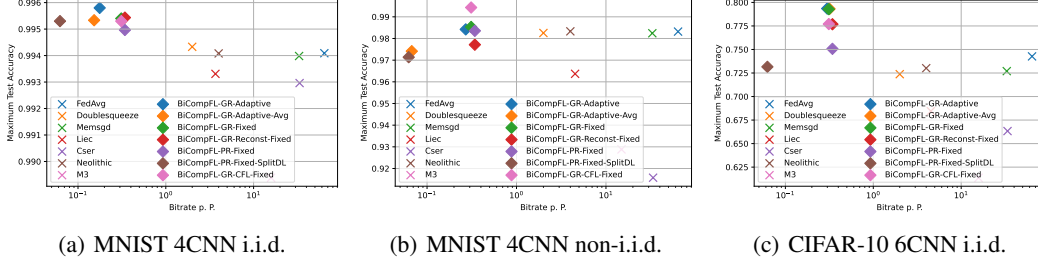


Figure 2: Maximum test accuracy over total communication cost measured by bitrate per parameter.

to full transmission. In Fig. 2(c), we provide the results for CIFAR-10 and uniform data allocation. BiCOMPFL-GR and BiCOMPFL-PR both achieve better results with a bitrate **smaller by a factor of 5** than the best baselines. More detailed numerical results can be found in Appendices H and I.

The adaptive block allocation (Adaptive) of Isik et al. [2024] saves communication costs in many settings and provides better performance than the fixed block allocation (Fixed), due to more accurate MRC tailored to the exact divergences. The proposed low complexity adaptive strategy based on the average KL-divergence (Adaptive-Avg) per block can additionally save in communication (and computation) with no or little performance degradation. We refer the reader to Appendix H for further extensive experiments, graphs for accuracies over epochs, separate studies of uplink and downlink costs, and comparisons for the case of an available broadcast channel from federator to the clients. Finally, we refer to Appendix I for various ablation studies analyzing the sensitivity of BiCOMPFL with respect to the choices of the priors,  $n$ ,  $n_{DL}$ ,  $n_{IS}$ , the block size  $d/B$ , and the learning rate  $\eta$ .

## 5 Theoretical Results

**Convergence.** In stochastic FL, the exact time dynamics of the system are challenging to analyze due to the round-dependent interplay of the learning procedure with the transmission noise. However, when using BiCOMPFL for conventional FL with stochastic quantization (cf. BiCOMPFL-GR-CFL), convergence guarantees can be given. We prove the convergence for a general and widely used class of stochastic quantizers  $Q_s(\cdot)$ , which are natively unbiased.  $Q_s(\cdot)$  takes as input the entry  $g_e$  of a gradient vector  $\mathbf{g} \in \mathbb{R}^d$  and operates as follows. Let  $s$  be the number of quantization intervals, and let  $0 \leq \tau_e < s$  be an integer such that  $\frac{\tau_e}{s} \leq \frac{|g_e|}{\|\mathbf{g}\|} \leq \frac{\tau_e+1}{s}$ , then  $Q_s(g_e)$  outputs  $\|\mathbf{g}\| \cdot \text{sign}(g_e)(\tau_e + 1)/s$  with probability  $s|g_e|/\|\mathbf{g}\| - \tau_e$ , and  $\|\mathbf{g}\| \cdot \text{sign}(g_e)\tau_e/s$  otherwise.  $Q_s(\cdot)$  is unbiased, i.e.,  $\mathbb{E}[Q_s(\mathbf{x})] = \mathbf{x}$ , and its variance satisfies  $\mathbb{E}[\|Q_s(\mathbf{x}) - \mathbf{x}\|^2] \leq \min\{d/s^2, \sqrt{d}/s\}\|\mathbf{x}\|_2^2$  [Alistarh et al., 2017]. See Remark 2 in Appendix B for a discussion on the choice of  $s$ .

Replacing stochastic SignSGD by  $Q_s(\cdot)$  in BiCOMPFL-GR-CFL, the posterior is given by a Bernoulli distribution with parameter  $q_{i,e}^t = s|g_{i,e}^t|/\|g_i^t\| - \tau_e$ . The values  $\|\mathbf{g}\|$ ,  $\text{sign}(\mathbf{g})$ , and  $\tau_e$  can be encoded independently, e.g., using Elias coding. With a slight abuse of notation, let  $\mathcal{C}_{\text{mrc}}(Q_s(\cdot), \cdot)$  denote the composition of  $Q_s(\cdot)$  and MRC with  $n_{IS}$  samples per entry. The compression  $\mathcal{C}_{\text{mrc}}(Q_s(g_i^t), \cdot)$  takes a gradient  $g_i^t$  and outputs samples from a distribution close to  $Q_s(g_i^t)$ , and falls in the class of biased compressors. We can prove the following contraction property for  $\mathcal{C}_{\text{mrc}}(Q_s(\cdot), \cdot)$ , which will facilitate convergence analysis for uni- and bi-directional compression. This constitutes a substantial improvement over [Isik et al., 2024], where such guarantees were missing, and hence no convergence guarantees were given. A prominent biased contractive compressor is TopK.

**Lemma 1.** *For any  $\mathbf{x} \in \mathbb{R}^d$  and corresponding posterior  $q$  following  $Q_s(\mathbf{x})$ , and a prior  $p \in [0, 1]^d$ , let  $\bar{\Delta} := \max_{e \in [d]} \frac{q_e}{p_e} - \frac{1-q_e}{1-p_e}$ ,  $\bar{\Delta}' := \max_{e \in [d]} q_e \left( \frac{p_e}{q_e} + \frac{1-p_e}{1-q_e} \right)$ , and  $\bar{p} := \max_{e \in [d]} p_e$ . The compressor  $\mathcal{C}_{\text{mrc}}(Q_s(\cdot))$  satisfies the following contraction property for  $n_{IS} = \mathcal{O}(\max\{\sqrt{2\bar{\Delta}}, \log(6\bar{p}(\bar{\Delta} + \bar{\Delta}^2))\sqrt{6\bar{p}(\bar{\Delta} + \bar{\Delta}^2)}\})$  and  $s \geq \sqrt{2d}$ :*

$$\mathbb{E}[\|\mathcal{C}_{\text{mrc}}(Q_s(\mathbf{x})) - \mathbf{x}\|^2] \leq (1 - \delta)\|\mathbf{x}\|^2,$$

$$\text{for } \delta = 1 - \frac{d}{s^2} \left( 1 + \frac{\bar{\Delta}'}{n_{IS}^2} + \mathcal{O} \left( (\bar{\Delta} + \bar{\Delta}^2) \sqrt{\frac{6\bar{p} \log(2n_{IS})}{n_{IS}}} \right) \right).$$



The underlying core result is a refinement of the MRC analysis, cf. Lemma 2 (Appendix A). Hence, for sufficiently large  $n_{\text{IS}}$ , the compressor  $\mathcal{C}_{\text{mrc}}(Q_s(\cdot), \cdot)$  can be used as an alternative to common compressors such as  $Q_s(\cdot)$ . The use of MRC introduces a bias into the otherwise unbiased stochastic quantization. Based on the contraction property in Lemma 1, standard convergence results (cf. Theorem 2) follow easily by a straightforward extension of our conventional FL algorithm BiCOMPFL-GR-CFL to error feedback (cf. Algorithm 3) as detailed in Appendix B.

**Communication Cost.** We analyze the communication cost in a specific iteration  $t$  and comment on the inter-round dependency later. When the latest global model estimate  $\hat{\theta}_{i,t}$  is chosen as a prior in MRC, the uplink cost is determined by how far the model evolves during the client's training, i.e.,  $d_{\text{KL}}(q_i^t || p_{i,u}^t) = d_{\text{KL}}(q_i^t || \hat{\theta}_{i,t})$ . After communicating samples of the posteriors, the federator obtains an estimate  $\hat{q}_i^t$  for all  $i \in [n]$ . The cost of communication on the downlink to client  $i$  is then determined by  $d_{\text{KL}}(\frac{1}{n} \sum_{i=1}^n \hat{q}_i^t || \hat{\theta}_{i,t})$ . While  $d_{\text{KL}}(q_i^t || \hat{\theta}_{i,t})$  depends on the progress during client training, the core challenge is to bound the expected KL-divergence of each model estimate  $d_{\text{KL}}(\hat{q}_i^t || \hat{\theta}_{i,t})$  in the presence of potentially different priors, i.e.,  $\hat{\theta}_{i,t} \neq \hat{\theta}_{j,t}, i \neq j$ . For each client  $i$ , the overall communication cost is in the order of  $n_{\text{DL}} \exp(d_{\text{KL}}(\frac{1}{n} \sum_{i=1}^n \hat{q}_i^t || p_{i,d}^t)) + n_{\text{UL}} \exp(d_{\text{KL}}(q_i^t || p_{i,u}^t))$ . We will next quantify  $d_{\text{KL}}(\frac{1}{n} \sum_{i=1}^n \hat{q}_i^t || \hat{\theta}_{i,t})$  for the case  $p_{i,u}^t = p_{i,d}^t$ , however, the analysis can be extended to  $p_{i,u}^t \neq p_{i,d}^t$  by an additional assumption on the divergence between the two priors.

For the theoretical analysis, we focus on the scalar case for a single iteration  $t$ , where client  $i \in [n]$  has a posterior  $Q_i$ , and the federator and client  $i$  share a common prior  $P_i$ , both are Bernoulli distributions with parameters  $q_i$  and  $p_i$ , respectively. In the context of FL, the client locally trains  $P_i$  and results with  $Q_i$ . According to Chatterjee and Diaconis [2018] and the multi-client extension of Isik et al. [2024], the communication cost in the uplink is determined by  $\exp(d_{\text{KL}}(Q_i || P_i))$ . After uplink transmission, the federator obtains an estimate  $\hat{q}_i$  of  $q_i$ ; and hence, the updated global model is given by  $\frac{1}{n} \sum_{i=1}^n \hat{q}_i$ . The downlink cost for client  $i$  is determined by  $d_{\text{KL}}(\frac{1}{n} \sum_{i=1}^n \hat{q}_i || p_i)$ .

We derive a new high probability upper bound on this quantity, refining previous MRC analysis for the special case of Bernoulli distributions. While a more general analysis can be conducted for other classes of distributions (cf. Remark 1 in Appendix A), the Bernoulli-based optimization method described earlier proves particularly efficient by enabling communication-regularized training, a unique property that renders our methods substantially more communication-efficient while preserving state-of-the-art performance.

Let  $X$  be a Bernoulli sample obtained through MRC. As an initial step, we bound the difference between  $q_i$  and the probability  $\Pr(X = 1)$  that the samples are drawn from, which vanishes when  $p_i = q_i$  (and hence  $d_{\text{KL}}(q_i || p_i) = 0$ ). We note that the bound of Chatterjee and Diaconis [2018, Theorem 1.1] does not satisfy this natural property. We formally state the result in Proposition 1 in Appendix A, which, however, does not yet capture the dependency on the number of samples  $n_{\text{IS}}$  used in MRC to sample an index. We refine Proposition 1 with Lemma 2 (cf. Appendix A), which additionally captures this dependency, and will allow us to derive an upper bound on  $d_{\text{KL}}(\frac{1}{n} \sum_{i=1}^n \hat{q}_i || p_i)$ . Lemma 2 is of independent interest and can be seen as a refinement of [Chatterjee and Diaconis, 2018] for Bernoulli distributions. It is required to prove Theorem 1.

For the statement of the following theorem, we assume that the progress by one local client training is bounded by  $|q_j - p_j| \leq \rho$  for all  $j \in [n]$ . Using Pinsker's inequality to bound  $|q_j - p_j| \leq \frac{1}{2} \sqrt{d_{\text{KL}}(q_j || p_j)}/2$ , this is a natural assumption given from the KL-proximity term of mirror descent (for one local iteration), and can be strictly enforced through the projection of  $q_j$  onto a KL ball around  $p_j$  of fixed divergence. We assume that the difference between the clients' priors, i.e., their global model estimates in our algorithms, are bounded as  $|p_i - p_j| \leq \zeta$  for all  $i, j \in [n]$ .

**Theorem 1.** Assume  $p_j > \zeta$  for all  $j \in [n]$ , for  $\Delta_j := \frac{q_j}{p_j - \zeta} - \frac{1 - q_j}{1 - p_j + \zeta}$  and  $\Delta'_j := q_j \left( \frac{p_j + \zeta}{q_j} + \frac{1 - p_j + \zeta}{1 - q_j} \right)$ , with probability  $1 - \delta'$ , the global model divergence  $d_{\text{KL}}(\frac{1}{n} \sum_{j=1}^n \hat{q}_j || p_i)$  is upper bounded by

$$\sum_{j=1}^n \frac{2}{n \min\{p_i, 1 - p_i\}} \left( \frac{\Delta'_j}{n_{\text{IS}}^2} + \sqrt{\frac{\ln(2/\delta')}{2n_{\text{UL}}}} + \rho + \zeta^2 + \mathcal{O}\left((\Delta_j + \Delta_j^2) \sqrt{\frac{6(p_i + \zeta) \log(2n_{\text{IS}})}{n_{\text{IS}}}}\right) \right).$$

By Chatterjee and Diaconis [2018], this provides an immediate bound on the cost of downlink transmission. The bound applies to both algorithms BiCOMPFL-PR and BiCOMPFL-GR. However, when all priors  $p_j$  are the same (such as in BiCOMPFL-GR-Reconst), i.e.,  $\zeta = 0$ , the bound simplifies

accordingly. The explicit dependency on the factor  $1/\sqrt{n_{UL}}$  reflects the interplay between uplink and downlink cost. The parameter  $\zeta$  gives rise to an inter-round dependency of the communication cost. The more accurate the estimation of the global model in the previous iteration (given the priors are chosen as  $\hat{\theta}_{i,t}$ ), the smaller  $\zeta$ , and hence the lower the transmission cost in the subsequent iteration. The proofs of Proposition 1, Lemma 2, and Theorem 1 can be found in Appendix A.

## 6 Related Work

Following the introduction of FL [McMahan et al., 2017], lossy compression of gradients or model updates has been a long-studied narrative in FL, with prominent representatives such as SignSGD, also known as 1-bit Stochastic Gradient Descent (SGD) [Seide et al., 2014], QSGD [Alistarh et al., 2017], TernGrad [Wen et al., 2017], SignSGD with error feedback [Karimireddy et al., 2019], vector-quantized SGD [Gandikota et al., 2021] and natural compression [Horvóth et al., 2022]. Such methods retain satisfactory final model accuracy even with aggressive quantization. Sparsification-based methods have also been considered as alternatives, e.g., TopK [Wangni et al., 2018, Shi et al., 2019]. The importance of bi-directional gradient compression in many settings was outlined by Philippenko and Dieuleveut [2020]. Many schemes were proposed that leverage combinations of gradient compression in the uplink and downlink, error-feedback, and momentum, e.g., Mem-SGD [Stich et al., 2018], DoubleSqueeze [Tang et al., 2019], block-wise SignSGD with momentum [Zheng et al., 2019], communication-efficient SGD with error reset (Cser) [Xie et al., 2020], Artemis [Philippenko and Dieuleveut, 2020], Neolithic [Huang et al., 2022], DoCoFL [Dorfman et al., 2023], EF21-P and friends [Gruntowska et al., 2023], 2Direction [Tyurin and Richtárik, 2023], M3 [Gruntowska et al., 2024], and LIEC [Cheng et al., 2024]. With the exception of the methods MCM [Philippenko and Dieuleveut, 2021] and M3 [Gruntowska et al., 2024], each client receives the same broadcast, potentially compressed, global gradient or model update. Isik et al. [2024] studied uplink compression for stochastic FL and showed significant communication reduction with competitive performance. Their framework, termed KLMS, applies to a variety of stochastic compressors and to Bayesian FL settings, e.g., QLSG Vono et al. [2022]. The compression is based on importance sampling and MRC, thoroughly studied by Chatterjee and Diaconis [2018] and Havasi et al. [2019]. Such methods, known as relative entropy coding, have been used in FL in conjunction with differential privacy, cf. DP-REC [Triastcyn et al., 2022].

Since the lottery ticket hypothesis [Frankle and Carbin, 2019], finding sparse subnetworks of neural networks that achieve satisfactory accuracy was investigated. Ramanujan et al. [2020] showed that randomly weighted networks contain suitable subnetworks of large neural networks capable of achieving competitive performance. Isik et al. [2023] formulated a probabilistic method of training neural network masks collaboratively in an FL context.

## 7 Conclusion

We illuminated bi-directional compression in stochastic FL via federated probabilistic mask training, which we showed to inherently optimize both the learning objective and the communication costs. By leveraging side information through carefully chosen prior distributions, the total communication costs are reduced by factors between 5 – 32 compared to non-stochastic FL baselines, while achieving state-of-the-art accuracies on classification tasks, for both homogeneous and heterogeneous data. We thus close the gap of downlink compression for stochastic FL and complement the existing literature on bi-directional compression for standard FL. Applying our methods to stochastic quantization in conventional FL, we paved the way to convergence analysis for MRC-based compression. Allowing different priors among all clients, this work opens the door to studying compression under side-information in *decentralized stochastic FL*, where a central coordinator is missing. Our theoretical results are of independent interest and may be applied in various scenarios where MRC is used.

Client privacy and fairness are important directions beyond the primary focus of this work. Compression methods commonly strengthen the clients’ privacy in FL. Particularly, the noise introduced through our sampling methods dilute the individual contributions of the clients observed by the federator, naturally enhancing the clients’ privacy and reducing the risk of privacy breaches through, e.g., membership inference attacks. Quantifying the privacy-utility trade-off that arises in the presence of bi-directional compression is left as an interesting direction for future work.

## Acknowledgments

The research of R. Bitar is supported by the German Research Foundation (DFG) under grant agreement number BI 2492/1-1. The research of N. Weinberger was partially supported by the Israel Science Foundation (ISF), grant no. 1782/22. D. Gündüz acknowledges funding from UKRI under the ERC-Consolidator project AI-R (EP/X030806/1) and INFORMED-AI Hub (EP/Y028732/1).

## References

- A. Abdi and F. Fekri. Nested dithered quantization for communication reduction in distributed training. *arXiv preprint arXiv:1904.01197*, 2019.
- D. Alistarh, D. Grubic, J. Li, R. Tomioka, and M. Vojnovic. QSGD: Communication-efficient SGD via gradient quantization and encoding. In *Advances in Neural Information Processing Systems*, volume 30, 2017.
- S. Chatterjee and P. Diaconis. The sample size required in importance sampling. *The Annals of Applied Probability*, 28(2):1099–1135, 2018.
- Y. Cheng, L. Shen, L. Xu, X. Qian, S. Wu, Y. Zhou, T. Zhang, D. Tao, and E. Chen. Communication-efficient distributed learning with local immediate error compensation. *arXiv preprint arXiv:2402.11857*, 2024.
- R. Dorfman, S. Vargaftik, Y. Ben-Itzhak, and K. Y. Levy. DoCoFL: Downlink compression for cross-device federated learning. In *International Conference on Machine Learning*, pages 8356–8388, 2023.
- J. Frankle and M. Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural networks. In *International Conference on Learning Representations*, 2019.
- V. Gandikota, D. Kane, R. Kumar Maity, and A. Mazumdar. vqSGD: Vector quantized stochastic gradient descent. In *International Conference on Artificial Intelligence and Statistics*, volume 130, pages 2197–2205, 2021.
- K. Gruntkowska, A. Tyurin, and P. Richtárik. EF21-P and friends: Improved theoretical communication complexity for distributed optimization with bidirectional compression. In *International Conference on Machine Learning*, pages 11761–11807, 2023.
- K. Gruntkowska, A. Tyurin, and P. Richtárik. Improving the worst-case bidirectional communication complexity for nonconvex distributed optimization under function similarity. *arXiv preprint arXiv:2402.06412*, 2024.
- M. Havasi, R. Peharz, and J. M. Hernández-Lobato. Minimal random code learning: Getting bits back from compressed model parameters. In *International Conference on Learning Representations*, 2019.
- S. Horvóth, C.-Y. Ho, L. Horvath, A. N. Sahu, M. Canini, and P. Richtarik. Natural compression for distributed deep learning. In *Proceedings of Mathematical and Scientific Machine Learning*, volume 190, pages 129–141, 2022.
- X. Huang, Y. Chen, W. Yin, and K. Yuan. Lower bounds and nearly optimal algorithms in distributed learning with communication compression. *Advances in Neural Information Processing Systems*, 35:18955–18969, 2022.
- B. Isik, F. Pase, D. Gunduz, T. Weissman, and Z. Michele. Sparse random networks for communication-efficient federated learning. In *International Conference on Learning Representations*, 2023.
- B. Isik, F. Pase, D. Gunduz, S. Koyejo, T. Weissman, and M. Zorzi. Adaptive compression in federated learning via side information. In *International Conference on Artificial Intelligence and Statistics*, pages 487–495, 2024.

- P. Kairouz and H. B. McMahan. Advances and open problems in federated learning. *Foundations and trends in machine learning*, 14(1-2):1–210, 2021.
- S. P. Karimireddy, Q. Rebjock, S. Stich, and M. Jaggi. Error feedback fixes SignSGD and other gradient compression schemes. In *International Conference on Machine Learning*, volume 97, pages 3252–3261, 2019.
- D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations*, 2015.
- Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- C. T. Li. Channel simulation: Theory and applications to lossy compression and differential privacy. *Foundations and Trends® in Communications and Information Theory*, 21(6):847–1106, 2024.
- B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y. Arcas. Communication-Efficient Learning of Deep Networks from Decentralized Data. In *International Conference on Artificial Intelligence and Statistics*, volume 54, pages 1273–1282, 2017.
- U. Milasheuski, L. Barbieri, S. Kianoush, M. Nicoli, and S. Savazzi. Bayesian federated learning for continual training. *arXiv preprint arXiv:2504.15328*, 2025.
- C. Philippenko and A. Dieuleveut. Bidirectional compression in heterogeneous settings for distributed or federated learning with partial participation: tight convergence guarantees. *arXiv preprint arXiv:2006.14591*, 2020.
- C. Philippenko and A. Dieuleveut. Preserved central model for faster bidirectional compression in distributed settings. *Advances in Neural Information Processing Systems*, 34:2387–2399, 2021.
- V. Ramanujan, M. Wortsman, A. Kembhavi, A. Farhadi, and M. Rastegari. What’s hidden in a randomly weighted neural network? In *IEEE/CVF conference on computer vision and pattern recognition*, pages 11893–11902, 2020.
- P. Richtárik, I. Sokolov, and I. Fatkhullin. Ef21: A new, simpler, theoretically better, and practically faster error feedback. *Advances in Neural Information Processing Systems*, 34:4384–4396, 2021.
- F. Seide, H. Fu, J. Droppo, G. Li, and D. Yu. 1-bit stochastic gradient descent and its application to data-parallel distributed training of speech DNNs. In *Interspeech*, pages 1058–1062, 2014.
- S. Shi, Q. Wang, K. Zhao, Z. Tang, Y. Wang, X. Huang, and X. Chu. A distributed synchronous SGD algorithm with global top-k sparsification for low bandwidth networks. In *International Conference on Distributed Computing Systems (ICDCS)*, pages 2238–2247, 2019.
- R. Srinivasan. *Importance sampling: Applications in communications and detection*. Springer Science & Business Media, 2002.
- S. U. Stich, J.-B. Cordonnier, and M. Jaggi. Sparsified SGD with memory. *Advances in Neural Information Processing Systems*, 31, 2018.
- H. Tang, C. Yu, X. Lian, T. Zhang, and J. Liu. Doublesqueeze: Parallel stochastic gradient descent with double-pass error-compensated compression. In *International Conference on Machine Learning*, pages 6155–6165, 2019.
- L. Theis and N. Y. Ahmed. Algorithms for the communication of samples. In *International Conference on Machine Learning*, pages 21308–21328, 2022.
- A. Triastcyn, M. Reisser, and C. Louizos. DP-REC: Private & communication-efficient federated learning, 2022.
- A. Tyurin and P. Richtárik. 2Direction: theoretically faster distributed training with bidirectional communication compression. In *Conference on Neural Information Processing Systems*, 2023.

- M. Vono, V. Plassier, A. Durmus, A. Dieuleveut, and E. Moulines. QLSD: quantised langevin stochastic dynamics for bayesian federated learning. In *International Conference on Artificial Intelligence and Statistics*, pages 6459–6500, 2022.
- J. Wangni, J. Wang, J. Liu, and T. Zhang. Gradient sparsification for communication-efficient distributed optimization. *Advances in Neural Information Processing Systems*, 31, 2018.
- N. Weinberger and M. Yemini. Multi-armed bandits with self-information rewards. *IEEE Transactions on Information Theory*, 2023.
- J. Wen, Z. Zhang, Y. Lan, Z. Cui, J. Cai, and W. Zhang. A survey on federated learning: challenges and applications. *International Journal of Machine Learning and Cybernetics*, 14(2):513–535, 2023.
- W. Wen, C. Xu, F. Yan, C. Wu, Y. Wang, Y. Chen, and H. Li. TernGrad: Ternary gradients to reduce communication in distributed deep learning. In *Advances in Neural Information Processing Systems*, volume 30, 2017.
- C. Xie, S. Zheng, S. Koyejo, I. Gupta, M. Li, and H. Lin. Cser: Communication-efficient sgd with error reset. *Advances in Neural Information Processing Systems*, 33:12593–12603, 2020.
- C. Zhang, Y. Xie, H. Bai, B. Yu, W. Li, and Y. Gao. A survey on federated learning. *Knowledge-Based Systems*, 216:106775, 2021.
- X. Zhang, Y. Li, W. Li, K. Guo, and Y. Shao. Personalized federated learning via variational bayesian inference. In *International Conference on Machine Learning*, pages 26293–26310, 2022.
- S. Zheng, Z. Huang, and J. Kwok. Communication-efficient distributed blockwise momentum SGD with error-feedback. *Advances in Neural Information Processing Systems*, 32, 2019.

## A Proofs and Intermediate Results

In the following, we provide the formal statements of Proposition 1 and Lemma 2 including their proofs. Parts of the proof of Proposition 1 will be used to prove Lemma 2. We prove Theorem 1 afterward.

**Proposition 1.** *For a sample  $X_\ell$  transmitted by MRC with posterior and prior Bernoulli distributions with parameters  $q$  and  $p$ , we have*

$$|\Pr(X_\ell = 1) - q| \leq q \left( \max \left\{ \frac{p}{q}, \frac{1-p}{1-q}, \frac{q}{p}, \frac{1-q}{1-p} \right\} - 1 \right).$$

*Proof of Proposition 1.* Assume a party wants to sample from a Bernoulli distribution  $Q$  with parameter  $q$ , which is held by another party. Both parties share a common prior  $P$  in the form of a Bernoulli distribution with parameter  $p$  and have access to shared randomness. Fix any sample index  $\ell$  for the moment (this index will be needed for the proof of Theorem 1). Both parties sample  $Kn_{\text{IS}}$  i.i.d. samples  $X_{\ell,i} \sim P$  for  $i \in [n_{\text{IS}}]$  independently and identically from  $P$ . The party holding  $Q$  constructs an auxiliary distribution

$$W_\ell(i) = \frac{Q(X_{\ell,i})/P(X_{\ell,i})}{\sum_{i=1}^{n_{\text{IS}}} Q(X_{\ell,i})/P(X_{\ell,i})},$$

from which it samples to obtain an index  $I_\ell$ . The index is transmitted to the other party, which reconstructs the corresponding sample  $X_{\ell,I_\ell}$ .

To bound the difference  $|\Pr(X_\ell = 1) - q|$ , i.e., the target Bernoulli parameter compared to the parameter which the sample is drawn from, by the independence of the samples  $X_{\ell,I_\ell}$  for different  $\ell$ ,

we focus on a single sample  $\ell \in [K]$ , for which it holds that

$$\begin{aligned}
& \Pr(X_{\ell, I_\ell} = 1) \\
&= \sum_{i=1}^{n_{\text{IS}}} \sum_{\{x_1, \dots, x_{n_{\text{IS}}} : x_i = i\}} \Pr(X_{\ell, 1} = x_1, \dots, X_{\ell, n_{\text{IS}}} = x_{n_{\text{IS}}}) \Pr(I_\ell = i \mid X_{\ell, 1} = x_1, \dots, X_{\ell, n_{\text{IS}}} = x_{n_{\text{IS}}}) \\
&\stackrel{(a)}{=} n_{\text{IS}} \sum_{\{x_2, \dots, x_{n_{\text{IS}}}\}} \Pr(X_{\ell, 1} = 1, X_{\ell, 2} = x_2, \dots, X_{\ell, n_{\text{IS}}} = x_{n_{\text{IS}}}) \\
&\quad \cdot \Pr(I_\ell = 1 \mid X_{\ell, 1} = 1, X_{\ell, 2} = x_2, \dots, X_{\ell, n_{\text{IS}}} = x_{n_{\text{IS}}}) \\
&\stackrel{(b)}{=} n_{\text{IS}} \sum_{L=0}^{n_{\text{IS}}-1} \sum_{\{x_2, \dots, x_{n_{\text{IS}}} : \sum_{i=2}^{n_{\text{IS}}} x_i = L\}} \Pr(X_{\ell, 1} = 1, X_{\ell, 2} = x_2, \dots, X_{\ell, n_{\text{IS}}} = x_{n_{\text{IS}}}) \\
&\quad \cdot \Pr(I_\ell = 1 \mid X_{\ell, 1} = 1, X_{\ell, 2} = x_2, \dots, X_{\ell, n_{\text{IS}}} = x_{n_{\text{IS}}}),
\end{aligned}$$

where (a) follows from symmetry, (b) follows since by permutation invariance, the inner probability only depends on the number of ones in  $\{x_2, \dots, x_{n_{\text{IS}}}\}$ .

The inner probability is given by the distribution  $W_\ell(i)$ . Given that  $X_{\ell, 1} = 1$  and that  $\sum_{i=2}^{n_{\text{IS}}} X_{\ell, i} = L$ , it holds that

$$\sum_{i=1}^{n_{\text{IS}}} Q(X_{\ell, i})/P(X_{\ell, i}) = (L+1) \cdot \frac{q}{p} + (n_{\text{IS}} - L - 1) \cdot \frac{1-q}{1-p}.$$

Hence,

$$\Pr(I_\ell = 1 \mid X_{\ell, 1} = 1, X_{\ell, 2} = x_2, \dots, X_{\ell, n_{\text{IS}}} = x_{n_{\text{IS}}}) = \frac{\frac{q}{p}}{(L+1) \cdot \frac{q}{p} + (n_{\text{IS}} - L - 1) \cdot \frac{1-q}{1-p}},$$

which is independent of the exact choice of  $\{x_2, \dots, x_{n_{\text{IS}}}\}$  given their sum  $\sum_{i=2}^{n_{\text{IS}}} X_{\ell, i} = L$ . Since  $\Pr(X_{\ell, 1} = 1, X_{\ell, 2} = x_2, \dots, X_{\ell, n_{\text{IS}}} = x_{n_{\text{IS}}}) = p^{L+1}(1-p)^{n_{\text{IS}}-L-1}$  by the Bernoulli distribution assumption, we have

$$\Pr(X_{\ell, I_\ell} = 1) = n_{\text{IS}} \sum_{L=0}^{n_{\text{IS}}-1} \binom{n_{\text{IS}}-1}{L} p^{L+1}(1-p)^{n_{\text{IS}}-L-1} \frac{\frac{q}{p}}{(L+1) \cdot \frac{q}{p} + (n_{\text{IS}} - L - 1) \cdot \frac{1-q}{1-p}},$$

Defining a binary random variable  $M$  with sample space  $\left\{\frac{q}{p}, \frac{1-q}{1-p}\right\}$ , for a Bernoulli distribution  $\text{Ber}\left(\frac{L+1}{n_{\text{IS}}}\right)$  with success probability parameter  $\frac{L+1}{n_{\text{IS}}}$ , where a success refers to the outcome  $M = \frac{q}{p}$ , we can write that

$$\begin{aligned}
\Pr(X_{\ell, I_\ell} = 1) &= q \cdot \sum_{L=0}^{n_{\text{IS}}-1} \binom{n_{\text{IS}}-1}{L} p^L (1-p)^{n_{\text{IS}}-L-1} \frac{1}{\frac{L+1}{n_{\text{IS}}} \frac{q}{p} + \frac{n_{\text{IS}}-L-1}{n_{\text{IS}}} \frac{1-q}{1-p}} \\
&= q \cdot \mathbb{E} \left[ \frac{1}{\frac{L+1}{n_{\text{IS}}} \frac{q}{p} + \frac{n_{\text{IS}}-L-1}{n_{\text{IS}}} \frac{1-q}{1-p}} \right] = q \mathbb{E} \left[ \frac{1}{\mathbb{E}_{\text{Ber}\left(\frac{L+1}{n_{\text{IS}}}\right)}[M]} \right] \\
&\stackrel{(a)}{\leq} q \mathbb{E} \left[ \mathbb{E}_{\text{Ber}\left(\frac{L+1}{n_{\text{IS}}}\right)} \left[ \frac{1}{M} \right] \right],
\end{aligned} \tag{1}$$

where the outer expectation is over the binomial distribution with  $n_{\text{IS}} - 1$  trials and success probability  $p$ , i.e.,  $L \sim \text{Binomial}(n_{\text{IS}} - 1, p)$ , and where (a) follows from Jensen's inequality over the inner expectation. Hence,

$$\begin{aligned}
\Pr(X_{\ell, I_\ell} = 1) - q &= q \left( \frac{\Pr(X_{\ell, I_\ell} = 1)}{q} - 1 \right) \\
&\leq q \left( \mathbb{E} \left[ \mathbb{E}_{\text{Ber}\left(\frac{L+1}{n_{\text{IS}}}\right)} \left[ \frac{1}{M} \right] \right] - 1 \right)
\end{aligned} \tag{2}$$

Since  $\frac{1}{\mathbb{E}_{\text{Ber}\left(\frac{L+1}{n_{\text{IS}}}\right)}[M]} \geq 2 - \mathbb{E}_{\text{Ber}\left(\frac{L+1}{n_{\text{IS}}}\right)}[M]$ , it also follows from (1) that

$$\begin{aligned} \Pr(X_{\ell, I_{\ell}} = 1) &= q \cdot \mathbb{E} \left[ \frac{1}{\frac{L+1}{n_{\text{IS}}} \frac{q}{p} + \frac{n_{\text{IS}} - L - 1}{n_{\text{IS}}} \frac{1-q}{1-p}} \right] = q \mathbb{E} \left[ \frac{1}{\mathbb{E}_{\text{Ber}\left(\frac{L+1}{n_{\text{IS}}}\right)}[M]} \right] \\ &\geq q \mathbb{E} \left[ 2 - \mathbb{E}_{\text{Ber}\left(\frac{L+1}{n_{\text{IS}}}\right)}[M] \right], \end{aligned}$$

from which we have

$$\Pr(X_{\ell, I_{\ell}} = 1) - q \geq q \left( 1 - \mathbb{E} \left[ \mathbb{E}_{\text{Ber}\left(\frac{L+1}{n_{\text{IS}}}\right)}[M] \right] \right). \quad (3)$$

Combining the upper and lower bound in (2) and (3), respectively, we derive

$$\begin{aligned} |\Pr(X_{\ell, I_{\ell}} = 1) - q| &\leq q \left( \max \left\{ \mathbb{E} \left[ 1 - \mathbb{E}_{\text{Ber}\left(\frac{L+1}{n_{\text{IS}}}\right)}[M] \right], \mathbb{E} \left[ \mathbb{E}_{\text{Ber}\left(\frac{L+1}{n_{\text{IS}}}\right)} \left[ \frac{1}{M} \right] \right] \right\} - 1 \right) \\ &\leq q \left( \mathbb{E} \left[ \max \left\{ \mathbb{E}_{\text{Ber}\left(\frac{L+1}{n_{\text{IS}}}\right)}[M], \mathbb{E}_{\text{Ber}\left(\frac{L+1}{n_{\text{IS}}}\right)} \left[ \frac{1}{M} \right] \right\} \right] - 1 \right) \\ &\leq q \left( \mathbb{E} \left[ \mathbb{E}_{\text{Ber}\left(\frac{L+1}{n_{\text{IS}}}\right)} \left[ \max \left\{ M, \frac{1}{M} \right\} \right] \right] - 1 \right) \\ &\leq q \left( \mathbb{E} \left[ \max \left\{ \frac{p}{q}, \frac{1-p}{1-q}, \frac{q}{p}, \frac{1-q}{1-p} \right\} \right] - 1 \right) \\ &= q \left( \max \left\{ \frac{p}{q}, \frac{1-p}{1-q}, \frac{q}{p}, \frac{1-q}{1-p} \right\} - 1 \right). \end{aligned}$$

This concludes the proof.  $\square$

**Lemma 2.** For a sample  $X_{\ell}$  transmitted via MRC with posterior and prior being Bernoulli distributions with parameters  $q$  and  $p$ ,  $\Delta := \frac{q}{p} - \frac{1-q}{1-p}$  and  $\Delta' := q \left( \frac{p}{q} + \frac{1-p}{1-q} \right)$ , we have

$$|\Pr(X_{\ell} = 1) - q| \leq \frac{\Delta'}{n_{\text{IS}}^2} + \mathcal{O} \left( (\Delta + \Delta^2) \sqrt{\frac{6p \log(2n_{\text{IS}})}{n_{\text{IS}}}} \right).$$

*Proof of Lemma 2.* The proof starts with the same derivations as for the proof of Proposition 1, which we follow until (1) to get

$$\Pr(X_{\ell, I_{\ell}} = 1) = q \mathbb{E} \left[ \frac{1}{\mathbb{E}_{\text{Ber}\left(\frac{L+1}{n_{\text{IS}}}\right)}[M]} \right]$$

Since  $L$  is a random quantity that follows a Binomial distribution, we bound  $|\Pr(X_{\ell, I_{\ell}} = 1) - q|$  using a concentration bound on  $L$ . The relative (multiplicative) Chernoff bound states that

$$\begin{aligned} \Pr(|L - \varepsilon(n_{\text{IS}}p)| \geq \varepsilon n_{\text{IS}}p) &= \Pr(L - \varepsilon(n_{\text{IS}}p) \geq \varepsilon n_{\text{IS}}p) + \Pr(L - \varepsilon(n_{\text{IS}}p) \leq -\varepsilon n_{\text{IS}}p) \\ &\leq 2 \exp \left( -\frac{\varepsilon^2 n_{\text{IS}}p}{3} \right) \end{aligned}$$

for any  $\varepsilon \in [0, 1]$ . Setting  $\varepsilon = \sqrt{\frac{3 \log(2/\delta)}{n_{\text{IS}}p}}$  implies that

$$|L - n_{\text{IS}}p| \geq \sqrt{3n_{\text{IS}}p \log(2/\delta)}$$

with probability at most  $\delta$ . Setting  $\delta = \frac{1}{n_{\text{IS}}^2}$ , we obtain for a concentration parameter<sup>2</sup>  $\eta_{\delta} :=$

$\sqrt{\frac{6p \log(2n_{\text{IS}})}{n_{\text{IS}}}}$  that

$$\mathcal{E} := \{|L - n_{\text{IS}}p| \geq n_{\text{IS}}\eta_{\delta}\}$$

---

<sup>2</sup>Note that we can assume  $p + \eta_{\delta} \leq 1$  and  $p - \eta_{\delta} \geq 0$ , otherwise the concentration can be trivially bounded.

with probability  $\Pr(\mathcal{E}) \leq \frac{1}{n_{\text{IS}}^2}$ .

Then, we can write

$$\begin{aligned} \Pr(X_{\ell, I_\ell} = 1) &= q \mathbb{E} \left[ \frac{1}{\mathbb{E}_{\text{Ber}\left(\frac{L+1}{n_{\text{IS}}}\right)}[M]} \right] \\ &= q \mathbb{E} \left[ \frac{1}{\mathbb{E}_{\text{Ber}\left(\frac{L+1}{n_{\text{IS}}}\right)}[M]} \cdot \mathbb{1}\{\mathcal{E}^c\} \right] + q \mathbb{E} \left[ \frac{1}{\mathbb{E}_{\text{Ber}\left(\frac{L+1}{n_{\text{IS}}}\right)}[M]} \cdot \mathbb{1}\{\mathcal{E}\} \right] \end{aligned} \quad (4)$$

Assume for now that  $q < p$  (we will later proof the opposite event), then  $\frac{1}{\mathbb{E}_{\text{Ber}\left(\frac{L+1}{n_{\text{IS}}}\right)}[M]}$  is strictly non-increasing in  $L$  since  $\frac{q}{p} < \frac{1-q}{1-p}$ , and hence, when  $\mathcal{E}^c$  holds and hence  $L$  concentration around the average that

$$\begin{aligned} \frac{1}{\mathbb{E}_{\text{Ber}\left(\frac{L+1}{n_{\text{IS}}}\right)}[M]} &\leq \frac{1}{\mathbb{E}_{\text{Ber}\left(\frac{(L+1) \cdot (p-\eta_\delta)}{n_{\text{IS}}}\right)}[M]} \\ &= \frac{1}{\frac{(n_{\text{IS}}-1)(p-\eta_\delta)+1}{n_{\text{IS}}} \frac{q}{p} + \frac{n_{\text{IS}}-1-(n_{\text{IS}}-1)(p-\eta_\delta)}{n_{\text{IS}}} \frac{1-q}{1-p}} \\ &= \frac{1}{\left(p - \frac{p}{n_{\text{IS}}} + \frac{\eta_\delta}{n_{\text{IS}}} - \eta_\delta + \frac{1}{n_{\text{IS}}}\right) \frac{q}{p} + \left(1 - p - \frac{1}{n_{\text{IS}}} + \frac{p}{n_{\text{IS}}} + \eta_\delta - \frac{\eta_\delta}{n_{\text{IS}}}\right) \frac{1-q}{1-p}} \\ &= \frac{1}{1 + \left(\frac{q}{p} - \frac{1-q}{1-p}\right) \left(\frac{1-p+\eta_\delta-n\eta_\delta}{n_{\text{IS}}}\right)} \\ &= 1 + \sum_{\kappa=1}^{\infty} (-1)^\kappa \left(\frac{q}{p} - \frac{1-q}{1-p}\right)^\kappa \left(\frac{1-p+\eta_\delta-n\eta_\delta}{n_{\text{IS}}}\right)^\kappa, \end{aligned}$$

where the last step is by Taylor expansion. Using (4) and the monotonicity of  $\frac{1}{\mathbb{E}_{\text{Ber}\left(\frac{L+1}{n_{\text{IS}}}\right)}[M]}$ , we write

$$\begin{aligned} \Pr(X_{\ell, I_\ell} = 1) &= q \mathbb{E} \left[ \frac{1}{\mathbb{E}_{\text{Ber}\left(\frac{L+1}{n_{\text{IS}}}\right)}[M]} \right] \\ &\leq q \left( 1 + \sum_{\kappa=1}^{\infty} (-1)^\kappa \left(\frac{q}{p} - \frac{1-q}{1-p}\right)^\kappa \left(\frac{1-p+\eta_\delta-n\eta_\delta}{n_{\text{IS}}}\right)^\kappa \right) + q \delta \frac{p}{q}, \end{aligned}$$

and hence

$$\Pr(X_{\ell, I_\ell} = 1) - q \leq \delta p + (1-\delta) \sum_{\kappa=1}^{\infty} (-1)^\kappa \left(\frac{q}{p} - \frac{1-q}{1-p}\right)^\kappa \left(\frac{1-p+\eta_\delta-n\eta_\delta}{n_{\text{IS}}}\right)^\kappa$$

Similarly, we get by bounding  $\frac{1}{\mathbb{E}_{\text{Ber}\left(\frac{L+1}{n_{\text{IS}}}\right)}[M]} \geq \frac{1}{\mathbb{E}_{\text{Ber}\left(\frac{(L+1) \cdot (p+\eta_\delta)}{n_{\text{IS}}}\right)}[M]}$  and using (4) that

$$\begin{aligned} \Pr(X_{\ell, I_\ell} = 1) - q &\geq \delta q \frac{1-p}{1-q} + (1-\delta) \sum_{\kappa=1}^{\infty} (-1)^\kappa \left(\frac{q}{p} - \frac{1-q}{1-p}\right)^\kappa \left(\frac{1-p-\eta_\delta+n\eta_\delta}{n_{\text{IS}}}\right)^\kappa \Leftrightarrow \\ q - \Pr(X_{\ell, I_\ell} = 1) &\leq -\delta q \frac{1-p}{1-q} + (1-\delta) \sum_{\kappa=1}^{\infty} (-1)^{\kappa+1} \left(\frac{q}{p} - \frac{1-q}{1-p}\right)^\kappa \left(\frac{1-p-\eta_\delta+n\eta_\delta}{n_{\text{IS}}}\right)^\kappa. \end{aligned}$$

When  $p \leq q$ , then  $\frac{1}{\mathbb{E}_{\text{Ber}\left(\frac{L+1}{n_{\text{IS}}}\right)}[M]}$  is strictly non-decreasing, hence, under  $\mathcal{E}$ , we have

$$\frac{1}{\mathbb{E}_{\text{Ber}\left(\frac{L+1}{n_{\text{IS}}}\right)}[M]} \leq \frac{1}{\mathbb{E}_{\text{Ber}\left(\frac{(L+1) \cdot (p+\eta_\delta)}{n_{\text{IS}}}\right)}[M]} = 1 + \sum_{\kappa=1}^{\infty} (-1)^\kappa \left(\frac{q}{p} - \frac{1-q}{1-p}\right)^\kappa \left(\frac{1-p-\eta_\delta+n\eta_\delta}{n_{\text{IS}}}\right)^\kappa,$$



and thus from (4) that

$$\Pr(X_{\ell, I_\ell} = 1) - q \leq q\delta \frac{1-p}{1-q} + (1-\delta) \sum_{\kappa=1}^{\infty} (-1)^\kappa \left( \frac{q}{p} - \frac{1-q}{1-p} \right)^\kappa \left( \frac{1-p-\eta_\delta+n\eta_\delta}{n_{\text{IS}}} \right)^\kappa.$$

Similarly, we bound  $\frac{1}{\mathbb{E}_{\text{Ber}\left(\frac{L+1}{n_{\text{IS}}}\right)}[\mathbf{M}]} \leq \frac{1}{\mathbb{E}_{\text{Ber}\left(\frac{(L+1) \cdot (p+\eta_\delta)}{n_{\text{IS}}}\right)}[\mathbf{M}]}$  to obtain

$$\begin{aligned} \Pr(X_{\ell, I_\ell} = 1) - q &\geq q\delta \frac{p}{q} + (1-\delta) \sum_{\kappa=1}^{\infty} (-1)^\kappa \left( \frac{q}{p} - \frac{1-q}{1-p} \right)^\kappa \left( \frac{1-p+\eta_\delta-n\eta_\delta}{n_{\text{IS}}} \right)^\kappa \Leftrightarrow \\ q - \Pr(X_{\ell, I_\ell} = 1) &\leq -q\delta \frac{p}{q} + (1-\delta) \sum_{\kappa=1}^{\infty} (-1)^{\kappa+1} \left( \frac{q}{p} - \frac{1-q}{1-p} \right)^\kappa \left( \frac{1-p+\eta_\delta-n\eta_\delta}{n_{\text{IS}}} \right)^\kappa \end{aligned}$$

Since  $0 \leq p + \eta_\delta \leq 1$  and  $1 \geq p - \eta_\delta \geq 0$  by an appropriate choice of the concentration intervals, we have by approximations up to second order terms that

$$\begin{aligned} |\Pr(X_{\ell, I_\ell} = 1) - q| &\leq q\delta \max \left\{ \frac{p}{q}, \frac{1-p}{1-q} \right\} + \eta_\delta \left( \frac{q}{p} - \frac{1-q}{1-p} \right) + \left( \frac{q}{p} - \frac{1-q}{1-p} \right)^2 \mathcal{O} \left( \frac{1}{n_{\text{IS}}^2} + \eta_\delta^2 \right) \\ &= \frac{q}{n_{\text{IS}}^2} \left( \frac{p}{q} + \frac{1-p}{1-q} \right) + \mathcal{O} \left( \left[ \left( \frac{q}{p} - \frac{1-q}{1-p} \right) + \left( \frac{q}{p} - \frac{1-q}{1-p} \right)^2 \right] \sqrt{\frac{6p \log(2n_{\text{IS}})}{n_{\text{IS}}}} \right). \end{aligned}$$

This concludes the proof.  $\square$

*Proof of Lemma 1.* Using Lemma 2, we can show the following. Recall the following probability law of the stochastic quantizer  $Q_s(\cdot)$  [Alistarh et al., 2017] using  $s > 0$  quantization intervals, which takes as input the entry  $x_e$  of a gradient  $\mathbf{x} \in \mathbb{R}^d$  vector. Let  $0 \leq \tau_e < s$  be an integer such that  $\frac{\tau_e}{s} \leq \frac{|x_e|}{\|\mathbf{x}\|} \leq \frac{\tau_e+1}{s}$ , then  $Q_s(x_e)$  is defined as  $\text{Ber}\left(\frac{|x_e|}{\|\mathbf{x}\|}s - \tau_e\right)$ , which outputs  $\|\mathbf{x}\| \cdot \text{sign}(x_e)(\tau_e + 1)/s$  in case of success, and  $\|\mathbf{x}\| \cdot \text{sign}(x_e)\tau_e/s$  otherwise.

Focusing on an entry  $x_e$ , we prove a contraction property for MRC with stochastic quantization with posterior  $q_e = \frac{|x_e|}{\|\mathbf{x}\|}s - \tau_e$ , and an arbitrary prior  $p_e$ . In fact, the MRC methodology  $\mathcal{C}_{\text{mrc}}(\cdot)$  leads to sampling from an approximate distribution with parameter  $\tilde{q}_e$ . To be more specific,  $\mathcal{C}_{\text{mrc}}(x_e)$  outputs  $\|\mathbf{x}\| \cdot \text{sign}(x_e)(\tau_e + 1)/s$  with probability  $\tilde{q}_e$ , and  $\|\mathbf{x}\| \cdot \text{sign}(x_e)\tau_e/s$  with probability  $1 - \tilde{q}_e$ . We established in Lemma 2 an upper bound on  $|q_e - \tilde{q}_e|$ , which will be useful in the following.

To prove a contraction property of the kind

$$\mathbb{E}[\|\mathcal{C}_{\text{mrc}}(\mathbf{x}) - \mathbf{x}\|_2^2] \leq (1-\delta)\|\mathbf{x}\|^2,$$

we can write

$$\begin{aligned} \mathbb{E}[\|\mathcal{C}_{\text{mrc}}(\mathbf{x}) - \mathbf{x}\|^2] &= \mathbb{E} \left[ \sum_{e=1}^d (\mathcal{C}_{\text{mrc}}(x_e) - x_e)^2 \right] \\ &= \|\mathbf{x}\|^2 \sum_{e=1}^d \mathbb{E} \left[ \left( \frac{\mathcal{C}_{\text{mrc}}(x_e)}{\|\mathbf{x}\|} - \frac{x_e}{\|\mathbf{x}\|} \right)^2 \right] \\ &= \|\mathbf{x}\|^2 \sum_{e=1}^d \left[ \tilde{q}_e \left( \frac{\text{sign}(x_e)(\tau_e + 1)}{s} - \frac{x_e}{\|\mathbf{x}\|} \right)^2 + (1 - \tilde{q}_e) \left( \frac{\text{sign}(x_e)\tau_e}{s} - \frac{x_e}{\|\mathbf{x}\|} \right)^2 \right] \\ &= \|\mathbf{x}\|^2 \sum_{e=1}^d \left[ (\tilde{q}_e - q_e + q_e) \left( \frac{\tau_e + 1}{s} - \frac{|x_e|}{\|\mathbf{x}\|} \right)^2 + (1 - \tilde{q}_e - q_e + q_e) \left( \frac{\tau_e}{s} - \frac{|x_e|}{\|\mathbf{x}\|} \right)^2 \right] \\ &= \|\mathbf{x}\|^2 \sum_{e=1}^d \left[ (q_e + \tilde{q}_e - q_e) \left( \left( \frac{\tau_e}{s} - \frac{|x_e|}{\|\mathbf{x}\|} \right)^2 + \frac{1}{s^2} + \frac{1}{s} \left( \frac{\tau_e}{s} - \frac{|x_e|}{\|\mathbf{x}\|} \right) \right) \right] \end{aligned}$$

$$\begin{aligned}
& + (1 - q_e + q_e - \tilde{q}_e) \left( \frac{\tau_e}{s} - \frac{|x_e|}{\|\mathbf{x}\|} \right)^2 \Big] \\
& = \|\mathbf{x}\|^2 \sum_{e=1}^d \left[ (\tilde{q}_e - q_e) \left( \frac{1}{s^2} + \frac{1}{s} \left( \frac{\tau_e}{s} - \frac{|x_e|}{\|\mathbf{x}\|} \right) \right) + q_e \left( \frac{1}{s^2} + \frac{1}{s} \left( \frac{\tau_e}{s} - \frac{|x_e|}{\|\mathbf{x}\|} \right) \right) + \left( \frac{\tau_e}{s} - \frac{|x_e|}{\|\mathbf{x}\|} \right)^2 \right], \tag{5}
\end{aligned}$$

where

$$\begin{aligned}
& q_e \left( \frac{1}{s^2} + \frac{1}{s} \left( \frac{\tau_e}{s} - \frac{|x_e|}{\|\mathbf{x}\|} \right) \right) \\
& = \left( \frac{|x_e|}{\|\mathbf{x}\|} s - \tau_e \right) \left( \frac{1}{s^2} + \frac{1}{s} \left( \frac{\tau_e}{s} - \frac{|x_e|}{\|\mathbf{x}\|} \right) \right) \\
& = -s \left( \frac{\tau_e}{s} - \frac{|x_e|}{\|\mathbf{x}\|} \right) \frac{1}{s} \left( \frac{1}{s} + \left( \frac{\tau_e}{s} - \frac{|x_e|}{\|\mathbf{x}\|} \right) \right) \\
& = - \left( \frac{\tau_e}{s} - \frac{|x_e|}{\|\mathbf{x}\|} \right)^2 - \frac{1}{s} \left( \frac{\tau_e}{s} - \frac{|x_e|}{\|\mathbf{x}\|} \right).
\end{aligned}$$

Substituting the result in (5), obtain

$$\begin{aligned}
\mathbb{E}[\|\mathcal{C}_{\text{mrc}}(\mathbf{x}) - \mathbf{x}\|^2] & = \mathbb{E} \left[ \sum_{e=1}^d (\mathcal{C}_{\text{mrc}}(x_e) - x_e)^2 \right] \\
& = \|\mathbf{x}\|^2 \sum_{e=1}^d \left[ (\tilde{q}_e - q_e) \left( \frac{1}{s^2} + \frac{1}{s} \left( \frac{\tau_e}{s} - \frac{|x_e|}{\|\mathbf{x}\|} \right) \right) - \frac{1}{s} \left( \frac{\tau_e}{s} - \frac{|x_e|}{\|\mathbf{x}\|} \right) \right] \\
& = \|\mathbf{x}\|^2 \sum_{e=1}^d \left[ (\tilde{q}_e - q_e) \frac{1}{s} \left( \frac{\tau_e + 1}{s} - \frac{|x_e|}{\|\mathbf{x}\|} \right) - \frac{1}{s} \left( \frac{\tau_e}{s} - \frac{|x_e|}{\|\mathbf{x}\|} \right) \right] \\
& \leq \|\mathbf{x}\|^2 \sum_{e=1}^d \left[ |\tilde{q}_e - q_e| \frac{1}{s} \left( \frac{\tau_e + 1}{s} - \frac{|x_e|}{\|\mathbf{x}\|} \right) + \frac{1}{s} \left( \frac{|x_e|}{\|\mathbf{x}\|} - \frac{\tau_e}{s} \right) \right] \\
& \leq \|\mathbf{x}\|^2 (|\tilde{q}_e - q_e| \frac{d}{s^2} + \frac{d}{s^2}),
\end{aligned}$$

where, by Lemma 2, we have for  $\Delta_e := \frac{q_e}{p_e} - \frac{1-q_e}{1-p_e}$  and  $\Delta'_e := q_e \left( \frac{p_e}{q_e} + \frac{1-p_e}{1-q_e} \right)$  that

$$|\tilde{q}_e - q_e| \leq \frac{\Delta'_e}{n_{\text{IS}}^2} + \mathcal{O} \left( (\Delta_e + \Delta_e^2) \sqrt{\frac{6p_e \log(2n_{\text{IS}})}{n_{\text{IS}}}} \right).$$

Let  $\bar{\Delta} := \max_{e \in [d]} \frac{q_e}{p_e} - \frac{1-q_e}{1-p_e}$ ,  $\bar{\Delta}' := \max_{e \in [d]} q_e \left( \frac{p_e}{q_e} + \frac{1-p_e}{1-q_e} \right)$ , and  $\bar{p} := \max_{e \in [d]} p_e$ . We will ensure that  $\frac{\bar{\Delta}'}{n_{\text{IS}}^2} + \mathcal{O} \left( (\bar{\Delta} + \bar{\Delta}^2) \sqrt{\frac{6\bar{p} \log(2n_{\text{IS}})}{n_{\text{IS}}}} \right) \leq 1$  by making each of the individual terms  $\leq \frac{1}{2}$ . By choosing  $n_{\text{IS}} \geq \sqrt{2\bar{\Delta}'}$ , we have  $\frac{\bar{\Delta}'}{n_{\text{IS}}^2} \leq \frac{1}{2}$ . To ensure that  $(\bar{\Delta} + \bar{\Delta}^2) \sqrt{\frac{6\bar{p} \log(2n_{\text{IS}})}{n_{\text{IS}}}} \leq \frac{1}{2}$ , we require  $\frac{\log(2n_{\text{IS}})}{n_{\text{IS}}} \leq \frac{1}{\sqrt{6\bar{p}(\bar{\Delta} + \bar{\Delta}^2)}}$ . By Weinberger and Yemini [2023, Lemma 15], this holds when  $n_{\text{IS}} = \mathcal{O}(\log(6\bar{p}(\bar{\Delta} + \bar{\Delta}^2)) \sqrt{6\bar{p}(\bar{\Delta} + \bar{\Delta}^2)})$ . Hence, choosing  $n_{\text{IS}} = \mathcal{O}(\max\{\sqrt{2\bar{\Delta}'}, \log(6\bar{p}(\bar{\Delta} + \bar{\Delta}^2)) \sqrt{6\bar{p}(\bar{\Delta} + \bar{\Delta}^2)}\})$ , we have  $\frac{\bar{\Delta}'}{n_{\text{IS}}^2} + \mathcal{O} \left( (\bar{\Delta} + \bar{\Delta}^2) \sqrt{\frac{6\bar{p} \log(2n_{\text{IS}})}{n_{\text{IS}}}} \right) \leq 1$ . Thus, we have  $0 \leq \delta \leq 1$  if  $\frac{2d}{s^2} \leq 1$ , and hence  $s \geq \sqrt{2d}$ . This concludes the proof.  $\square$

*Proof of Theorem 1.* Assume a party estimates the Bernoulli distributions  $Q_j$  with parameters  $q_j$  held by parties  $j \in [n]$ . The estimating party shares with each of the other parties a common prior  $P_j$  in the form of a Bernoulli distribution with parameter  $p_j$  and access to unlimited shared randomness. To help estimate  $Q_j$ , the  $j$ -th party sends  $K$  samples to the estimator through MRC. Therefore, both parties sample  $Kn_{\text{IS}}$  i.i.d. samples  $X_{\ell,i} \sim P_j$  for  $\ell \in [K], i \in [n_{\text{IS}}]$ , independently and identically from  $P_j$ . The party holding  $Q_j$  constructs for each  $\ell \in [K]$  an auxiliary distribution

$$W_\ell(i) = \frac{Q_j(X_{\ell,i})/P_j(X_{\ell,i})}{\sum_{i=1}^{n_{\text{IS}}} Q_j(X_{\ell,i})/P_j(X_{\ell,i})},$$

from which it samples to obtain an index  $I_\ell$ . The index is transmitted to the estimating party, which reconstructs the corresponding sample  $X_{\ell,I_\ell}$ . Averaging the samples for all  $\ell \in [K]$  gives an estimate  $\hat{q}_j$  of  $q_j$ , i.e.,  $\hat{q}_j = \frac{1}{K} \sum_{\ell=1}^K X_{\ell,I_\ell}$ . This process is repeated for all  $j \in [n]$ .

We assume that  $|q_j - p_j| \leq \rho$  for all  $i, j \in [n]$ , and that the difference between the priors, is bounded as  $|p_i - p_j| \leq \zeta$  for all  $i, j \in [n]$ . The goal is to bound  $d_{\text{KL}}\left(\frac{1}{n} \sum_{j=1}^n \hat{q}_j \| p_i\right)$  from above for any  $i \in [n]$ .

By the convexity of KL-divergence, we have

$$d_{\text{KL}}\left(\frac{1}{n} \sum_{j=1}^n \hat{q}_j \| p_i\right) \leq \frac{1}{n} \sum_{i=1}^n d_{\text{KL}}(\hat{q}_j \| p_i).$$

To bound  $d_{\text{KL}}(\hat{q}_j \| p_i)$  for any  $i, j \in [n]$ , by the triangle inequality, we can write

$$|\hat{q}_j - p_i| \leq |\hat{q}_j - \Pr(X_\ell = 1)| + |\Pr(X_\ell = 1) - q_j| + |q_j - p_j| + |p_j - p_i|,$$

where  $|\hat{q}_j - \Pr(X_\ell = 1)|$  is bounded by Lemma 2. By Hoeffding's inequality, we have with probability at least  $1 - \delta'$  that

$$|\hat{q}_j - \Pr(X_\ell = 1)| \leq \sqrt{\frac{-\ln(\delta'/2)}{2n_{\text{IS}}}}.$$

Thus, with probability at least  $1 - \delta'$ , since  $p_j \leq p_i + \zeta$ , we have with  $\Delta_j := \frac{q_j}{p_j - \zeta} - \frac{1 - q_j}{1 - p_j + \zeta}$  and  $\Delta'_j := q_j \left( \frac{p_j + \zeta}{q_j} + \frac{1 - p_j + \zeta}{1 - q_j} \right)$  that

$$|\hat{q}_j - p_i| \leq \frac{\Delta'_j}{n_{\text{IS}}^2} + \mathcal{O}\left((\Delta_j + \Delta_j^2) \sqrt{\frac{6(p_i + \zeta) \log(2n_{\text{IS}})}{n_{\text{IS}}}}\right) + \sqrt{\frac{-\ln(\delta'/2)}{2n_{\text{IS}}}} + \rho + \zeta.$$

This holds under the assumption that  $p_j > \zeta$  for all  $j \in [n]$ . By the reversed Pinsker's inequality, we obtain

$$\begin{aligned} D_{\text{KL}}(\hat{q}_j \| p_i) &\leq \frac{2}{\min\{p_i, 1 - p_i\}} \left( \frac{\Delta'_j}{n_{\text{IS}}^2} + \mathcal{O}\left((\Delta_j + \Delta_j^2) \sqrt{\frac{6(p_i + \zeta) \log(2n_{\text{IS}})}{n_{\text{IS}}}}\right) \right. \\ &\quad \left. + \sqrt{\frac{-\ln(\delta'/2)}{2n_{\text{IS}}}} + \rho + \zeta \right)^2. \end{aligned}$$

The statement of the theorem follows by the convexity of KL-divergence.  $\square$

**Remark 1.** Note that our analysis can be extended to other parametric distributions, such as multi-variate Gaussians. The key ingredient is to replace the specific upper bound on the bias in Lemma 2 with the generic results from Chatterjee and Diaconis [2018], which holds for all classes of distributions. Using this upper bound, one can follow our derivations to prove the communication costs in Theorem 1, i.e., using the convexity of KL divergence and decomposing the error in parameter estimation into a bias term and a concentration term to bound the sampling error. Standard concentration results can be utilized to bound the latter from above, e.g., Hoeffding's inequality for sub-Gaussian random variables. The remaining steps follow analogously to our proof. Similar adaptations apply to the contraction property in Lemma 1 necessary to establish convergence guarantees for conventional FL.

## B Convergence Analysis

Using the contraction property derived in Lemma 1, we can show that a straightforward extension of BiCOMPFL-GR-CFL to error-feedback as used in [Richtárik et al., 2021] leads to the following convergence guarantee. The algorithmic details of the extension can be found in Algorithm 3. Therefore, assume that for all  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$  and  $i \in [n]$ , the following Lipschitz property holds:

$$\|\nabla F(\mathbf{x}, \mathcal{D}_i) - \nabla F(\mathbf{y}, \mathcal{D}_i)\| \leq L_i \|\mathbf{x} - \mathbf{y}\|$$

Let  $F(\theta) := \frac{1}{n} \sum_{i=1}^n \nabla F(\theta, \mathcal{D}_i)$  be the global loss function and  $L' := \sqrt{\frac{1}{n} \sum_{i=1}^n L_i}$ .

**Theorem 2.** *If  $F^* := \inf_{\theta \in \mathbb{R}^d} \{F(\theta)\} > -\infty$  and  $\mathbb{E}[\|\mathbf{g}^t - \nabla F(\theta_t)\|^2] \leq \sigma^2$ , then with  $\eta \leq \left(L + L' \sqrt{\frac{1-\delta}{(1-\sqrt{1-\delta})^2}}\right)^{-1}$ ,  $L = 1$ ,  $s \geq \sqrt{2d}$ , and  $n_{\text{IS}}$  satisfying Lemma 1 in every iteration  $t$ , we have for Algorithm 3 that*

$$\sum_{t=1}^T \mathbb{E} [\|F(\theta_t)\|^2] \leq \frac{2(F(\theta_0) - F^*)}{\eta T} + \frac{\sigma^2}{(1 - \sqrt{1 - \delta})T}.$$

Similarly, guarantees can be derived for other algorithms, such as modified versions of BiCOMPFL-PR with error-feedback and momentum, using Lemma 1. However, we emphasize the generality of BiCOMPFL, reaching beyond conventional FL with stochastic compression to pure stochastic narratives.

**Remark 2.** *The choice of  $s$  depends on the model architecture and dataset complexities, affecting the number of iterations required until convergence. The simpler the learning task, the fewer quantization intervals are sufficient for convergence. Beyond generic arguments, we particularly highlight that  $s$  should be carefully selected, respecting the expected variance of the gradients. If gradients exhibit inherently large variance, e.g., through small mini-batch sizes, large variance of the local datasets, or substantially over-parameterized networks, fewer quantization intervals might be efficient, and increasing that number would only negligibly improve the performance. We also note that one-bit quantization has been shown to be often remarkably effective despite its simplicity [Seide et al., 2014, Karimireddy et al., 2019].*

## C Gradient Descent with a KL-Proximity

Mirror descent employs point-wise optimization in the form of a first-order approximation of  $F(\hat{\theta}_t, \mathcal{D}_i)$  with proximity term  $D_F(p, q)$ , where  $D_F$  is the Bregman divergence associated with function  $F(\cdot)$ . When  $F(x) = \|x\|^2$ , and hence the Bregman divergence is the Euclidean distance, this is known as gradient descent. Let now  $p$  and  $q$  be vectors with the entries corresponding to independent Bernoulli parameters. When we choose  $F(x) = x \log(x) + (1 - x) \log(1 - x)$ , the Bregman divergence becomes  $D_F(p, q) = \sum_{k=1}^d D_{\text{KL}}(p_k \| q_k)$ . Hence, we are optimizing with respect to a KL-proximity constraint. The mapping between dual and primal spaces is then given by  $\nabla F(x) = \log(x) - \log(1 - x)$  and  $(\nabla F(x))^{-1} = \frac{1}{e^{-x} + 1}$ , respectively; also known as the inverse sigmoid and the sigmoid functions.

## D Block Allocation

The simplest yet effective strategy for block allocation is to partition the model into equally-sized blocks of size  $d/B$  for MRC (Fixed). The partitioning into blocks is required to make MRC practically feasible in this setting. It is known that for vanishing MRC error, the number of samples  $n_{\text{IS}}$  from a block  $p_{i,u,b}^t$  of the prior is supposed to be in the order of  $\exp\left(D_{\text{KL}}\left(q_{i,b}^t \| p_{i,u,b}^t\right)\right)$ , where  $q_{i,b}^t$  is the  $b$ -th block of posterior  $q_i^t$ . It was observed by Isik et al. [2024] that the KL-divergence decreases as the training progresses with the global model used as a prior, which is intuitive since the local training will change the posterior less and less as training converges. To adapt the block size according to the

---

**Algorithm 3** BiCOMPFL-GR-CFL with stochastic quantization  $Q_s(\cdot)$  and EF21 from Richtárik et al. [2021]

---

**Require:** Both clients and federator initialize the same global model  $\theta_0$  using a shared seed

**Ensure:** Set  $t = 0$ , clients set prior  $p^t = \hat{\theta}_0 = \theta_0, \forall i \in [n]$ , clients compute and broadcast  $\mathbf{v}_i^0 = \mathcal{C}_{\text{mrc}}(Q_s(g_i^t), p^t)$ , with  $g_i^t$  the local gradient for  $\theta_0$ ; hence,  $\mathbf{v}^0 = \frac{1}{n} \sum_{i=1}^n \mathbf{v}_i^0$  public

```

1: Update  $\forall i : \hat{\theta}_{t+1} = \hat{\theta}_t - \eta \mathbf{v}^{t+1}$ 
2: repeat
3:   for Client  $i \in [n]$  do
4:     Compute gradient  $g_i^t$  by local training over  $L$  local iterations
5:     Stochastic compression  $q_i^t \leftarrow Q_s(g_i^t - \mathbf{v}_i^t)$ 
6:     Sample indices  $I_{i,\ell}^b, \ell \in [n_{\text{UL}}], b \in [B]$  from  $q_i^t$  with prior  $p^t$  and transmit to federator to
       reconstruct  $\hat{q}_i^t = \mathcal{C}_{\text{mrc}}(q_i^t, p^t)$ 
7:     Update  $\mathbf{v}_i^{t+1} = \mathbf{v}_i^t + \hat{q}_i^t$ 
8:   end for
9:   Federator reconstructs and computes  $\mathbf{v}^{t+1} = \mathbf{v}^t + \frac{1}{n} \sum_{j=1}^n \hat{q}_j^t$  from  $\{I_{i,\ell}^b\}$ 
10:  Federator updates  $\theta_{t+1} = \theta_t - \eta \mathbf{v}^{t+1}$ 
11:  Federator relays to client  $j$  the other clients' indices  $\{I_{i,\ell}^b\}_{\ell \in [n_{\text{UL}}], b \in [B], i \in [n] \setminus \{j\}}$ 
12:  for Clients  $i \in [n]$  do
13:    Reconstruct and compute  $\mathbf{v}^{t+1} = \mathbf{v}^t + \frac{1}{n} \sum_{j=1}^n \hat{q}_j^t$  from  $\{I_{i,\ell}^b\}$ 
14:    Update  $\hat{\theta}_{t+1} = \hat{\theta}_t - \eta \mathbf{v}^{t+1}$  from  $\{I_{i,\ell}^b\}$ 
15:  end for
16:  Clients and federator set prior  $p^t = \hat{\theta}_{t+1}$ 
17:   $t \leftarrow t + 1$ 
18: until Convergence

```

---

divergence from the posterior with respect to the prior, Isik et al. [2024] proposed an adaptive block allocation strategy (Adaptive), where upon realizing a large deviation from the target KL-divergence per block, clients partition their model into blocks with equal sums of parameter-wise KL-divergences and transmit the block intervals to the federator. The federator aggregates the indices of all the clients, and broadcasts the updated block allocation.

We propose in this work a low complexity solution that adapts the block size according to the average KL-divergence per block (Adaptive-Avg). This alleviates the cost of computing and transmitting the exact block partitions, where the transmission of each block size requires  $\log_2(b_{\text{max}})$  bits, with  $b_{\text{max}}$  the maximum pre-defined block size. Instead, the transmission of one size is enough in our solution. If the average KL per block  $D_{\text{KL}}(q_{i,b}^t \| p_{i,u,b}^t)$  deviates more than a given factor, the clients request to update the blocks. In the next iteration, each client proposes a block size, and the federator averages and broadcasts an updated size.

## E Additional Experimental Details

We use the cross-entropy loss and a batch size of 128 in all our experiments. We use Adam [Kingma and Ba, 2015] as an optimizer with learning rate  $\eta = 0.0003$  for all non-stochastic methods, and  $\eta = 0.1$  for probabilistic mask training. For non-stochastic FL, we use a federator (server) learning rate of 0.1, i.e., the clients' gradients are averaged, and the federator updates the global model with learning rate 0.1, and with a learning rate of 0.005 for BiCOMPFL-GR with SignSGD. For M3, we use a federator learning rate of 0.02 to obtain reliable results. For LIEC and CSER, we use an average period of 50 global iterations (cf. [Cheng et al., 2024, Xie et al., 2020]). For M3, we use TopK with  $K = \lfloor d/n \rfloor$ . To run the simulations, we use a cluster of different architectures, which we list in the following table.

Table 1: System specifications of our simulation cluster.

CPU(s)	RAM	GPU(s)	VRAM
2x Intel Xeon Platinum 8176 (56 cores)	256 GB	2x NVIDIA GeForce GTX 1080 Ti	11 GB
2x AMD EPYC 7282 (32 cores)	512 GB	NVIDIA GeForce RTX 4090	24 GB
2x AMD EPYC 7282 (32 cores)	640 GB	NVIDIA GeForce RTX 4090	24 GB
2x AMD EPYC 7282 (32 cores)	448 GB	NVIDIA GeForce RTX 4080	16 GB
2x AMD EPYC 7282 (32 cores)	256 GB	NVIDIA GeForce RTX 4080	16 GB
HGX-A100 (96 cores)	1 TB	4x NVIDIA A100	80 GB
DGX-A100 (252 cores)	2 TB	8x NVIDIA Tesla A100	80 GB
DGX-1-V100 (76 cores)	512 GB	8x NVIDIA Tesla V100	16 GB
DGX-1-P100 (76 cores)	512 GB	8x NVIDIA Tesla P100	16 GB
HPE-P100 (28 cores)	256 GB	4x NVIDIA Tesla P100	16 GB

The details of the CNN architectures used in our experiments are summarized in the following. The parameter count is 61706 for LeNet5, 1933258 for 4CNN, and 2262602 for 6CNN.

Table 2: LeNet5 Architecture Overview

Layer	Specification	Activation
5x5 Conv	6 filters, stride 1	ReLU, AvgPool (2x2)
5x5 Conv	16 filters, stride 1	ReLU, AvgPool (2x2)
Linear	120 units	ReLU
Linear	84 units	ReLU
Linear	10 units	Softmax

Table 3: 4-layer CNN (4CNN) Architecture Overview

Layer	Specification	Activation
3x3 Conv	64 filters, stride 1	ReLU
3x3 Conv	64 filters, stride 1	ReLU, MaxPool (2x2)
3x3 Conv	128 filters, stride 1	ReLU
3x3 Conv	128 filters, stride 1	ReLU, MaxPool (2x2)
Linear	256 units	ReLU
Linear	256 units	ReLU
Linear	10 units	Softmax

Table 4: 6-layer CNN (6CNN) Architecture Overview

Layer	Specification	Activation
3x3 Conv	64 filters, stride 1	ReLU
3x3 Conv	64 filters, stride 1	ReLU, MaxPool (2x2)
3x3 Conv	128 filters, stride 1	ReLU
3x3 Conv	128 filters, stride 1	ReLU, MaxPool (2x2)
3x3 Conv	256 filters, stride 1	ReLU
3x3 Conv	256 filters, stride 1	ReLU, MaxPool (2x2)
Linear	256 units	ReLU
Linear	256 units	ReLU
Linear	10 units	Softmax

For the sake of clarity, in the paper we restrict the analysis to a fixed number of importance samples  $n_{IS}$ , block sizes  $B$ , and choice of priors  $p_{i,u}^t, p_{i,d}^t$ . Our experiments have shown that, while increasing  $n_{IS}$  beyond the ones used in our algorithms slightly improves the convergence over the number of epochs, the convergence with respect to the communication cost did not significantly improve. The block size is mainly limited by the system resources at hand, and one would choose the largest possible for best efficiency while complying with memory resources. We investigated many different

---

**Algorithm 4** Local Training at Client  $i$ 

---

**Require:** Model  $\hat{\theta}_{i,t}$

- 1: Map model to scores in the dual space:  $\mathbf{s}_{i,t}^{(0)} = \sigma^{-1}(\hat{\theta}_{i,t}) = \log\left(\frac{\hat{\theta}_{i,t}}{1-\hat{\theta}_{i,t}}\right)$
  - 2: **for** Local iterations  $m \in [L]$  **do**
  - 3:  $\mathbf{s}_{i,t}^{(m)} = \mathbf{s}_{i,t}^{(m-1)} - \eta \nabla_{\mathbf{s}_{i,t}^{(m-1)}} F(\hat{\theta}_{i,t}^{(m-1)}, \mathcal{D}_i)$ , where  $\hat{\theta}_{i,t}^{(m-1)} = \sigma(\mathbf{s}_{i,t}^{(m-1)})$
  - 4: **end for**
  - 5: Map back to primal space:  $q_i^t = \sigma(\mathbf{s}_{i,t}^{(L)})$
- 

prior choices and found the former global model to be reasonably good in almost all cases. With high heterogeneity, it might be beneficial to use different convex combinations as priors, which mix the former global model with the latest posterior estimate of a certain client, but the gains we experienced were minor. Hence, we settled on the former global estimate for simplicity in presenting the algorithm.

## F Federated Probabilistic Mask Training

The idea in federated probabilistic mask training (FedPM) Isik et al. [2023] is to collaboratively train a probabilistic mask that determines which weights to maintain from a randomly initialized network. The motivation stems from the *lottery-ticket hypothesis* [Frankle and Carbin, 2019], which claims that randomly initialized networks contain sub-networks capable of reaching accuracy comparable to that of the full network. The weights  $w$  of the network are randomly initialized at the start of training, and remain fixed. The federator and clients only train a mask, which determines for each parameter whether it is activated or not, i.e., identifying an efficient subnetwork within the given fixed network. The probabilistic masks  $\theta_t$  are described by Bernoulli distributions, i.e.,  $\theta_t \in [0, 1]^d$  contains a Bernoulli parameter to be trained for each weight of the network. These parameters determine the probability of retaining the corresponding weights. During inference, the weights  $w$  are masked with samples  $x^t \in \{0, 1\}^d \sim \theta_t$  from the distribution  $\theta_t$ , i.e., the inference is conducted on a network with weights  $w \odot x^t$ . In FedPM, clients sample from their locally trained models, and send these samples to the federator, which, in turn, updates the global model by averaging these samples. The communication cost of this scheme is fixed for all iterations, even though the communication cost can be reduced since the KL-divergence between the global model and the locally trained models diminishes as the training progresses.

We adopt the following federated learning procedure for collaboratively learning network masks, and highlight in the following the parallels to mirror descent by referring to primal and dual spaces. Starting from a common model  $\theta_0$ , at iteration  $t$ , each client  $i$  locally trains the model  $\hat{\theta}_{i,t}$  in  $L$  local iterations. To enable gradient descent, the model  $\hat{\theta}_{i,t}$  is mapped to scores  $\mathbf{s}_{i,t}^{(0)}$  in a dual space by the inverse Sigmoid function  $\mathbf{s}_{i,t}^{(0)} = \sigma^{-1}(\hat{\theta}_{i,t}) = \log(\hat{\theta}_{i,t}) - \log(1 - \hat{\theta}_{i,t})$ . The scores are then trained for  $L$  local iterations  $m \in [L]$  by computing the gradient  $\nabla_{\mathbf{s}_{i,t}^{(m-1)}} F(\hat{\theta}_{i,t}^{(m-1)}, \mathcal{D}_i)$ , where the straight-through estimator is used to compute the gradient of the non-differentiable Bernoulli sampling operation based on the distribution  $\hat{\theta}_{i,t}^{(m-1)} = \sigma(\mathbf{s}_{i,t}^{(m-1)})$ , i.e., the gradient equals the Bernoulli parameter. By mapping the model back to the primal space, each client  $i$  obtains a model update in terms of a posterior  $q_i^t = \sigma(\mathbf{s}_{i,t}^{(L)})$ . The client training process is summarized in Algorithm 4.

## G Minimal Random Coding (MRC)

Isik et al. [2024] proposed a method, called KL minimization with side information (KLMS), to reduce the cost of transmitting the local models  $q_i^t$  to the federator. Consequently, the communication cost depends on the KL-divergence between the desired distribution and the common prior. This method utilizes the common side information available at both the clients and the federator, as well as shared randomness. The idea is that instead of sampling locally and sending the samples to the

federator, the federator in the KLMS method samples from the desired distribution through MRC. In a nutshell, MRC [Havasi et al., 2019] is based on importance sampling [Srinivasan, 2002] and makes use of a common prior to sample from a desired distribution. Consider two distributions  $P$  and  $Q$ , where  $P$  is known to both parties, and  $Q$  is only known to the client. To make the federator sample from  $Q$ , both parties sample  $n_{IS}$  samples  $\{X_i\}_{i \in [n_{IS}]}$  from  $P$ . The client forms an auxiliary distribution  $W(i) = \frac{Q(X_i)/P(X_i)}{\sum_{i=1}^{n_{IS}} Q(X_i)/P(X_i)}$  capturing the importance of the samples. A sample from  $W$  is fully described by its index  $i$ , which can be transmitted with  $\log_2(n_{IS})$  bits, and approximates a sample from  $Q$ . Chatterjee and Diaconis [2018] shown that importance sampling with posterior  $Q$  and prior  $P$  requires  $n_{IS}$  to be in the order of  $\Theta(\exp(D_{KL}(Q||P)))$ , where  $D_{KL}(Q||P)$  denotes the KL-divergence between distributions  $Q$  and  $P$ . In what follows, we will also denote the KL-divergence between two Bernoulli distributions  $Q$  and  $P$  with parameters  $q$  and  $p$  by  $d_{KL}(q||p)$ .

## H Additional Experiments

We provide in the following experiments for both uniform (i.i.d.) and heterogeneous (non-i.i.d.) data distributions for training LeNet5 and a 4-layer CNN on MNIST, a 4-layer CNN on Fashion MNIST, and a 6-layer CNN on CIFAR-10. The details of the neural networks can be found in Tables 2 to 4. For each setting and method depicted, we show the average of three simulation runs with different seeds. We plot for each setting the test accuracies over the communication cost in bits, and the maximum test accuracy over the bitrate. We provide tables summarizing the maximum test accuracies with their standard deviation over multiple runs, the total bitrates and the bitrates split into uplink and downlink. The overall bitrates per parameter (bpp) are computed assuming point-to-point links between all participants, i.e., uplink and downlink costs have equal weight. For the case when a broadcast (BC) link between the federator and the clients is available, the bitrate per parameter for all baseline schemes reduces by a factor of  $n$ . BiCompPFL-GR profits similarly from the broadcast link, but BiCompPFL-PR cannot profit due to the absence of shared randomness, giving the same overall bitrate compared to the point-to-point link scenario. We highlight for each of the measures the scheme with the best result. Consistently throughout all experiments, BiCompPFL achieves order-wise savings in the bitrates per parameter while reaching state-of-the-art accuracies in the classification task. While the sampling can introduce an additional computational overhead depending on the implementation, the storage cost is similar to the baselines. Since we leverage as priors the former global model, the additional storage cost incurred is limited to storing until the next iteration the estimate of the former global model at each client, i.e., where the training started, which is usually not a bottleneck. This can be cheaper than some baselines, which require storing data for momentum and error-feedback.

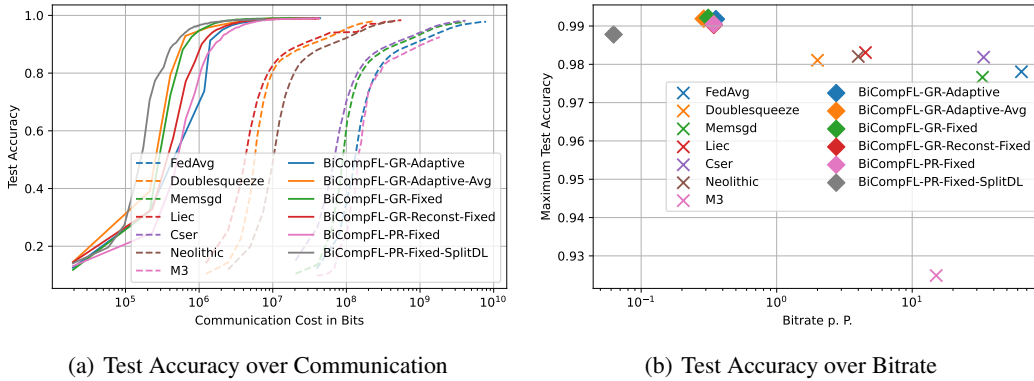


Figure 3: MNIST LeNet i.i.d.

For LeNet5 on MNIST, it can be observed that all our proposed methods converge significantly faster to satisfying accuracies with respect to the communication cost, while achieving higher maximum accuracies after 200 epochs than the non-stochastic baselines. Partitioning the model on the downlink can help to further reduce the communication cost with only a minor loss in performance, especially in the i.i.d. setting. For non-i.i.d. data distribution, the loss in performance is larger than for i.i.d.



distribution. However, at the beginning of the training, the model improves faster with respect to the communication cost than all other schemes. The bitrates are comparable for all our methods, with the exception of BiCompFL-PR-Fixed-SplitDL. Further, BiCompFL-GR-Reconst-Fixed does not suffer notable performance degradation from employing an additional MRC step (especially for i.i.d. data allocation).

Table 5: MNIST LeNet i.i.d.

Method	Acc (mean $\pm$ std)	bpp	bpp (BC)	Uplink	Downlink
FedAvg	$0.978 \pm 0.1$	64.0	35.0	32.0	32.0
Doublesqueeze	$0.981 \pm 0.1$	2.0	1.1	1.0	1.0
Memsgd	$0.977 \pm 0.1$	33.0	4.2	1.0	32.0
Liec	$0.983 \pm 0.1$	4.5	2.5	2.3	2.3
Cser	$0.982 \pm 0.09$	34.0	4.3	1.0	33.0
Neolithic	$0.982 \pm 0.1$	4.0	2.2	2.0	2.0
M3	$0.925 \pm 0.2$	15.0	2.2	8.0	7.1
BiCompFL-GR-Adaptive	<b><math>0.992 \pm 0.0006</math></b>	0.36	0.068	0.036	0.32
BiCompFL-GR-Adaptive-Avg	$0.992 \pm 0.0003$	0.29	<b>0.055</b>	<b>0.029</b>	0.26
BiCompFL-GR-Fixed	$0.992 \pm 0.0002$	0.31	0.059	0.031	0.28
BiCompFL-GR-Reconst-Fixed	$0.99 \pm 0.0002$	0.34	0.063	0.031	0.31
BiCompFL-PR-Fixed	$0.99 \pm 0.0004$	0.34	0.34	0.031	0.31
BiCompFL-PR-Fixed-SplitDL	$0.988 \pm 0.0009$	<b>0.063</b>	0.063	0.031	<b>0.031</b>

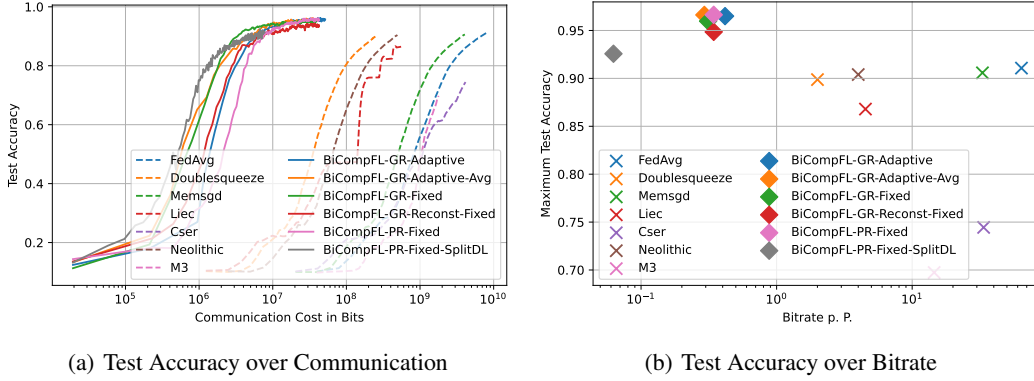


Figure 4: MNIST LeNet non-i.i.d.

Table 6: MNIST LeNet non-i.i.d.

Method	Acc (mean $\pm$ std)	bpp	bpp (BC)	Uplink	Downlink
FedAvg	$0.911 \pm 0.2$	64.0	35.0	32.0	32.0
Doublesqueeze	$0.899 \pm 0.2$	2.0	1.1	1.0	1.0
Memsgd	$0.906 \pm 0.2$	33.0	4.2	1.0	32.0
Liec	$0.866 \pm 0.2$	4.5	2.5	2.3	2.3
Cser	$0.744 \pm 0.2$	34.0	4.3	1.0	33.0
Neolithic	$0.904 \pm 0.2$	4.0	2.2	2.0	2.0
M3	$0.697 \pm 0.2$	15.0	2.2	7.3	7.2
BiCompFL-GR-Adaptive	$0.965 \pm 0.02$	0.42	0.079	0.042	0.37
BiCompFL-GR-Adaptive-Avg	<b><math>0.966 \pm 0.02</math></b>	0.29	<b>0.056</b>	<b>0.029</b>	0.26
BiCompFL-GR-Fixed	$0.96 \pm 0.03$	0.31	0.059	0.031	0.28
BiCompFL-GR-Reconst-Fixed	$0.949 \pm 0.03$	0.34	0.063	0.031	0.31
BiCompFL-PR-Fixed	$0.966 \pm 0.02$	0.34	0.34	0.031	0.31
BiCompFL-PR-Fixed-SplitDL	$0.926 \pm 0.04$	<b>0.063</b>	0.063	0.031	<b>0.031</b>

For 4CNN trained on MNIST, the differences between the proposed approaches become more visible. In the i.i.d. setting, we can observe that the adaptive block allocations (both Adaptive and Adaptive-Avg) can drastically reduce the average bitrate in BiCompFL-GR. Partitioning the model in the downlink (BiCompFL-PR-Fixed-SplitDL) improves the accuracy over bitrate significantly compared to BiCompFL-PR-Fixed.

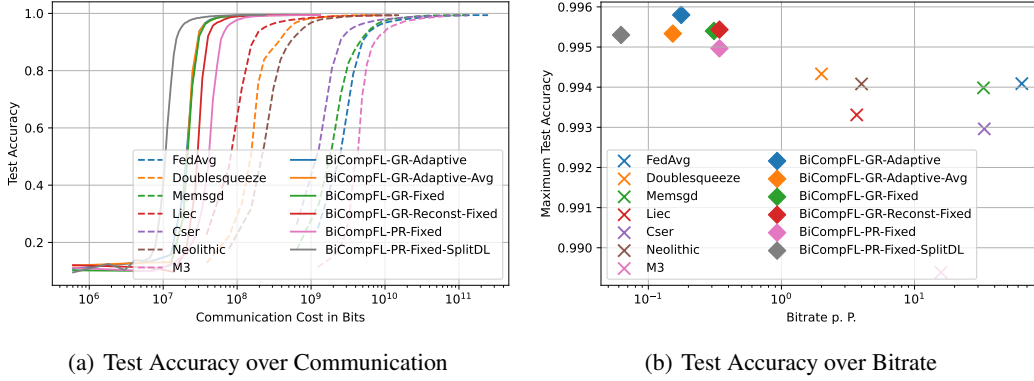


Figure 5: MNIST 4CNN i.i.d.

Table 7: MNIST 4CNN i.i.d.

Method	Acc (mean $\pm$ std)	bpp	bpp (BC)	Uplink	Downlink
FedAvg	$0.994 \pm 0.06$	64.0	35.0	32.0	32.0
Doublesqueeze	$0.994 \pm 0.1$	2.0	1.1	1.0	1.0
Memsgd	$0.994 \pm 0.08$	33.0	4.2	1.0	32.0
Liec	$0.993 \pm 0.07$	3.7	2.0	1.8	1.8
Cser	$0.993 \pm 0.06$	33.0	4.3	1.0	32.0
Neolithic	$0.994 \pm 0.08$	4.0	2.2	2.0	2.0
M3	$0.989 \pm 0.2$	16.0	2.2	8.4	7.4
BiCompFL-GR-Adaptive	<b><math>0.996 \pm 0.0001</math></b>	0.18	0.034	0.018	0.16
BiCompFL-GR-Adaptive-Avg	$0.995 \pm 0.0001$	0.15	<b>0.029</b>	<b>0.015</b>	0.14
BiCompFL-GR-Fixed	$0.995 \pm 0.0002$	0.31	0.059	0.031	0.28
BiCompFL-GR-Reconst-Fixed	$0.995 \pm 0.0001$	0.34	0.062	0.031	0.31
BiCompFL-PR-Fixed	$0.995 \pm 0.0002$	0.34	0.34	0.031	0.31
BiCompFL-PR-Fixed-SplitDL	$0.995 \pm 0.0002$	<b>0.062</b>	0.062	0.031	<b>0.031</b>

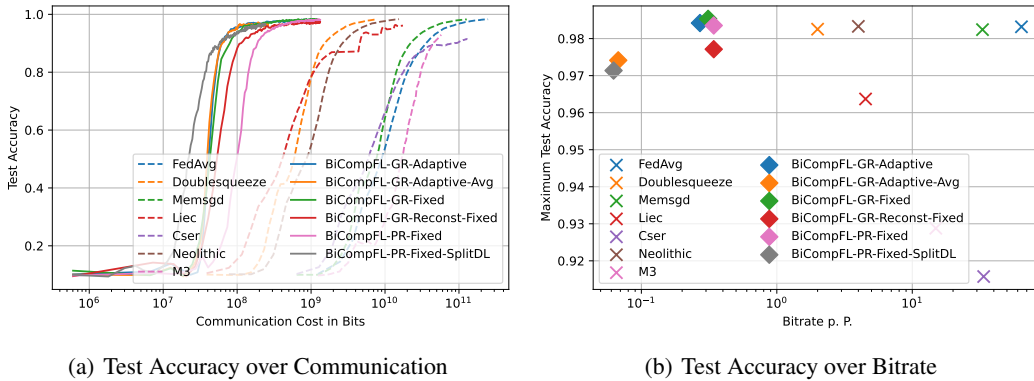


Figure 6: MNIST 4CNN non-i.i.d.

In the non-i.i.d. case of 4CNN on MNIST, the adaptive average allocation strategy provides a significant reduction in the bitrate for BiCompFL-GR, with similar loss in the accuracy as SplitDL for BiCompFL-PR. In this setting, it is also apparent that the reconstruction in BiCompFL-GR degrades the performance without gains in the bitrate compared to the proposed Algorithm 1.

Table 8: MNIST 4CNN non-i.i.d.

Method	Acc (mean $\pm$ std)	bpp	bpp (BC)	Uplink	Downlink
FedAvg	$0.983 \pm 0.1$	64.0	35.0	32.0	32.0
Doublesqueeze	$0.982 \pm 0.2$	2.0	1.1	1.0	1.0
Memsgd	$0.982 \pm 0.2$	33.0	4.2	1.0	32.0
Liec	$0.963 \pm 0.2$	4.5	2.5	2.3	2.3
Cser	$0.915 \pm 0.1$	34.0	4.3	1.0	33.0
Neolithic	$0.983 \pm 0.2$	4.0	2.2	2.0	2.0
M3	$0.929 \pm 0.3$	15.0	2.2	7.8	7.1
BiCompFL-GR-Adaptive	$0.984 \pm 0.009$	0.27	0.051	0.026	0.24
BiCompFL-GR-Adaptive-Avg	$0.974 \pm 0.02$	0.067	<b>0.013</b>	<b>0.0068</b>	0.061
BiCompFL-GR-Fixed	<b>0.985 <math>\pm</math> 0.008</b>	0.31	0.059	0.031	0.28
BiCompFL-GR-Reconst-Fixed	$0.977 \pm 0.01$	0.34	0.062	0.031	0.31
BiCompFL-PR-Fixed	$0.984 \pm 0.009$	0.34	0.34	0.031	0.31
BiCompFL-PR-Fixed-SplitDL	$0.971 \pm 0.02$	<b>0.062</b>	0.062	0.031	<b>0.031</b>

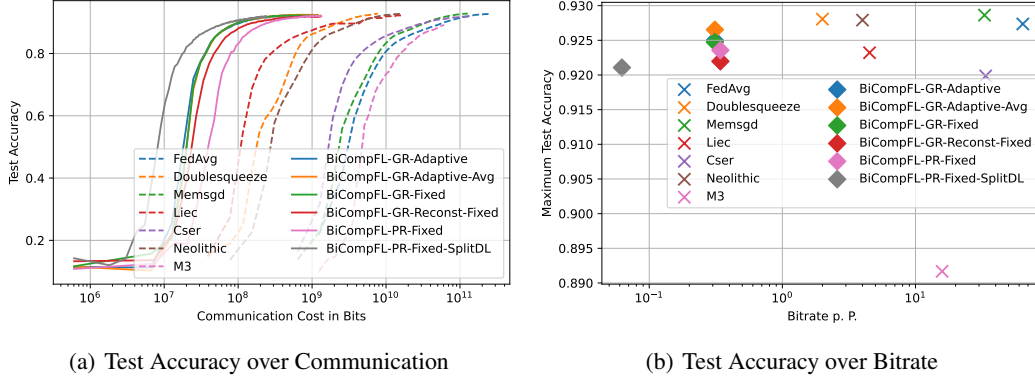


Figure 7: Fashion MNIST 4CNN i.i.d.

Table 9: Fashion MNIST 4CNN i.i.d.

Method	Acc (mean $\pm$ std)	bpp	bpp (BC)	Uplink	Downlink
FedAvg	$0.927 \pm 0.07$	64.0	35.0	32.0	32.0
Doublesqueeze	<b>0.928 <math>\pm</math> 0.1</b>	2.0	1.1	1.0	1.0
Memsgd	$0.928 \pm 0.09$	33.0	4.2	1.0	32.0
Liec	$0.923 \pm 0.08$	4.5	2.5	2.3	2.3
Cser	$0.92 \pm 0.08$	34.0	4.3	1.0	33.0
Neolithic	$0.928 \pm 0.09$	4.0	2.2	2.0	2.0
M3	$0.892 \pm 0.2$	16.0	2.2	8.3	7.6
BiCompFL-GR-Adaptive	$0.925 \pm 0.001$	0.31	<b>0.059</b>	<b>0.031</b>	0.28
BiCompFL-GR-Adaptive-Avg	$0.927 \pm 0.0007$	0.31	<b>0.059</b>	<b>0.031</b>	0.28
BiCompFL-GR-Fixed	$0.925 \pm 0.0007$	0.31	<b>0.059</b>	<b>0.031</b>	0.28
BiCompFL-GR-Reconst-Fixed	$0.922 \pm 0.001$	0.34	0.062	<b>0.031</b>	0.31
BiCompFL-PR-Fixed	$0.924 \pm 0.002$	0.34	0.34	<b>0.031</b>	0.31
BiCompFL-PR-Fixed-SplitDL	$0.921 \pm 0.002$	<b>0.062</b>	0.062	<b>0.031</b>	<b>0.031</b>

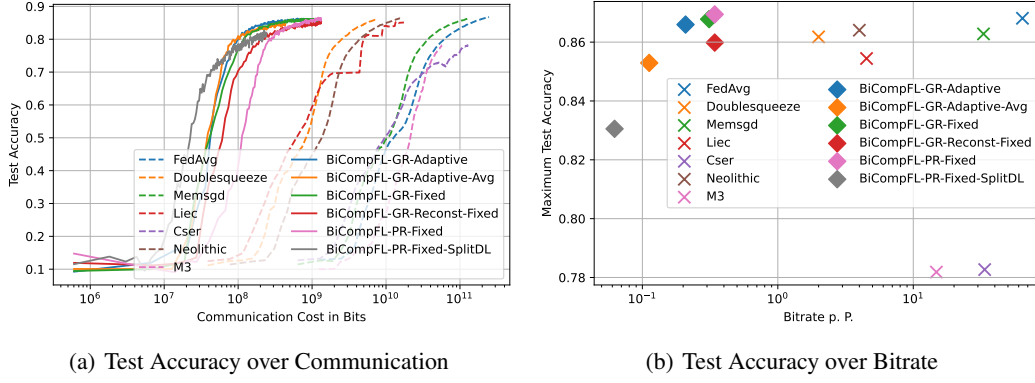


Figure 8: Fashion MNIST 4CNN non-i.i.d.

The results for Fashion MNIST are similar compared to the MNIST case. However, it becomes clear that BiCOMPFL-PR can significantly suffer from the unavailability of shared randomness in terms of the achieved accuracy when data is highly heterogeneous.

Table 10: Fashion MNIST 4CNN non-i.i.d.

Method	Acc (mean $\pm$ std)	bpp	bpp (BC)	Uplink	Downlink
FedAvg	$0.867 \pm 0.1$	64.0	35.0	32.0	32.0
Doublesqueeze	$0.861 \pm 0.2$	2.0	1.1	1.0	1.0
Memsgd	$0.863 \pm 0.2$	33.0	4.2	1.0	32.0
Liec	$0.853 \pm 0.1$	4.5	2.5	2.3	2.3
Cser	$0.781 \pm 0.1$	34.0	4.3	1.0	33.0
Neolithic	$0.864 \pm 0.2$	4.0	2.2	2.0	2.0
M3	$0.782 \pm 0.2$	15.0	2.2	8.0	6.9
BiCompFL-GR-Adaptive	$0.866 \pm 0.03$	0.21	0.04	0.021	0.19
BiCompFL-GR-Adaptive-Avg	$0.853 \pm 0.04$	0.11	<b>0.021</b>	<b>0.011</b>	0.1
BiCompFL-GR-Fixed	$0.868 \pm 0.03$	0.31	0.059	0.031	0.28
BiCompFL-GR-Reconst-Fixed	$0.86 \pm 0.02$	0.34	0.062	0.031	0.31
BiCompFL-PR-Fixed	<b><math>0.869 \pm 0.03</math></b>	0.34	0.34	0.031	0.31
BiCompFL-PR-Fixed-SplitDL	$0.831 \pm 0.03$	<b>0.062</b>	0.062	0.031	<b>0.031</b>

For 6CNN trained on CIFAR-10, the negative effects of missing global shared randomness and reconstructing in the case of BiCOMPFL-GR are prominent. For non-i.i.d. data distributions, the adaptive average allocation shows improvements over the fixed or the average block allocation. Partitioning the model is not a viable option in this setting, especially under non-i.i.d. data.

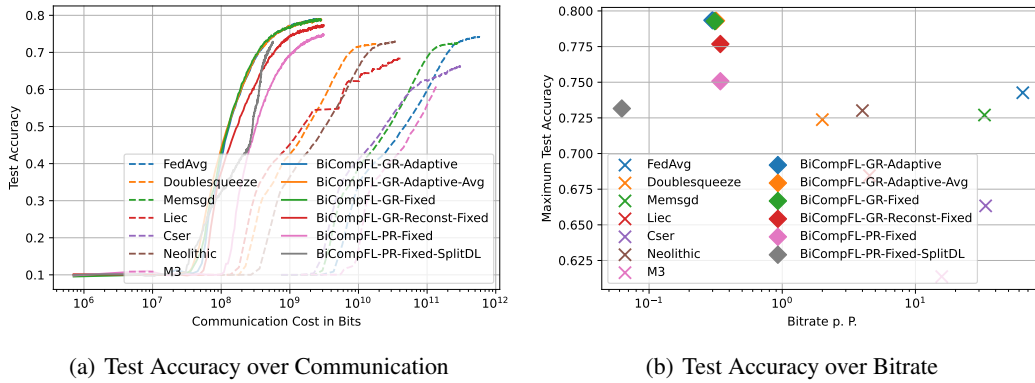


Figure 9: CIFAR-10 6CNN i.i.d.

Table 11: CIFAR-10 6CNN i.i.d.

Method	Acc (mean $\pm$ std)	bpp	bpp (BC)	Uplink	Downlink
FedAvg	$0.742 \pm 0.1$	64.0	35.0	32.0	32.0
Doublesqueeze	$0.723 \pm 0.1$	2.0	1.1	1.0	1.0
Memsgd	$0.727 \pm 0.1$	33.0	4.2	1.0	32.0
Liec	$0.684 \pm 0.09$	4.5	2.5	2.3	2.3
Cser	$0.663 \pm 0.08$	34.0	4.3	1.0	33.0
Neolithic	$0.73 \pm 0.1$	4.0	2.2	2.0	2.0
M3	$0.614 \pm 0.1$	16.0	2.2	8.3	7.5
BiCompFL-GR-Adaptive	<b><math>0.793 \pm 0.002</math></b>	0.3	<b>0.057</b>	<b>0.03</b>	0.27
BiCompFL-GR-Adaptive-Avg	$0.793 \pm 0.002$	0.32	0.061	0.032	0.29
BiCompFL-GR-Fixed	$0.793 \pm 0.004$	0.31	0.059	0.031	0.28
BiCompFL-GR-Reconst-Fixed	$0.777 \pm 0.002$	0.34	0.062	0.031	0.31
BiCompFL-PR-Fixed	$0.751 \pm 0.003$	0.34	0.34	0.031	0.31
BiCompFL-PR-Fixed-SplitDL	$0.732 \pm 0.02$	<b>0.062</b>	0.062	0.031	<b>0.031</b>

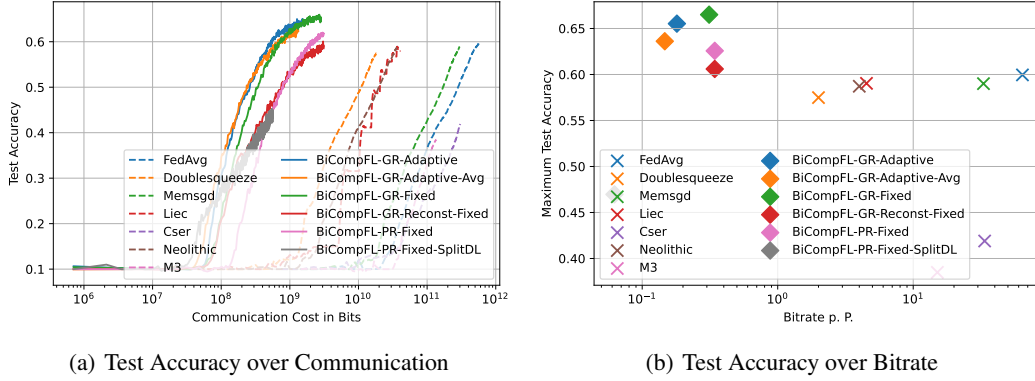
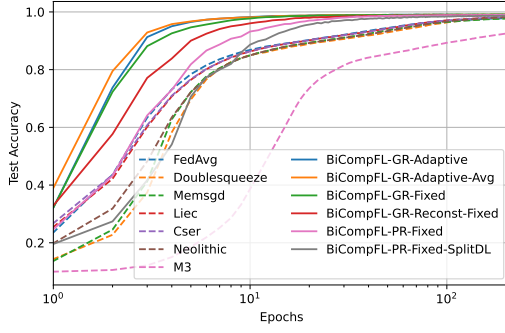


Figure 10: CIFAR-10 6CNN non-i.i.d.

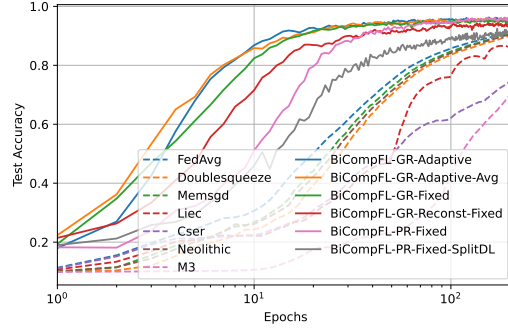
Table 12: CIFAR-10 6CNN non-i.i.d.

Method	Acc (mean $\pm$ std)	bpp	bpp (BC)	Uplink	Downlink
FedAvg	$0.599 \pm 0.1$	64.0	35.0	32.0	32.0
Doublesqueeze	$0.575 \pm 0.1$	2.0	1.1	1.0	1.0
Memsgd	$0.589 \pm 0.1$	33.0	4.2	1.0	32.0
Liec	$0.589 \pm 0.2$	4.5	2.5	2.3	2.3
Cser	$0.419 \pm 0.09$	34.0	4.3	1.0	33.0
Neolithic	$0.587 \pm 0.1$	4.0	2.2	2.0	2.0
M3	$0.385 \pm 0.1$	15.0	2.2	8.3	6.7
BiCompFL-GR-Adaptive	$0.655 \pm 0.04$	0.18	0.034	0.018	0.16
BiCompFL-GR-Adaptive-Avg	$0.636 \pm 0.05$	0.15	<b>0.028</b>	<b>0.015</b>	0.13
BiCompFL-GR-Fixed	<b><math>0.665 \pm 0.03</math></b>	0.31	0.059	0.031	0.28
BiCompFL-GR-Reconst-Fixed	$0.606 \pm 0.05$	0.34	0.062	0.031	0.31
BiCompFL-PR-Fixed	$0.626 \pm 0.03$	0.34	0.34	0.031	0.31
BiCompFL-PR-Fixed-SplitDL	$0.47 \pm 0.07$	<b>0.062</b>	0.062	0.031	<b>0.031</b>

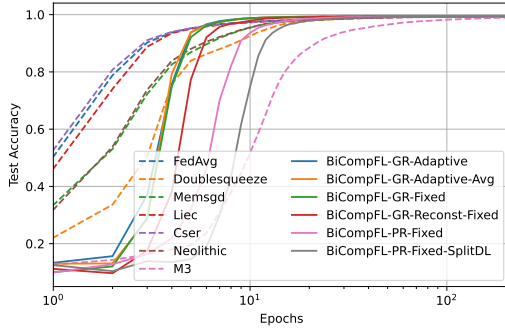
For completeness, we present in Fig. 11 the test accuracies over the number of trained epochs for all scenarios considered above. The setting of interest to this work is that of limited communication cost, and in particular, which performance is achievable given a fixed communication budget. Nonetheless, we can find that our proposed methods are not inferior in convergence speed over epochs compared to the baselines.



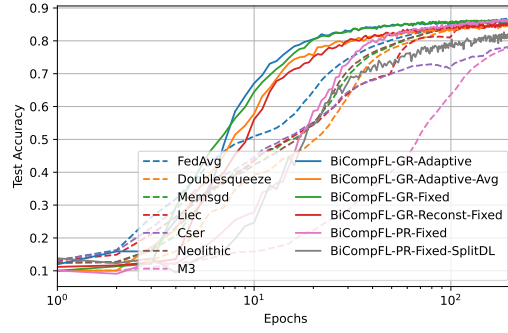
(a) MNIST LeNet i.i.d.



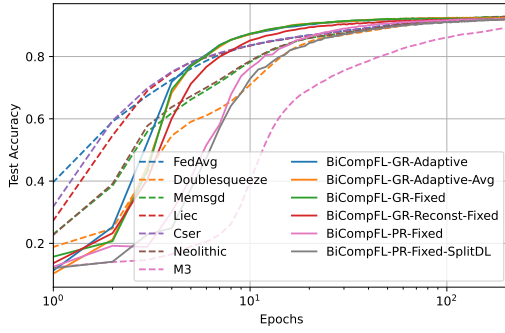
(b) MNIST LeNet non-i.i.d.



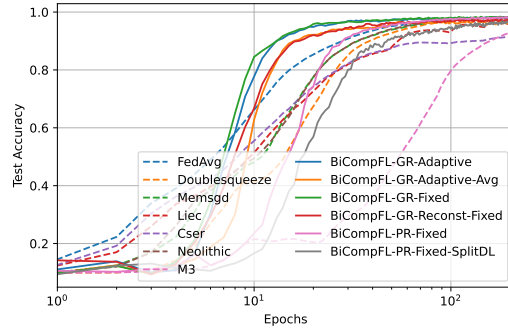
(c) MNIST 4CNN i.i.d.



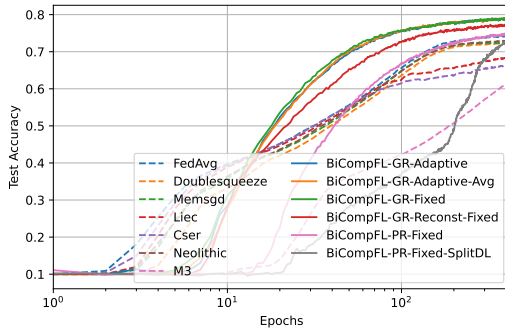
(d) MNIST 4CNN non-i.i.d.



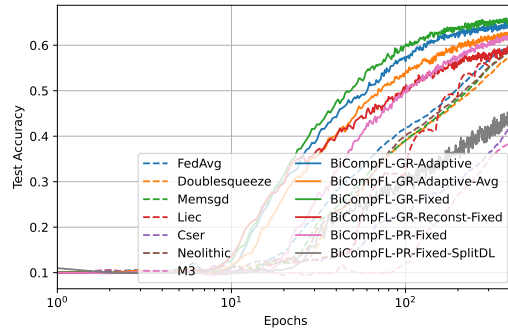
(e) Fashion MNIST 4CNN i.i.d.



(f) Fashion MNIST 4CNN non-i.i.d.



(g) Cifar-10 6CNN i.i.d.



(h) Cifar-10 6CNN non-i.i.d.

Figure 11: Test Accuracy over Epochs

## I Ablation Studies

### I.1 Number of Clients

We study in what follows the sensitivity to various hyperparameters of our algorithms. For comparability, we conduct all experiments on the model 4CNN, Fashion MNIST, and i.i.d.data. We plot for all experiments the accuracies over the number of epochs, and over the communication cost in bits. We first evaluate in Fig. 12 the effectiveness of BiCompFL-PR and BiCompFL-GR for different numbers of clients. It can be found that both algorithms exhibit satisfying performance even for  $n = 50$ , given that the same data is now distributed on more clients. The overall communication cost increases by roughly the factor of the increase in the number of  $n$ . To illustrate this further, we additionally plot in Fig. 13 the bitrates per parameter.

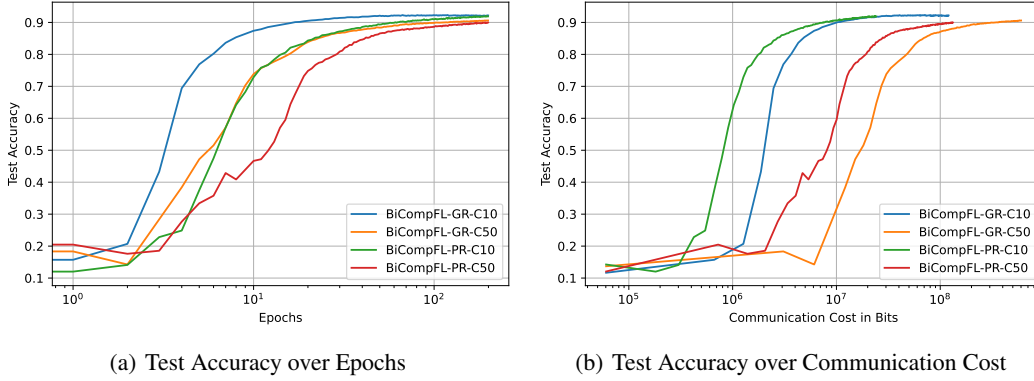


Figure 12: BiCompFL-GR and BiCompFL-GR With Different Number of Clients

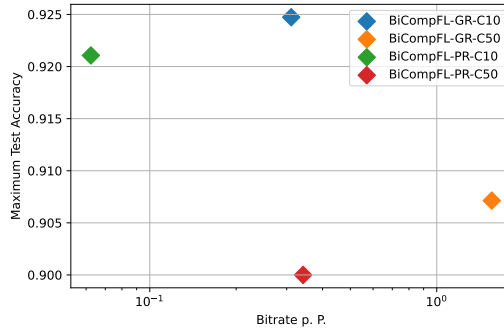


Figure 13: Bitrates for BiCompFL-GR and BiCompFL-GR With Different Number of Clients

### I.2 Optimization of the Prior

As described in the main body of the paper, BiCompFL-PR allows for optimizing the choice of the prior at the clients by optimizing the convexity parameter  $\lambda$  that mixes the global model estimate with the posterior transmitted by the client an iteration ahead, i.e.,  $p_{i,u}^t = \lambda \hat{\theta}_{i,t} + (1 - \lambda) \hat{q}_i^t$  to reduce the communication cost. To evaluate the potential of this method, we optimize  $\lambda$  so that it minimized the KL-divergence between the current posterior  $q_i^t$  (to be transmitted) and the prior  $p_{i,u}^t$ , representative for the uplink communication cost. The KL-minimizing  $\lambda$  is transmitted to the federator, which is necessary for the federator to reconstruct the importance samples. This optimization is conducted at each iteration individually at the clients. We present in Fig. 14 the performance of this method compared with the algorithms that use as priors exclusively the global model estimates of the clients. Note that optimizing the prior individually at the clients is only possible for BiCompFL-PR. We plot



the performance of BiCOMPFL-GR for reference only. To assess the potential, we ignore for the moment the cost of transmitting  $\lambda$ , which could be reduced by further compression techniques and leveraging the inter-round dependencies of the choice of  $\lambda$ .

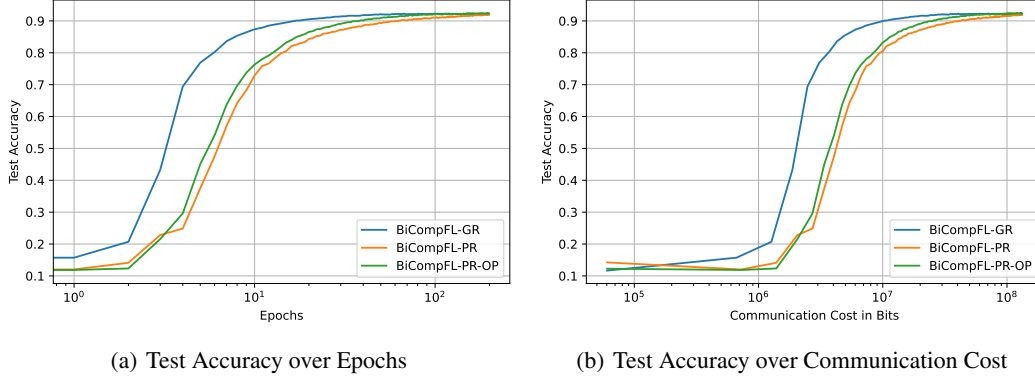


Figure 14: BiCOMPFL-PR With and Without Optimization over the Prior. Optimization over the Priors is denoted by OP.

It can be found that, while optimizing the prior improves the accuracy over epochs and with respect to the communication cost compared to BiCOMPFL-PR the improvements are rather insignificant. We therefore present for clarity the algorithm with a fixed choice of the prior as the former global model estimate, which additionally reduce the computation overhead at the clients by avoiding the optimization over  $\lambda$ . Nonetheless, we note that in certain edge cases, there can be merit in the optimization approach, for instance when the number  $n_{DL}$  of samples on the downlink is very small, and hence the global model estimate is inaccurate.

### I.3 Number of Samples

We continue to assess the impact of the number  $n_{DL}$  of samples on the downlink. We therefore evaluate the performance of BiCOMPFL-PR for  $n_{DL} \in \{5, 10, 20\}$ . We evaluate the differences on BiCOMPFL-PR. The results in Fig. 15 reflect the obvious: the larger  $n_{DL}$ , the better the accuracy when plotted over the number of epochs. On the contrary, the larger  $n_{DL}$ , the larger the communication cost per epoch. The final accuracies do not show substantial differences, and hence,  $n_{DL} = 5$  is sufficient in this setting. To avoid assessing our method overly optimistic and provide a fair comparison to other methods, we choose  $n_{DL} = 10$  in all our experiments, noting that the communication can further be reduced in certain scenarios by lowering  $n_{DL}$  without notable performance loss.

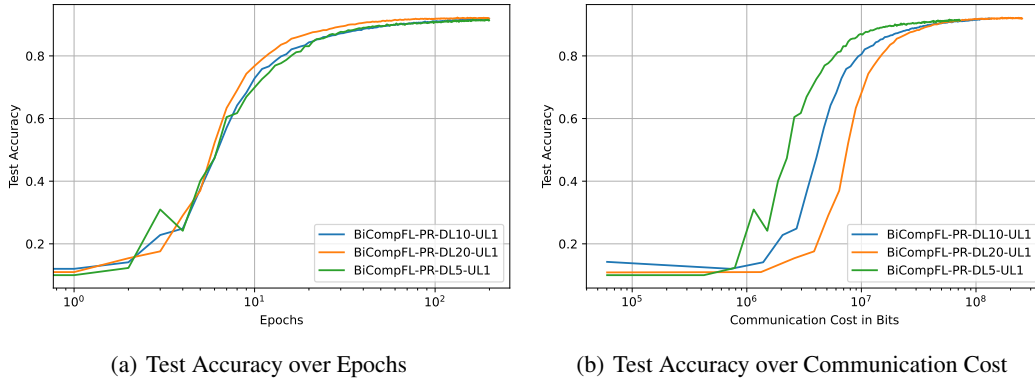


Figure 15: BiCOMPFL-PR for Different Number of Downlink Samples and a Single Uplink Sample.



#### I.4 Block Size

We compare in Fig. 16 the performance of BiCOMPFL-GR for different block sizes  $BS = d/B \in \{128, 256, 512\}$ . As expected, fixing  $n_{IS}$ , larger block sizes worsen the performance of the algorithm when evaluated over the number of epochs. However, larger block sizes simultaneously reduce the communication cost, and can hence be beneficial in many scenarios. However, we also note that larger block sizes comes at the expense of increases sampling complexities, and hence, the maximum block sizes are also dominated by the resources of the clients and the federator.

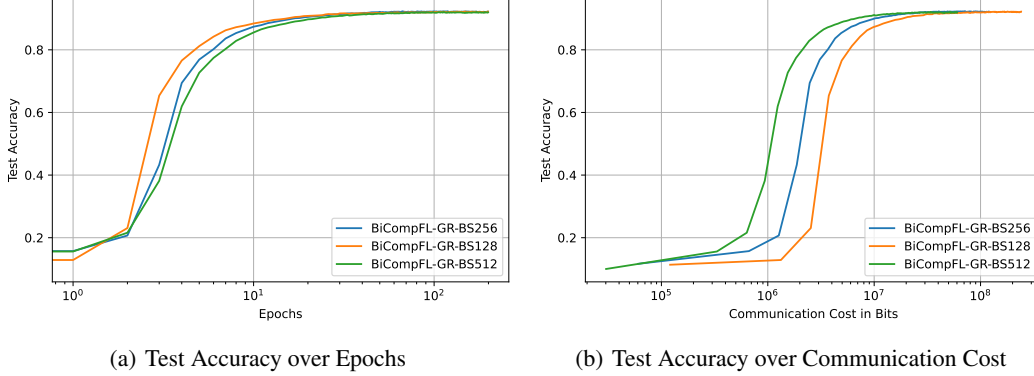


Figure 16: BiCOMPFL-GR With Fixed Block Allocation for Varying Block Sizes (BS)  $d/B$ .

#### I.5 Number of Importance Samples

In Fig. 17, we study the sensitivity of our algorithms with respect to the number of importance samples  $n_{IS}$  at the example of BiCOMPFL-GR. While larger number of  $n_{IS}$  slightly improves the performance as of the epoch number, the improvements do not outweigh the additional communication costs. Overall, our algorithm proves rather stable within reasonable ranges for  $n_{IS}$ . We fix in all our experiments  $n_{IS} = 256$ , presenting a good trade-off between performance and efficiency.

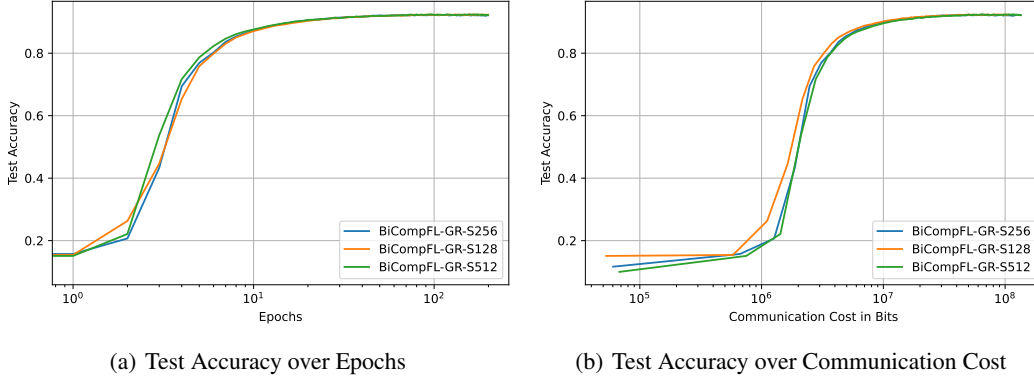


Figure 17: BiCOMPFL-GR with Varying Number of Importance Samples  $n_{IS}$  per Block.

#### I.6 Learning Rate

Our main claims are centered around the per-client bitrates per parameter, rendering the choices of learning rate secondary to our reasoning. Nonetheless, we tune the learning rates of all methods so that the baselines and BiCOMPFL achieve roughly the same final accuracies, allowing a fair comparison of resulting communication costs. We analyze the impact of the learning rate choice

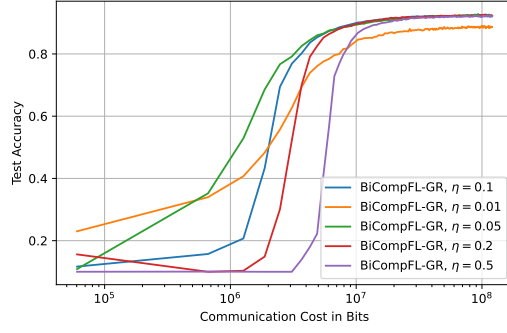


Figure 18: BiCompFL-GR with Varying Number of Importance Samples  $n_{IS}$  per Block.

on BiCompFL in Fig. 18, for  $\eta \in \{0.01, 0.05, 0.1, 0.2, 0.5\}$ . It is particularly noteworthy that BiCompFL exhibits stable performance across most learning rates we study, which we attribute to the regularization effects that occur in stochastic FL, detailed in the main body of the paper. Only for  $\eta = 0.01$ , the final performance is decreased, indicating that BiCompFL is not able to escape local optima in this setting. Although  $\eta = 0.05$  provides the best communication efficiency, we choose a moderate learning rate of  $\eta = 0.1$  not to overestimate our method compared to other approaches.

## NeurIPS Paper Checklist

### 1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [\[Yes\]](#)

Justification: All claims made in the abstract and the introduction are later justified in detail by the introduction of the schemes in Section 4 together with additional details in Appendices C and D, the various experiments in Section 4 and Appendix H, the formal theoretical guarantees (cf. Section 5), and the proofs provided in Appendix A.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

### 2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [\[Yes\]](#)

Justification: The theoretical and experimental limitations are discussed in the paper. For instance, we discuss in Sections 3 and 4 how the availability of shared randomness affects our protocols. Further, we discuss in Section 5 the assumptions required for our main results to hold. We provide in Appendix I details on the impact of various hyperparameters, including the limitations of optimizing the prior distributions, and how the block size is affected by the computational resources. The space constraints didn't allow for adding a separate section.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover

limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

### 3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [\[Yes\]](#)

Justification: All assumptions are provided along with the statements of our theoretical results. We provide the full set of proofs in Appendix A.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

### 4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [\[Yes\]](#)

Justification: The paper provides all experimental details to reproduce the results in Sections 3 and 4 and the corresponding Appendices D to F. The code to reproduce our experiments is included in the supplementary material. This includes parameter files for the various experiments and recommended environments.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
  - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
  - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.

- (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
- (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

## 5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: The source code required to reproduce our experiments is included in the supplementary material, including parameter files for the various experiments and instructions on how to use the provided code and how to prepare the publicly available datasets. We also provide the Python environment used to generate the results.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

## 6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: The paper specifies all training and test details necessary to understand the results in Section 4 and Appendix E.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

## 7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: Each experiment is repeated for multiple runs to provide statistical confidence. The depicted results are averaged performance measures, consistently shown together with their respective standard deviations.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

## 8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: The exact details of our simulation clusters are provided in Appendix E.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

## 9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

Answer: [Yes]

Justification: The paper conforms with the NeurIPS Code of Ethics. The research does not involve any human subjects or participants, and all datasets used are standard and available to the public.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.

- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

#### 10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [No]

Justification: This paper presents a method for bi-directional compression in stochastic (Bayesian) FL, with the aim of reducing communication overhead and improving scalability. The contribution is methodological and does not target a specific downstream application. While the work does not explicitly discuss societal impacts, such efficiency improvements may enable more practical and energy-efficient deployment of machine learning systems in resource-constrained environments, such as edge devices, or indirectly support privacy-preserving machine learning in settings such as medical healthcare. We do not foresee any immediate negative societal impacts arising from the work as presented.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

#### 11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: The paper poses no such risks.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.

- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

## 12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [\[Yes\]](#)

Justification: All assets are properly credited in the paper.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

## 13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [\[Yes\]](#)

Justification: Alongside the paper, we introduce new source code that we provide as supplementary material. This material is well documented with particular instructions on how to use the provided framework.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

## 14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [\[NA\]](#)

Justification: The paper does not involve crowdsourcing experiments nor research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.



- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

**15. Institutional review board (IRB) approvals or equivalent for research with human subjects**

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: The paper does not involve crowdsourcing experiments nor research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

**16. Declaration of LLM usage**

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigor, or originality of the research, declaration is not required.

Answer: [NA]

Justification: The core method development in this research does not involve LLMs as any important, original, or non-standard components.

Guidelines:

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>) for what should or should not be described.