

FEDSKETCH: COMMUNICATION-EFFICIENT FEDERATED LEARNING VIA SKETCHING

Anonymous authors

Paper under double-blind review

ABSTRACT

Communication complexity and data privacy are the two key challenges in Federated Learning (FL) where the goal is to perform a distributed learning through a large volume of devices. In this work, we introduce two new algorithms, namely FedSKETCH and FedSKETCHGATE, to address jointly both challenges and which are, respectively, intended to be used for homogeneous and heterogeneous data distribution settings. Our algorithms are based on a key and novel sketching technique, called HEAPRIX that is unbiased, compresses the accumulation of local gradients using count sketch, and exhibits communication-efficiency properties leveraging low-dimensional sketches. We provide sharp convergence guarantees of our algorithms and validate our theoretical findings with various sets of experiments.

1 INTRODUCTION

Federated Learning (FL) is a recently emerging framework for distributed large scale machine learning problems. In FL, data is distributed across devices (Konečný et al., 2016; McMahan et al., 2017) and due to privacy concerns, users are only allowed to communicate with the parameter server. Formally, the optimization problem across p distributed devices is defined as follows:

$$\min_{\mathbf{x} \in \mathbb{R}^d, \sum_{j=1}^p q_j = 1} f(\mathbf{x}) \triangleq \sum_{j=1}^p q_j F_j(\mathbf{x}), \quad (1)$$

where $F_j(\mathbf{x}) = \mathbb{E}_{\xi \in \mathcal{D}_j} [L_j(\mathbf{x}, \xi)]$ is the local cost function at device j , $q_j \triangleq \frac{n_j}{n}$, n_j is the number of data shards at device j and $n = \sum_{j=1}^p n_j$ is the total number of data samples, ξ is a random variable distributed according to probability distribution \mathcal{D}_j , and L_j is a loss function that measures the performance of model \mathbf{x} at device j . We note that, while for the homogeneous setting we assume $\{\mathcal{D}_j\}_{j=1}^p$ have the same distribution across devices and $L_i = L_j$, $1 \leq (i, j) \leq p$, in the heterogeneous setting, these distributions and loss functions L_j can vary from a device to another.

There are several challenges that need to be addressed in FL in order to efficiently learn a global model that performs well in average for all devices:

– *Communication-efficiency*: There are often many devices communicating with the server, thus incurring immense communication overhead. One approach to reduce communication round is using *local SGD with periodic averaging* (Zhou and Cong, 2018; Stich, 2019; Yu et al., 2019b; Wang and Joshi, 2018) which periodically averages models after a few local updates, contrary to baseline SGD (Bottou and Bousquet, 2008) where gradient averaging is performed at each iteration. Local SGD has been proposed in McMahan et al. (2017); Konečný et al. (2016) under the FL setting and its convergence analysis is studied in Stich (2019); Wang and Joshi (2018); Zhou and Cong (2018); Yu et al. (2019b), later on improved in the followup references (Basu et al., 2019; Haddadpour and Mahdavi, 2019; Khaled et al., 2020; Stich and Karimireddy, 2019) for homogeneous setting. It is further extended to heterogeneous setting (Haddadpour and Mahdavi, 2019; Karimireddy et al., 2019; Yu et al., 2019a; Li et al., 2020c; Sahu et al., 2018; Liang et al., 2019). The second approach to deal with communication cost aims at reducing the size of communicated message per communication round, such as local gradient quantization (Alistarh et al., 2017; Bernstein et al., 2018; Tang et al., 2018; Wen et al., 2017; Wu et al., 2018) or sparsification (Alistarh et al., 2018; Lin et al., 2018; Stich et al., 2018; Stich and Karimireddy, 2019).

– *Data heterogeneity*: Since locally generated data in each device may come from different distribution, local computations involved in FL setting can lead to poor convergence error in practice (Li et al.,

2020a; Liang et al., 2019). To mitigate the negative impact of data heterogeneity, (Haddadpour et al., 2020; Horváth et al., 2019; Liang et al., 2019; Karimireddy et al., 2019) suggest applying variance reduction or gradient tracking techniques along local computations.

–*Privacy* (Geyer et al., 2017; Hardy et al., 2017): Privacy has been widely addressed by injecting an additional layer of randomness to respect differential-privacy property (McMahan et al., 2018) or using cryptography-based approaches under secure multi-party computation (Bonawitz et al., 2017). Further study of challenges can be found in recent surveys Li et al. (2020a) and Kairouz et al. (2019).

To tackle the aforementioned challenges in FL jointly, sketching based algorithms (Charikar et al., 2004; Cormode and Muthukrishnan, 2005; Kleinberg, 2003; Li et al., 2008) are promising approaches. For instance, to reduce communication cost, (Ivkin et al., 2019) develops a distributed SGD algorithm using sketching along providing its convergence analysis in the homogeneous setting, and establish a communication complexity of order $\mathcal{O}(\log(d))$ per round, where d is the dimension of the vector of parameters compared to $\mathcal{O}(d)$ complexity per round of baseline mini-batch SGD. Yet, the proposed sketching scheme in Ivkin et al. (2019), built from a communication-efficiency perspective, is based on a deterministic procedure which requires access to the exact information of the gradients, thus not meeting the privacy-preserving criteria. This systemic issue is partially addressed in Rothchild et al. (2020).

Focusing on privacy, (Li et al., 2019) derives a single framework in order to tackle these issues jointly and introduces `DiffSketch` algorithm, based on the Count Sketch operator, yet does not provide its convergence analysis. Additionally, the estimation error of `DiffSketch` is higher than the sketching scheme in Ivkin et al. (2019) which may end up in poor convergence.

Our main contributions are summarized as follows:

- We provide a new algorithm – `HEAPRIX` – and theoretically show that it reduces the cost of communication between devices and server, based on unbiased sketching without requiring the broadcast of exact values of gradients to the server. Based on `HEAPRIX`, we develop general algorithms for communication-efficient and sketch-based FL, namely `FedSKETCH` and `FedSKETCHGATE` for homogeneous and heterogeneous data distribution settings respectively.
- We establish non-asymptotic convergence bounds for convex, Polyak-Łojasiewicz (PL) and non-convex functions in Theorems 1 and 2 in both homogeneous and heterogeneous cases, and highlight an improvement in the number of iteration to reach a stationary point. We also provide a convergence analysis for the `PRIVIX/DiffSketch`¹ algorithm proposed in Li et al. (2019).
- We illustrate the benefits of `FedSKETCH` and `FedSKETCHGATE` over baseline methods through a set of experiments. The latter shows the advantages of the `HEAPRIX` compression method achieving comparable test accuracy as Federated SGD (`FedSGD`) while compressing the information exchanged between devices and server.

Notation: We denote the number of communication rounds and bits per round and per device by R and B respectively. The count sketch of vector \mathbf{x} is designated by $\mathbf{S}(\mathbf{x})$. $[p]$ denotes the set $\{1, \dots, p\}$.

2 COMPRESSION USING COUNT SKETCH

In this paper, we exploit the commonly used `Count Sketch` (Charikar et al., 2004) which uses two sets of functions that encode any input vector \mathbf{x} into a **hash table** $\mathbf{S}_{m \times t}(\mathbf{x})$. Pairwise independent hash functions $\{h_{j,1 \leq j \leq t} : [d] \rightarrow m\}$ are used along with another set of pairwise independent sign hash functions $\{\text{sign}_{j,1 \leq j \leq t} : [d] \rightarrow \{+1, -1\}\}$ to map entries of \mathbf{x} (x_i , $1 \leq i \leq d$) into t different columns of $\mathbf{S}_{m \times t}$, wherein to lower the dimension of the input vector we usually have $d \gg mt$. The final update reads $\mathbf{S}[j][h_j(i)] = \mathbf{S}[j][h_j(i)] + \text{sign}_j(i)x_i$ for any $1 \leq j \leq t$. There are various types of sketching algorithms which are developed based on count sketching that we develop in the following subsections. See the Appendix for the detailed Count Sketch algorithm.

2.1 SKETCHING BASED UNBIASED COMPRESSOR

We define an unbiased compressor as follows:

¹We use `PRIVIX` and `DiffSketch` Li et al. (2019) interchangeably throughout the paper.

Definition 1 (Unbiased compressor). We call randomized function, $C : \mathbb{R}^d \rightarrow \mathbb{R}^d$ an unbiased compression operator with $\Delta \geq 1$, if

$$\mathbb{E}[C(\mathbf{x})] = \mathbf{x} \quad \text{and} \quad \mathbb{E}[\|C(\mathbf{x})\|_2^2] \leq \Delta \|\mathbf{x}\|_2^2.$$

We denote this class of compressors by $\mathbb{U}(\Delta)$.

This definition leads to the following property

$$\mathbb{E}[\|C(\mathbf{x}) - \mathbf{x}\|_2^2] \leq (\Delta - 1) \|\mathbf{x}\|_2^2.$$

Note that if we let $\Delta = 1$ then our algorithm reduces to the case of no compression. This property allows us to control the noise of the compression.

An instance of such unbiased compressor is PRIVIX which obtains an estimate of input \mathbf{x} from a count sketch noted $\mathbf{S}(\mathbf{x})$. In this algorithm, to query the quantity x_i , the i -th element of the vector \mathbf{x} , we compute the median of t approximated values specified by the indices of $h_j(i)$ for $1 \leq j \leq t$, see (Li et al., 2019), or Algorithm 6 in the Appendix (for more details). The following property of count sketch would be useful for our theoretical analysis.

Property 1 (Li et al. (2019)). For any $\mathbf{x} \in \mathbb{R}^d$, we have:

Unbiased estimation: As in Li et al. (2019), we have $\mathbb{E}_{\mathbf{S}}[\text{PRIVIX}[\mathbf{S}(\mathbf{x})]] = \mathbf{x}$.

Bounded variance: For the given $m < d$, $t = \mathcal{O}(\ln(\frac{d}{\delta}))$ with probability $1 - \delta$ we have:

$$\mathbb{E}_{\mathbf{S}}[\|\text{PRIVIX}[\mathbf{S}(\mathbf{x})] - \mathbf{x}\|_2^2] \leq \frac{c \times d}{m} \|\mathbf{x}\|_2^2,$$

where c ($e \leq c < m$) is a positive constant independent of the dimension of the input, d .

We note that bounded variance assumption does not necessary implies any compression as d could be relatively large. Thus, with probability $1 - \delta$ we obtain $\text{PRIVIX} \in \mathbb{U}(1 + c\frac{d}{m})$. $\Delta = 1 + c\frac{d}{m}$ implies that if $m \rightarrow d$, then $\Delta \rightarrow 1 + c$, indicating a noisy reconstruction. The reference Li et al. (2019) shows that if the data is normally distributed, PRIVIX is differentially private (Dwork, 2006), up to additional assumptions and algorithmic design.

2.2 SKETCHING BASED BIASED COMPRESSOR

A biased compressor is defined as follows:

Definition 2 (Biased compressor). A (randomized) function, $C : \mathbb{R}^d \rightarrow \mathbb{R}^d$ belongs to $\mathbb{C}(\Delta, \alpha)$, a class of compression operators with $\alpha > 0$ and $\Delta \geq 1$, if

$$\mathbb{E}[\|\alpha \mathbf{x} - C(\mathbf{x})\|_2^2] \leq \left(1 - \frac{1}{\Delta}\right) \|\mathbf{x}\|_2^2,$$

The reference (Horváth and Richtárik, 2020) proves that $\mathbb{U}(\Delta) \subset \mathbb{C}(\Delta, \alpha)$. An example of biased compression via sketching and using top_m operation is given below:

Following Ivkin et al. (2019), HEAVYMIX with sketch size $\Theta(m \log(\frac{d}{\delta}))$ is a biased compressor with $\alpha = 1$ and $\Delta = d/m$ with probability $\geq 1 - \delta$, meaning that it reconstruct the $\tilde{\mathbf{g}}$ from input vector \mathbf{g} . In other words, with probability $1 - \delta$, $\text{HEAVYMIX} \in \mathbb{C}(\frac{d}{m}, 1)$. We note that Algorithm 1 is a variation of the sketching algorithm developed in Ivkin et al. (2019) with distinction that HEAVYMIX does not require a second round of communication to obtain the exact values of top_m . This is mainly because in SKETCHED-SGD Ivkin et al. (2019) the server has to obtain the exact values of the average of sketches; however HEAVYMIX obtains exact value

Algorithm 1 HEAVYMIX

- 1: **Inputs:** $\mathbf{S}(\mathbf{g})$; parameter m
 - 2: Query the vector $\tilde{\mathbf{g}} \in \mathbb{R}^d$ from $\mathbf{S}(\mathbf{g})$:
 - 3: Query $\ell_2^2 = (1 \pm 0.5) \|\mathbf{g}\|^2$ from sketch $\mathbf{S}(\mathbf{g})$
 - 4: $\forall j$ query $\hat{\mathbf{g}}_j^2 = \tilde{\mathbf{g}}_j^2 \pm \frac{1}{2m} \|\mathbf{g}\|^2$ from sketch $\mathbf{S}(\mathbf{g})$
 - 5: $H = \{j | \hat{\mathbf{g}}_j \geq \frac{\ell_2^2}{m}\}$ and $NH = \{j | \hat{\mathbf{g}}_j < \frac{\ell_2^2}{m}\}$
 - 6: $\text{Top}_m = H \cup \text{rand}_{\ell}(NH)$, where $\ell = m - |H|$
 - 7: Get exact values of Top_m
 - 8: **Output:** $\tilde{\mathbf{g}} : \forall j \in \text{Top}_m : \tilde{\mathbf{g}}_j = \mathbf{g}_j$ else $\mathbf{g}_i = 0$
-

locally, thus does not require a second round of communication. Additionally, while a sketching algorithm implementing HEAVYMIX has smaller estimation error compared to PRIVIX, it requires having access to the exact values of top_m , therefore not benefiting from privacy properties contrary to PRIVIX. In the following we introduce HEAPRIX which is built upon HEAVYMIX and PRIVIX methods.

2.3 SKETCHING BASED INDUCED COMPRESSOR

Due to Theorem 3 in Horváth and Richtárik (2020), which illustrates that we can convert the biased compressor into an unbiased one such that, for $C_1 \in \mathbb{C}(\Delta_1)$ with $\alpha = 1$, if you choose $C_2 \in \mathbb{U}(\Delta_2)$, then induced compressor $C : x \mapsto C_1(x) + C_2(x - C_1(x))$ belongs to $\mathbb{U}(\Delta)$ with $\Delta = \Delta_2 + \frac{1-\Delta_2}{\Delta_1}$.

Based on this notion, Algorithm 2 proposes an induced sketching algorithm by utilizing HEAVYMIX and PRIVIX for C_1 and C_2 respectively where the reconstruction of input x is performed using hash table \mathbf{S} and \mathbf{x} , similar to PRIVIX and HEAVYMIX. Note that if $m \rightarrow d$, then $C(x) \rightarrow x$, implying that the convergence rate can be improved by decreasing the size of compression m .

Algorithm 2 HEAPRIX

- 1: **Inputs:** $x \in \mathbb{R}^d, t, m, \mathbf{S}_{m \times t}, h_j(1 \leq i \leq t), \text{sign}_j(1 \leq i \leq t)$, parameter m
 - 2: Approximate $\mathbf{S}(x)$ using HEAVYMIX
 - 3: Approximate $\mathbf{S}(x - \text{HEAVYMIX}[\mathbf{S}(x)])$ with PRIVIX
 - 4: **Output:**
 $\text{HEAVYMIX}[\mathbf{S}(x)] + \text{PRIVIX}[\mathbf{S}(x - \text{HEAVYMIX}[\mathbf{S}(x)])]$.
-

Corollary 1. *Based on Theorem 3 of (Horváth and Richtárik, 2020), HEAPRIX in Algorithm 2 satisfies $C(x) \in \mathbb{U}(c \frac{d}{m})$.*

Benefits of HEAPRIX: Corollary 1 states that, unlike PRIVIX, HEAPRIX compression noise can be made as small as possible using larger hash size. In the distributed setting, contrary to SKETCHED-SGD Ivkin et al. (2019) where decompressing is happening at the server, HEAPRIX does not require having access to exact top_m values of the input as it is based on HEAVYMIX, which helps preserving privacy. In other words, HEAPRIX leverages the best of both: the *unbiasedness* of PRIVIX while using *heavy hitters* as in HEAVYMIX.

3 FEDSKETCH AND FEDSKETCHGATE

We introduce two new algorithms for both homogeneous and heterogeneous settings.

3.1 HOMOGENEOUS SETTING

In FedSKETCH, the number of local updates, between two consecutive communication rounds, at device j is denoted by τ . Unlike Haddadpour et al. (2020), server node does not store any global model, rather, device j has two models: $x^{(r)}$ and $x_j^{(\ell, r)}$, which are respectively the local and global models. We develop FedSKETCH in Algorithm 3. A variant of this algorithm implementing HEAPRIX is also described in Algorithm 3. We remark that for this variant, we need to have an additional communication round between server and worker j to aggregate $\delta_j^{(r)} \triangleq \mathbf{S}_j[\text{HEAVYMIX}(\mathbf{S}^{(r)})]$ (Lines 3 and 3) to compute $\mathbf{S}^{(r)} = \frac{1}{k} \sum_{j \in \mathcal{K}} \mathbf{S}_j^{(r)}$. The main difference between FedSKETCH and DiffSketch in Li et al. (2019) is that we use distinct local and global learning rates. Furthermore, unlike Li et al. (2019), we do not add local Gaussian noise.

Algorithm 3 FedSKETCH(R, τ, η, γ)

- 1: **Inputs:** $x^{(0)}$: initial model shared by local devices, global and local learning rates γ and η , respectively
- 2: **for** $r = 0, \dots, R - 1$ **do**
- 3: **parallel for device** $j \in \mathcal{K}^{(r)}$ **do:**
- 4: **if PRIVIX variant:**
 $\Phi^{(r)} \triangleq \text{PRIVIX}[\mathbf{S}^{(r-1)}]$
- 5: **if HEAPRIX variant:**
 $\Phi^{(r)} \triangleq \text{HEAVYMIX}[\mathbf{S}^{(r-1)}] + \text{PRIVIX}[\mathbf{S}^{(r-1)} - \tilde{\mathbf{S}}^{(r-1)}]$
- 6: Set $x^{(r)} = x^{(r-1)} - \gamma \Phi^{(r)}$ and $x_j^{(0, r)} = x^{(r)}$
- 7: **for** $\ell = 0, \dots, \tau - 1$ **do**
- 8: Sample a mini-batch $\xi_j^{(\ell, r)}$ and compute $\tilde{g}_j^{(\ell, r)}$
- 9: Update $x_j^{(\ell+1, r)} = x_j^{(\ell, r)} - \eta \tilde{g}_j^{(\ell, r)}$
- 10: **end for**
- 11: Device j broadcasts $\mathbf{S}_j^{(r)} \triangleq \mathbf{S}_j(x_j^{(0, r)} - x_j^{(\tau, r)})$.
- 12: Server computes $\mathbf{S}^{(r)} = \frac{1}{k} \sum_{j \in \mathcal{K}} \mathbf{S}_j^{(r)}$.
- 13: Server broadcasts $\mathbf{S}^{(r)}$ to devices in randomly drawn devices $\mathcal{K}^{(r)}$.
- 14: **if HEAPRIX variant:**
- 15: Second round of communication: $\delta_j^{(r)} \triangleq \mathbf{S}_j[\text{HEAVYMIX}(\mathbf{S}^{(r)})]$ and broadcasts $\tilde{\mathbf{S}}^{(r)} \triangleq \frac{1}{k} \sum_{j \in \mathcal{K}} \delta_j^{(r)}$ to devices in set $\mathcal{K}^{(r)}$
- 16: **end parallel for**
- 17: **end**
- 18: **Output:** $x^{(R-1)}$

Algorithmic comparison with Haddad-

pour et al. (2020) An important feature of our algorithm is that due to a lower dimension of the count sketch, the resulting averages ($\mathbf{S}^{(r)}$ and $\tilde{\mathbf{S}}^{(r)}$) received by the server are also of lower dimension. Therefore, these algorithms exploit a bidirectional compression during the communication from server to device back and forth. As a result, for the case of large quantization error $\omega = \theta(\frac{d}{m})$ as shown in Haddadpour et al. (2020), our algorithms can outperform FedCOM and FedCOMGATE developed in Haddadpour et al. (2020) if sufficiently large hash tables are used and the uplink communication cost is high. Furthermore, while, in Haddadpour et al. (2020), server stores a global model and aggregates the partial gradients from devices which can enable the server to extract some information regarding the device's data, in contrast, in our algorithms server does not store the global model and only broadcasts the average sketches. Thus, sketching-based server-devices communication algorithms such as ours do not reveal the exact values of the inputs, to preserve privacy as a by-product.

Remark 1. As pointed out in Horváth and Richtárik (2020), while induced compressors transform a biased compressor into unbiased one, as a drawback it doubles communication cost since the devices need to send $C_1(\mathbf{x})$ and $C_2(\mathbf{x} - C_1(\mathbf{x}))$ separately. We note that in the special case of HEAPRIX, due to the use of sketching, the extra communication round cost is compensated with lower number of bits per round thanks to the lower dimension of sketching.

3.2 HETEROGENEOUS SETTING

In this section, we focus on the optimization problem of (1) in the special case of $q_1 = \dots = q_p = \frac{1}{p}$ with full device participation ($k = p$). These results can be extended to the scenario where devices are sampled. For non i.i.d. data, the FedSKETCH algorithm, designed for homogeneous setting, may fail to perform well in practice. The main reason is that in FL, devices are using local stochastic descent direction which could be different than global descent direction when the data distribution are non-identical. Therefore, to mitigate the effect of data heterogeneity, we introduce a new algorithm called FedSKETCHGATE described in Algorithm 4. This algorithm leverages the idea of gradient tracking applied in Haddadpour et al. (2020) (with compression) and a special case of $\gamma = 1$ without compression (Liang et al., 2019). The main idea is that using an approximation of global gradient, $\mathbf{c}_j^{(r)}$ allows to correct the local gradient direction. For the FedSKETCHGATE with PRIVIX variant, the correction vector $\mathbf{c}_j^{(r)}$ at device j and communication round r is computed in Line 4. While using HEAPRIX compression, FedSKETCHGATE also updates $\tilde{\mathbf{S}}^{(r)}$ via Line 4.

Remark 2. Most of the existing communication-efficient algorithms with compression only consider communication-efficiency from devices to server. However,

Algorithm 4 FedSKETCHGATE(R, τ, η, γ)

- 1: **Inputs:** $\mathbf{x}^{(0)} = \mathbf{x}_j^{(0)}$ shared by all local devices, global and local learning rates γ and η .
- 2: **for** $r = 0, \dots, R - 1$ **do**
- 3: **parallel for device** $j = 1, \dots, p$ **do:**
- 4: **if PRIVIX variant:**

$$\mathbf{c}_j^{(r)} = \mathbf{c}_j^{(r-1)} - \frac{1}{\tau} [\text{PRIVIX}(\mathbf{S}^{(r-1)}) - \text{PRIVIX}(\mathbf{S}_j^{(r-1)})]$$

where $\Phi^{(r)} \triangleq \text{PRIVIX}(\mathbf{S}^{(r-1)})$
- 5: **if HEAPRIX variant:**

$$\mathbf{c}_j^{(r)} = \mathbf{c}_j^{(r-1)} - \frac{1}{\tau} (\Phi^{(r)} - \Phi_j^{(r)})$$
- 6: Set $\mathbf{x}^{(r)} = \mathbf{x}^{(r-1)} - \gamma \Phi^{(r)}$ and $\mathbf{x}_j^{(0,r)} = \mathbf{x}^{(r)}$
- 7: **for** $\ell = 0, \dots, \tau - 1$ **do**
- 8: Sample mini-batch $\xi_j^{(\ell,r)}$ and compute $\tilde{\mathbf{g}}_j^{(\ell,r)}$
- 9: $\mathbf{x}_j^{(\ell+1,r)} = \mathbf{x}_j^{(\ell,r)} - \eta (\tilde{\mathbf{g}}_j^{(\ell,r)} - \mathbf{c}_j^{(r)})$
- 10: **end for**
- 11: Device j broadcasts $\mathbf{S}_j^{(r)} \triangleq \mathbf{S}(\mathbf{x}_j^{(0,r)} - \mathbf{x}_j^{(\tau,r)})$.
- 12: Server computes $\mathbf{S}^{(r)} = \frac{1}{p} \sum_{j=1}^p \mathbf{S}_j^{(r)}$ and **broadcasts** $\mathbf{S}^{(r)}$ to all devices.
- 13: **if HEAPRIX variant:**
- 14: Device j computes $\Phi_j^{(r)} \triangleq \text{HEAPRIX}[\mathbf{S}_j^{(r)}]$
- 15: Second round of communication to obtain $\delta_j^{(r)} := \mathbf{S}_j(\text{HEAVYMIX}[\mathbf{S}^{(r)}])$
- 16: Broadcasts $\tilde{\mathbf{S}}^{(r)} \triangleq \frac{1}{p} \sum_{j=1}^p \delta_j^{(r)}$ to devices
- 17: **end parallel for**
- 18: **end**
- 19: **Output:** $\mathbf{x}^{(R-1)}$

Algorithms 3 and 4 also improve the communication efficiency from server to devices since it exploits low-dimensional sketches (and averages), communicated from the server to devices.

For both FedSKETCH and FedSKETCHGATE algorithms, unlike PRIVIX, HEAPRIX variant requires a second round of communication. Therefore, in Cross-Device FL setting, where there could be millions of devices, HEAPRIX variant may not be practical, and we note that it could be more suitable for Cross-Silo FL setting.

4 CONVERGENCE ANALYSIS

We first state commonly used assumptions required in the following convergence analysis (reminder of our notations can be found Table 1 of the Appendix).

Assumption 1 (Smoothness and Lower Boundedness). *The local objective function $f_j(\cdot)$ of device j is differentiable for $j \in [p]$ and L -smooth, i.e., $\|\nabla f_j(\mathbf{x}) - \nabla f_j(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|$, $\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d$. Moreover, the optimal objective function $f(\cdot)$ is bounded below by $f^* := \min_{\mathbf{x}} f(\mathbf{x}) > -\infty$.*

Assumption 1 is common in stochastic optimization. We present our results for PL, convex and general non-convex objectives. (Karimi et al., 2016) show that PL condition implies strong convexity property with same module (PL objectives can also be non-convex, hence strong convexity does not imply PL condition necessarily).

4.1 CONVERGENCE OF FEDSKETCH

We now focus on the homogeneous case where data is i.i.d. among local devices, and therefore, the stochastic local gradient of each worker is an unbiased estimator of the global gradient. We have:

Assumption 2 (Bounded Variance). *For all $j \in [m]$, we can sample an independent mini-batch ℓ_j of size $|\Xi_j^{(\ell, r)}| = b$ and compute an unbiased stochastic gradient $\tilde{\mathbf{g}}_j = \nabla f_j(\mathbf{x}; \Xi_j)$, $\mathbb{E}_{\Xi_j}[\tilde{\mathbf{g}}_j] = \nabla f(\mathbf{x}) = \mathbf{g}$ with the variance bounded is bounded by a constant σ^2 , i.e., $\mathbb{E}_{\Xi_j}[\|\tilde{\mathbf{g}}_j - \mathbf{g}\|^2] \leq \sigma^2$.*

Theorem 1. *Suppose Assumptions 1-2 hold. Given $0 < m \leq d$ and considering Algorithm 3 with sketch size $B = O(m \log(\frac{dR}{\delta}))$ and $\gamma \geq k$, with probability $1 - \delta$ we have:*

*In the **non-convex** case, $\{\mathbf{x}^{(r)}\}_{r=0}^R$ satisfies $\frac{1}{R} \sum_{r=0}^{R-1} \mathbb{E}[\|\nabla f(\mathbf{x}^{(r)})\|_2^2] \leq \epsilon$ if:*

- *FS-PRIVIX, for $\eta = \frac{1}{L\gamma} \sqrt{\frac{k}{R\tau(\frac{cd}{mk}+1)}}$: $R = O(1/\epsilon)$ and $\tau = O((d+m)/(mk\epsilon))$.*
- *FS-HEAPRIX, for $\eta = \frac{1}{L\gamma} \sqrt{\frac{k}{R\tau(\frac{cd-m}{mk}+1)}}$: $R = O(1/\epsilon)$ and $\tau = O(d/(mk\epsilon))$.*

*In the **PL or strongly convex** case, $\{\mathbf{x}^{(r)}\}_{r=0}^R$ satisfies $\mathbb{E}[f(\mathbf{x}^{(R-1)}) - f(\mathbf{x}^{(*)})] \leq \epsilon$ if we set:*

- *FS-PRIVIX, for $\eta = \frac{1}{2L(cd/mk+1)\tau\gamma}$: $R = O((d/mk+1)\kappa \log(1/\epsilon))$ and $\tau = O((d/m+1)/(d/m+k)\epsilon)$.*
- *FS-HEAPRIX, for $\eta = \frac{1}{2L((cd-m)/mk+1)\tau\gamma}$: $R = O(((d-m)/mk+1)\kappa \log(1/\epsilon))$ and $\tau = O(d/m/(((d-m-1)+k)\epsilon))$.*

*In the **Convex** case, $\{\mathbf{x}^{(r)}\}_{r=0}^R$ satisfies $\mathbb{E}[f(\mathbf{x}^{(R-1)}) - f(\mathbf{x}^{(*)})] \leq \epsilon$ if we set:*

- *FS-PRIVIX, for $\eta = \frac{1}{2L(cd/mk+1)\tau\gamma}$: $R = O(L(1+d/mk)/\epsilon \log(1/\epsilon))$ and $\tau = O((d/m+1)^2/(k(d/mk+1)^2\epsilon^2))$.*
- *FS-HEAPRIX, for $\eta = \frac{1}{2L((cd-m)/mk+1)\tau\gamma}$: $R = O(L(1+(d-m)/mk)/\epsilon \log(1/\epsilon))$ and $\tau = O((d/m)^2/(k([d-m]/mk+1)^2\epsilon^2))$.*

The bounds in Theorem 1 suggest that in homogeneous setting if we set $d = m$ (no compression), the number of communication rounds to achieve the ϵ error matches with the number of iterations

required to achieve the same error under a centralized setting. Additionally, computational complexity scales down with number of sampled devices. To stress on the further impact of using sketching, we also compare our results with prior works in terms of total number of communicated bits per device.

Comparison with Ivkin et al. (2019) From privacy aspect, we note Ivkin et al. (2019) requires for server to have access to exact values of top_m gradients, hence do not preserve privacy, whereas our schemes do not need those exact values. From communication cost point of view, for strongly convex objective and compared to Ivkin et al. (2019), we improve the total communication per worker from $RB = O\left(\frac{d}{\epsilon} \log\left(\frac{d}{\delta\sqrt{\epsilon}} \max\left(\frac{d}{m}, \frac{1}{\sqrt{\epsilon}}\right)\right)\right)$ to

$$RB = O\left(\kappa\left(\frac{d-m}{k} + m\right) \log \frac{1}{\epsilon} \log\left(\frac{\kappa d}{\delta}\left(\frac{d-m}{mk} + 1\right) \log \frac{1}{\epsilon}\right)\right).$$

We note that while reducing communication cost, our scheme requires $\tau = O(d/m(k(\frac{d}{mk} + 1)\epsilon)) > 1$, which scales down with the number of sampled devices, k . Moreover, unlike Ivkin et al. (2019), we do not use bounded gradient assumption. Therefore, we obtain stronger result with weaker assumptions. Regarding general non-convex objectives, our result improves the total communication cost per worker in Ivkin et al. (2019) from $RB = O\left(\max(\frac{1}{\epsilon^2}, \frac{d^2}{k^2\epsilon}) \log(\frac{d}{\delta} \max(\frac{1}{\epsilon^2}, \frac{d^2}{k^2\epsilon}))\right)$ for *only one device* to $RB = O(\frac{m}{\epsilon} \log(\frac{d}{\epsilon\delta}))$. We also highlight that we can obtain similar rates for Algorithm 3 in heterogeneous environment if we make the additional assumption of uniformly bounded gradient.

Note: Such improved communication cost over prior related works is due to joint exploitation of *sketching*, to reduce the dimension of communicated messages, and the use of *local updates*, to reduce the total number of communication rounds leading to a specific convergence error.

4.2 CONVERGENCE OF FEDSKETCHGATE

We start with bounded local variance assumption:

Assumption 3 (Bounded Local Variance). *For all $j \in [p]$, we can sample an independent mini-batch Ξ_j of size $|\xi_j| = b$ and compute an unbiased stochastic gradient $\tilde{\mathbf{g}}_j = \nabla f_j(\mathbf{x}; \Xi_j)$ with $\mathbb{E}_{\xi}[\tilde{\mathbf{g}}_j] = \nabla f_j(\mathbf{x}) = \mathbf{g}_j$. Moreover, the variance of local stochastic gradients is bounded such that $\mathbb{E}_{\Xi}[\|\tilde{\mathbf{g}}_j - \mathbf{g}_j\|^2] \leq \sigma^2$.*

Theorem 2. *Suppose Assumptions 1 and 3 hold. Given $0 < m \leq d$, and considering FedSKETCHGATE in Algorithm 4 with sketch size $B = O(m \log(\frac{dR}{\delta}))$ and $\gamma \geq p$ with probability $1 - \delta$ we have*

*In the **non-convex** case, $\eta = \frac{1}{L\gamma} \sqrt{\frac{mp}{R\tau(cd)}}$, $\{\mathbf{x}^{(r)}\}_{r=0}^{\infty}$ satisfies $\frac{1}{R} \sum_{r=0}^{R-1} \mathbb{E}[\|\nabla f(\mathbf{x}^{(r)})\|_2^2] \leq \epsilon$ if:*

- **FS-PRIVIX:**

$$R = O((d+m)/m\epsilon) \quad \text{and} \quad \tau = O(1/(p\epsilon)).$$

- **FS-HEAPRIX:** $R = O(d/m\epsilon)$ and $\tau = O(1/(p\epsilon))$.

*In the **PL or Strongly convex** case, $\{\mathbf{x}^{(r)}\}_{r=0}^{\infty}$ satisfies $\mathbb{E}[f(\mathbf{x}^{(R-1)}) - f(\mathbf{x}^{(*)})] \leq \epsilon$ if:*

- **FS-PRIVIX**, for $\eta = 1/(2L(\frac{cd}{m} + 1)\tau\gamma)$: $R = O((\frac{d}{m} + 1)\kappa \log(1/\epsilon))$ and $\tau = O(1/(p\epsilon))$
- **FS-HEAPRIX**, for $\eta = m/(2cLd\tau\gamma)$: $R = O((\frac{d}{m})\kappa \log(1/\epsilon))$ and $\tau = O(1/(p\epsilon))$.

*In the **convex** case, $\{\mathbf{x}^{(r)}\}_{r=0}^{\infty}$ satisfies $\mathbb{E}[f(\mathbf{x}^{(R-1)}) - f(\mathbf{x}^{(*)})] \leq \epsilon$ if:*

- **FS-PRIVIX**, for $\eta = 1/(2L(cd/m + 1)\tau\gamma)$: $R = O(L(d/m + 1)\epsilon \log(1/\epsilon))$ and $\tau = O(1/(p\epsilon^2))$.
- **FS-HEAPRIX**, for $\eta = m/(2Lcd\tau\gamma)$: $R = O(L(d/m)\epsilon \log(1/\epsilon))$ and $\tau = O(1/(p\epsilon^2))$.

Theorem 2 implies that the number of communication rounds and local updates are similar to the corresponding quantities in homogeneous setting except for the non-convex case where the number of rounds also depends on the compression rate (summarized Table 2-3 of the Appendix).

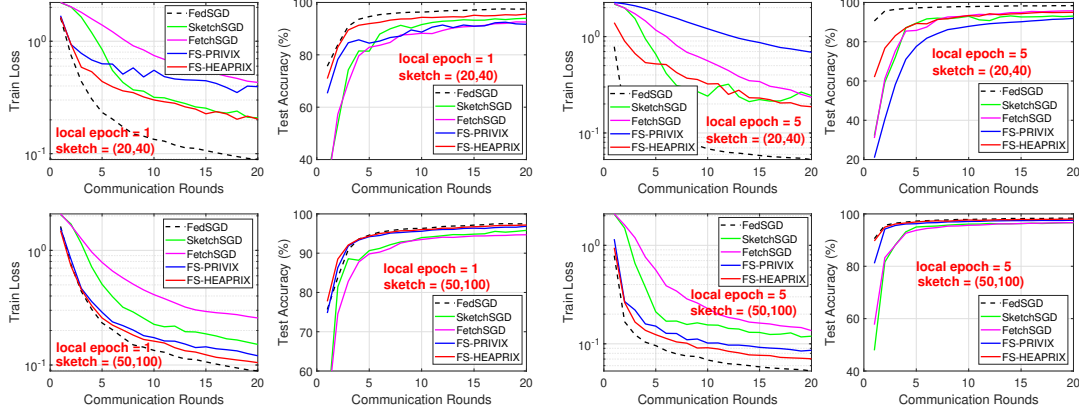


Figure 1: Homogeneous case: Comparison of compressed optimization methods on LeNet CNN.

4.3 COMPARISON WITH PRIOR METHODS

Before comparing with prior works, we highlight that privacy is another purpose of using unbiased sketching in addition to communication efficiency. Therefore, our main competing schemes are distributed algorithms based on sketching. Nonetheless, for the sake of showing the effectiveness of our algorithms, we also compare with prior non-sketching based distributed algorithms (Karimireddy et al. (2019); Basu et al. (2019); Reiszadeh et al. (2020); Haddadpour et al. (2020)) in Section B of Appendix.

Comparison with Li et al. (2019). Note that our convergence analysis does not rely on the bounded gradient assumption. We also improve both the number of communication rounds R and the size of transmitted bits B per communication round. Additionally, we highlight that, while (Li et al., 2019) provides a convergence analysis for convex objectives, our analysis holds for PL (thus strongly convex case), general convex and general non-convex objectives.

Comparison with Rothchild et al. (2020). Due to gradient tracking, our algorithm tackles data heterogeneity issue, while algorithms in Rothchild et al. (2020) does not particularly. As a consequence, in FedSKETCHGATE each device has to store an additional state vector compared to Rothchild et al. (2020). Yet, as our method is built upon an unbiased compressor, server does not need to store any additional error correction vector. The convergence results for both of two variants of FetchSGD in Rothchild et al. (2020) rely on the uniform bounded gradient assumption which may not be applicable with L -smoothness assumption when data distribution is highly heterogeneous, as in FL, see (Khaled et al., 2020), while our bounds do not assume such boundedness. Besides, Theorem 1 (Rothchild et al., 2020) assumes that *Contraction Holds* for the sequence of gradients which may not hold in practice, yet based on this strong assumption, their total communication cost (RB) in order to achieve ϵ error is $RB = O\left(m \max\left(\frac{1}{\epsilon^2}, \frac{d^2 - dm}{m^2 \epsilon}\right) \log\left(\frac{d}{\delta} \max\left(\frac{1}{\epsilon^2}, \frac{d^2 - dm}{m^2 \epsilon}\right)\right)\right)$. For the sake of comparison we let the compression ratio in Rothchild et al. (2020) to be $\frac{m}{d}$. In contrast, without any extra assumptions, our results in Theorem 2 for PRIVIX and HEAPRUX are respectively $RB = O\left(\frac{(d+m)}{\epsilon} \log\left(\frac{d^2}{\epsilon \delta} + d\right)\right)$ and $RB = O\left(\frac{d}{\epsilon} \log\left(\frac{d^2}{\epsilon m \delta}\right)\right)$ which improves the total communication cost of Theorem 1 in Rothchild et al. (2020) under regimes such that $\frac{1}{\epsilon} \geq d$ or $d \gg m$. Theorem 2 in Rothchild et al. (2020) is based the *Sliding Window Heavy Hitters* assumption, which is similar to the gradient diversity assumption in Li et al. (2020b); Haddadpour and Mahdavi (2019). Under that assumption the total communication cost is shown to be $RB = O\left(\frac{m \max(I^{2/3}, 2 - \alpha)}{\epsilon^3 \alpha} \log\left(\frac{d \max(I^{2/3}, 2 - \alpha)}{\epsilon^3 \delta}\right)\right)$ where I is a constant related to the window of gradients. We improve this bound under weaker assumptions in a regime where $\frac{I^{2/3}}{\epsilon^2} \geq d$. We also provide bounds for PL, convex and non-convex objectives contrary to Rothchild et al. (2020). Finally, we note that algorithms in Rothchild et al. (2020) are using momentum at server. While we do not use it explicitly, we can modify our algorithms to include momentum easily.

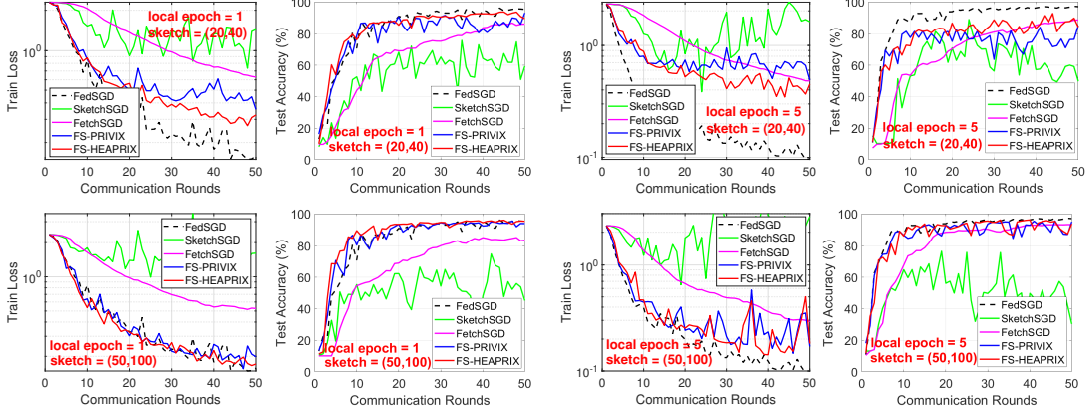


Figure 2: Heterogeneous case: Comparison of compressed optimization algorithms on LeNet CNN.

5 NUMERICAL STUDY

In this section, we provide empirical results on MNIST benchmark dataset to demonstrate the effectiveness of our proposed algorithms. We train LeNet-5 Convolutional Neural Network (CNN) architecture introduced in LeCun et al. (1998), with 60 000 parameters. We compare Federated SGD (FedSGD) as the full-precision baseline, along with four sketching methods SketchSGD (Ivkin et al., 2019), FetchSGD Rothchild et al. (2020), and two FedSketch variants FS-PRIVIX and FS-HEAPRIX. Note that in Algorithm 3, FS-PRIVIX with global learning rate $\gamma = 1$ is equivalent to the DiffSketch algorithm proposed in Li et al. (2020b). Also, SketchSGD is slightly modified to compress the change in local weights (instead of local gradient in every iteration), and FetchSGD is implemented with second round of communication for fairness. (The original proposal does not include second round of communication, which performs worse with small sketch size.) As suggested in Rothchild et al. (2020), the momentum factor of FetchSGD is set to 0.9, and we also follow some recommended implementation tricks to improve its performance, which are detailed in the Appendix. The number of workers is set to 50 and we report the results for 1 and 5 local epochs. A local epoch is finished when all workers go through their local data samples once. The local batch size is 30. In each round, we randomly choose half of the devices to be active. We tune the learning rates (η and γ , if applicable) over log-scale and report the best results, for both *homogeneous* and *heterogeneous* setting. In the former case, each device receives uniformly drawn data samples, and in the latter, it only receives samples from one or two classes among ten.

Homogeneous case. In Figure 1, we provide the training loss and test accuracy with different number of local epochs and sketch size, $(t, k) = (20, 40)$ and $(50, 100)$. Note that, these two choices of sketch size correspond to a $75\times$ and $12\times$ compression ratio, respectively. We conclude

- In general, increasing compression ratio would sacrifice learning performance. In all cases, FS-HEAPRIX performs the best in terms of both training objective and test accuracy, among all compressed methods.
- FS-HEAPRIX is better than FS-PRIVIX, especially with small sketches (high compression ratio). FS-HEAPRIX yields acceptable extra test error compared to full-precision FedSGD, particularly when considering the high compression ratio (e.g., $75\times$).
- From the training loss, we see that the performance of FS-HEAPRIX improves when the number of local updates increases. *That is, the proposed method is able to further reduce the communication cost by reducing the number of rounds required for communication.* This is also consistent with our theoretical findings.

In general, our proposed FS-HEAPRIX outperforms all competing methods, and a sketch size of $(50, 100)$ is sufficient to approach the accuracy of full-precision FedSGD.

Heterogeneous case. We plot similar set of results in Figure 2 for non-i.i.d. data distribution, which leads to more twists and turns in the training curves. We see that SketchSGD performs very poorly in the heterogeneous case, which is improved by error tracking and momentum in FetchSGD,

as expected. However, both of these methods are worse than our proposed FedSketchGATE methods, which can achieve similar generalization accuracy as full-precision FedSGD, even with small sketch size (i.e., $75\times$ compression with 1 local epoch). Note that, slower convergence and worse generalization of FedSGD in non-i.i.d. data distribution case is also reported in e.g. McMahan et al. (2017); Chen et al. (2020).

We also notice in Figure 2 the edge of FS-HEAPRIX over FS-PRIVIX in terms of training loss and test accuracy. However, we see that in the heterogeneous setting, more local updates tend to undermine the learning performance, especially with small sketch size. Nevertheless, when the sketch size is not too small, i.e., (50, 100), FS-HEAPRIX can still provide comparable test accuracy as FedSGD in both cases. Our empirical study demonstrates that FedSketch (and FedSketchGATE) frameworks are able to perform well in homogeneous (resp. heterogeneous) settings, with high compression rate. In particular, FedSketch methods are beneficial over SketchedSGD (Ivkin et al., 2019) and FetchSGD Rothchild et al. (2020) in all cases. FS-HEAPRIX performs the best among all the tested compressed algorithms, which in many cases achieves similar generalization accuracy as full-precision FedSGD with small sketch size.

6 CONCLUSION

In this paper, we introduced FedSKETCH and FedSKETCHGATE algorithms for homogeneous and heterogeneous data distribution setting respectively for Federated Learning wherein communication between server and devices is only performed using count sketch. Our algorithms, thus, provide communication-efficiency and privacy, through random hashes based sketches. We analyze the convergence error for *non-convex*, *PL* and *general convex* objective functions in the scope of Federated Optimization. We provide insightful numerical experiments showcasing the advantages of our FedSKETCH and FedSKETCHGATE methods over current federated optimization algorithm. The proposed algorithms outperform competing compression method and can achieve comparable test accuracy as Federated SGD, with high compression ratio.

REFERENCES

- 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 (NIPS)*, pages 1709–1720, Long Beach, 2017.
- D. Alistarh, T. Hoefer, M. Johansson, N. Konstantinov, S. Khirirat, and C. Renggli. The convergence of sparsified gradient methods. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 5973–5983, Montréal, Canada, 2018.
- D. Basu, D. Data, C. Karakus, and S. N. Diggavi. Qsparse-local-sgd: Distributed SGD with quantization, sparsification and local computations. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 14668–14679, Vancouver, Canada, 2019.
- J. Bernstein, Y. Wang, K. Azizzadenesheli, and A. Anandkumar. SIGNSGD: compressed optimisation for non-convex problems. In *Proceedings of the 35th International Conference on Machine Learning (ICML)*, pages 559–568, Stockholmsmässan, Stockholm, Sweden, 2018.
- K. Bonawitz, V. Ivanov, B. Kreuter, A. Marcedone, H. B. McMahan, S. Patel, D. Ramage, A. Segal, and K. Seth. Practical secure aggregation for privacy-preserving machine learning. In *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security (CCS)*, pages 1175–1191, Dallas, TX, 2017.
- L. Bottou and O. Bousquet. The tradeoffs of large scale learning. In *Advances in Neural Information Processing Systems (NIPS)*, pages 161–168, Vancouver, Canada, 2008.
- M. Charikar, K. C. Chen, and M. Farach-Colton. Finding frequent items in data streams. *Theoretical Computer Science*, 312(1):3–15, 2004. doi: 10.1016/S0304-3975(03)00400-6. URL [https://doi.org/10.1016/S0304-3975\(03\)00400-6](https://doi.org/10.1016/S0304-3975(03)00400-6).
- X. Chen, X. Li, and P. Li. Toward communication efficient adaptive gradient method. In *ACM-IMS Foundations of Data Science Conference (FODS)*, Seattle, WA, 2020.
- G. Cormode and S. Muthukrishnan. An improved data stream summary: the count-min sketch and its applications. *Journal of Algorithms*, 55(1):58–75, 2005.
- C. Dwork. Differential privacy. In *Automata, Languages and Programming, 33rd International Colloquium, ICALP 2006, Venice, Italy, July 10-14, 2006, Proceedings, Part II*, volume 4052 of *Lecture Notes in Computer Science*, pages 1–12. Springer, 2006.
- R. C. Geyer, T. Klein, and M. Nabi. Differentially private federated learning: A client level perspective. *arXiv preprint arXiv:1712.07557*, 2017.
- F. Haddadpour and M. Mahdavi. On the convergence of local descent methods in federated learning. *arXiv preprint arXiv:1910.14425*, 2019.
- F. Haddadpour, M. M. Kamani, A. Mokhtari, and M. Mahdavi. Federated learning with compression: Unified analysis and sharp guarantees. *arXiv preprint arXiv:2007.01154*, 2020.
- S. Hardy, W. Henecka, H. Ivey-Law, R. Nock, G. Patrini, G. Smith, and B. Thorne. Private federated learning on vertically partitioned data via entity resolution and additively homomorphic encryption. *arXiv preprint arXiv:1711.10677*, 2017.
- S. Horváth and P. Richtárik. A better alternative to error feedback for communication-efficient distributed learning. *arXiv preprint arXiv:2006.11077*, 2020.
- S. Horváth, D. Kovalev, K. Mishchenko, S. Stich, and P. Richtárik. Stochastic distributed learning with gradient quantization and variance reduction. *arXiv preprint arXiv:1904.05115*, 2019.
- N. Iykin, D. Rothchild, E. Ullah, V. Braverman, I. Stoica, and R. Arora. Communication-efficient distributed SGD with sketching. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 13144–13154, Vancouver, Canada, 2019.

- P. Kairouz, H. B. McMahan, B. Avent, A. Bellet, M. Bennis, A. N. Bhagoji, K. Bonawitz, Z. Charles, G. Cormode, R. Cummings, et al. Advances and open problems in federated learning. *arXiv preprint arXiv:1912.04977*, 2019.
- H. Karimi, J. Nutini, and M. Schmidt. Linear convergence of gradient and proximal-gradient methods under the polyak-tojasiewicz condition. In *Proceedings of European Conference on Machine Learning and Knowledge Discovery in Databases (ECML-PKDD)*, pages 795–811, Riva del Garda, Italy, 2016.
- S. P. Karimireddy, S. Kale, M. Mohri, S. J. Reddi, S. U. Stich, and A. T. Suresh. Scaffold: Stochastic controlled averaging for on-device federated learning. *arXiv preprint arXiv:1910.06378*, 2019.
- A. Khaled, K. Mishchenko, and P. Richtárik. Tighter theory for local SGD on identical and heterogeneous data. In *The 23rd International Conference on Artificial Intelligence and Statistics (AISTATS)*, pages 4519–4529, Online [Palermo, Sicily, Italy], 2020.
- J. Kleinberg. Bursty and hierarchical structure in streams. *Data Mining and Knowledge Discovery*, 7(4):373–397, 2003.
- J. Konečný, H. B. McMahan, F. X. Yu, P. Richtárik, A. T. Suresh, and D. Bacon. Federated learning: Strategies for improving communication efficiency. *arXiv preprint arXiv:1610.05492*, 2016.
- 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.
- P. Li, K. W. Church, and T. Hastie. One sketch for all: Theory and application of conditional random sampling. In *Advances in Neural Information Processing Systems (NIPS)*, pages 953–960, Vancouver, Canada, 2008.
- T. Li, Z. Liu, V. Sekar, and V. Smith. Privacy for free: Communication-efficient learning with differential privacy using sketches. *arXiv preprint arXiv:1911.00972*, 2019.
- T. Li, A. K. Sahu, A. Talwalkar, and V. Smith. Federated learning: Challenges, methods, and future directions. *IEEE Signal Process. Mag.*, 37(3):50–60, 2020a.
- T. Li, A. K. Sahu, M. Zaheer, M. Sanjabi, A. Talwalkar, and V. Smith. Federated optimization in heterogeneous networks. In *Proceedings of Machine Learning and Systems (MLSys)*, Austin, TX, 2020b.
- X. Li, K. Huang, W. Yang, S. Wang, and Z. Zhang. On the convergence of fedavg on non-iid data. In *Proceedings of the 8th International Conference on Learning Representations (ICLR)*, Addis Ababa, Ethiopia, 2020c.
- X. Liang, S. Shen, J. Liu, Z. Pan, E. Chen, and Y. Cheng. Variance reduced local sgd with lower communication complexity. *arXiv preprint arXiv:1912.12844*, 2019.
- Y. Lin, S. Han, H. Mao, Y. Wang, and B. Dally. Deep gradient compression: Reducing the communication bandwidth for distributed training. In *Proceedings of the 6th International Conference on Learning Representations (ICLR)*, Vancouver, Canada, 2018.
- B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas. Communication-efficient learning of deep networks from decentralized data. In *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics (AISTATS)*, pages 1273–1282, Fort Lauderdale, FL, 2017.
- H. B. McMahan, D. Ramage, K. Talwar, and L. Zhang. Learning differentially private recurrent language models. In *Proceedings of the 6th International Conference on Learning Representations (ICLR)*, Vancouver, Canada, 2018.
- C. Philippenko and A. Dieuleveut. Artemis: tight convergence guarantees for bidirectional compression in federated learning. *arXiv preprint arXiv:2006.14591*, 2020.

- A. Reisizadeh, A. Mokhtari, H. Hassani, A. Jadbabaie, and R. Pedarsani. Fedpaq: A communication-efficient federated learning method with periodic averaging and quantization. In *The 23rd International Conference on Artificial Intelligence and Statistics (AISTATS)*, pages 2021–2031, Online [Palermo, Sicily, Italy], 2020.
- D. Rothchild, A. Panda, E. Ullah, N. Ivkin, I. Stoica, V. Braverman, J. Gonzalez, and R. Arora. FetchSGD: Communication-efficient federated learning with sketching. *arXiv preprint arXiv:2007.07682*, 2020.
- A. K. Sahu, T. Li, M. Sanjabi, M. Zaheer, A. Talwalkar, and V. Smith. On the convergence of federated optimization in heterogeneous networks. *arXiv preprint arXiv:1812.06127*, 2018.
- S. U. Stich. Local sgd converges fast and communicates little. In *Proceedings of the 7th International Conference on Learning Representations (ICLR)*, New Orleans, LA, 2019.
- S. U. Stich and S. P. Karimireddy. The error-feedback framework: Better rates for sgd with delayed gradients and compressed communication. *arXiv preprint arXiv:1909.05350*, 2019.
- S. U. Stich, J.-B. Cordonnier, and M. Jaggi. Sparsified sgd with memory. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 4447–4458, Montréal, Canada, 2018.
- H. Tang, S. Gan, C. Zhang, T. Zhang, and J. Liu. Communication compression for decentralized training. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 7652–7662, Montréal, Canada, 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. PMLR, 2019.
- J. Wang and G. Joshi. Cooperative sgd: A unified framework for the design and analysis of communication-efficient sgd algorithms. *arXiv preprint arXiv:1808.07576*, 2018.
- 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 (NIPS)*, pages 1509–1519, Long Beach, CA, 2017.
- J. Wu, W. Huang, J. Huang, and T. Zhang. Error compensated quantized sgd and its applications to large-scale distributed optimization. *arXiv preprint arXiv:1806.08054*, 2018.
- H. Yu, R. Jin, and S. Yang. On the linear speedup analysis of communication efficient momentum SGD for distributed non-convex optimization. In *Proceedings of the 36th International Conference on Machine Learning (ICML)*, pages 7184–7193, Long Beach, CA, 2019a.
- H. Yu, S. Yang, and S. Zhu. Parallel restarted SGD with faster convergence and less communication: Demystifying why model averaging works for deep learning. In *The Thirty-Third AAAI Conference on Artificial Intelligence (AAAI)*, pages 5693–5700, Honolulu, HI, 2019b.
- S. Zheng, Z. Huang, and J. T. Kwok. Communication-efficient distributed blockwise momentum sgd with error-feedback. *arXiv preprint arXiv:1905.10936*, 2019.
- F. Zhou and G. Cong. On the convergence properties of a k-step averaging stochastic gradient descent algorithm for nonconvex optimization. In *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence (IJCAI)*, pages 3219–3227, Stockholm, Sweden, 2018.

Appendix for FedSKETCH: Communication-Efficient Federated Learning via Sketching

The appendix is organized as follows: Section A recalls important notations used throughout the paper and provides the formulation of related algorithms used in the main paper and omitted for the sake of the page limit. We present in Section B of this supplementary file, a through comparison with notable related works. Section C contains the proofs of our results and Section D presents additional numerical runs.

A NOTATIONS AND DEFINITIONS

Notation. Here we denote the count sketch of the vector \mathbf{x} by $\mathbf{S}(\mathbf{x})$ and with an abuse of notation, we indicate the expectation over the randomness of count sketch with $\mathbb{E}_{\mathbf{S}}[\cdot]$. We illustrate the random subset of the devices selected by the central server with \mathcal{K} with size $|\mathcal{K}| = k \leq p$, and we represent the expectation over the device sampling with $\mathbb{E}_{\mathcal{K}}[\cdot]$.

Table 1: Table of Notations

p	\triangleq	Number of devices
k	\triangleq	Number of sampled devices for homogeneous setting
$\mathcal{K}^{(r)}$	\triangleq	Set of sampled devices in communication round r
d	\triangleq	Dimension of the model
τ	\triangleq	Number of local updates
R	\triangleq	Number of communication rounds
B	\triangleq	Size of transmitted bits
$R \times B$	\triangleq	Total communication cost per device
κ	\triangleq	Condition number
ϵ	\triangleq	Target accuracy
μ	\triangleq	PL constant
m	\triangleq	Number of bins of hash tables
$\mathbf{S}(\mathbf{x})$	\triangleq	Count sketch of the vector \mathbf{x}
$\mathbb{U}(\Delta)$	\triangleq	Class of unbiased compressor, see Definition 1

Definition 3 (Polyak-Łojasiewicz). *A function $f(\mathbf{x})$ satisfies the Polyak-Łojasiewicz(PL) condition with constant μ if $\frac{1}{2}\|\nabla f(\mathbf{x})\|_2^2 \geq \mu(f(\mathbf{x}) - f(\mathbf{x}^*))$, $\forall \mathbf{x} \in \mathbb{R}^d$ with \mathbf{x}^* is an optimal solution.*

A.1 COUNT SKETCH

In this paper, we exploit the commonly used `Count Sketch` (Charikar et al., 2004) which is described in Algorithm 5.

Algorithm 5 Count Sketch (CS) (Charikar et al., 2004)

```

1: Inputs:  $\mathbf{x} \in \mathbb{R}^d, t, k, \mathbf{S}_{m \times t}, h_j(1 \leq i \leq t), \text{sign}_j(1 \leq i \leq t)$ 
2: Compress vector  $\mathbf{x} \in \mathbb{R}^d$  into  $\mathbf{S}(\mathbf{x})$ :
3: for  $x_i \in \mathbf{x}$  do
4:   for  $j = 1, \dots, t$  do
5:      $\mathbf{S}[j][h_j(i)] = \mathbf{S}[j-1][h_{j-1}(i)] + \text{sign}_j(i) \cdot x_i$ 
6:   end for
7: end for
8: return  $\mathbf{S}_{m \times t}(\mathbf{x})$ 

```

A.2 PRIVIX METHOD AND COMPRESSION ERROR OF HEAPRIX

For the sake of completeness we review PRIVIX algorithm that is also mentioned in Li et al. (2019) as follows:

Algorithm 6 PRIVIX/DiffSketch (Li et al., 2019): Unbiased compressor based on sketching.

```

1: Inputs:  $\mathbf{x} \in \mathbb{R}^d, t, m, \mathbf{S}_{m \times t}, h_j(1 \leq i \leq t), \text{sign}_j(1 \leq i \leq t)$ 
2: Query  $\tilde{\mathbf{x}} \in \mathbb{R}^d$  from  $\mathbf{S}(\mathbf{x})$ :
3: for  $i = 1, \dots, d$  do
4:    $\tilde{\mathbf{x}}[i] = \text{Median}\{\text{sign}_j(i) \cdot \mathbf{S}[j][h_j(i)] : 1 \leq j \leq t\}$ 
5: end for
6: Output:  $\tilde{\mathbf{x}}$ 

```

Regarding the compression error of sketching we restate the following Corollary from the main body of this paper:

Corollary 2. *Based on Theorem 3 of (Horváth and Richtárik, 2020) and using Algorithm 2, we have $C(x) \in \mathbb{U}(c \frac{d}{m})$. This shows that unlike PRIVIX (Algorithm 6) the compression noise can be made as small as possible using large size of hash table.*

Proof. The proof simply follows from Theorem 3 in Horváth and Richtárik (2020) and Algorithm 2 by setting $\Delta_1 = c \frac{d}{m}$ and $\Delta_2 = 1 + c \frac{d}{m}$ we obtain $\Delta = \Delta_2 + \frac{1-\Delta_2}{\Delta_1} = c \frac{d}{m} = O\left(\frac{d}{m}\right)$ for the compression error of HEAPRIX. \square

B SUMMARY OF COMPARISON OF OUR RESULTS WITH PRIOR WORKS

For the purpose of further clarification, we summarize the comparison of our results with related works. We recall that p is the number of devices, d is the dimension of the model, κ is the condition number, ϵ is the target accuracy, R is the number of communication rounds, and τ is the number of local updates. We start with the homogeneous setting comparison. Comparison of our results and existing ones for homogeneous and heterogeneous setting are given respectively Table 2 and Table 3.

Table 2: Comparison of results with compression and periodic averaging in the homogeneous setting. UG and PP stand for Unbounded Gradient and Privacy Property respectively.

Reference	PL/Strongly Convex	UG	PP
Ivkin et al. (Ivkin et al., 2019)	$R = O\left(\max\left(\frac{d}{m\sqrt{\epsilon}}, \frac{1}{\epsilon}\right)\right)$, $\tau = 1$, $B = O\left(m \log\left(\frac{dR}{\delta}\right)\right)$ $pRB = O\left(\frac{pd}{m\epsilon} \log\left(\frac{d}{\delta\sqrt{\epsilon}} \max\left(\frac{d}{m}, \frac{1}{\sqrt{\epsilon}}\right)\right)\right)$	✗	✗
Theorem 1	$R = O\left(\kappa\left(\frac{d-m}{mk} + 1\right) \log\left(\frac{1}{\epsilon}\right)\right)$, $\tau = O\left(\frac{d}{k(\frac{d}{k}+m)\epsilon}\right)$, $B = O\left(m \log\left(\frac{dR}{\delta}\right)\right)$ $kRB = O\left(m\kappa(d-m+mk) \log\frac{1}{\epsilon} \log\left(\frac{\kappa(d\frac{d-m}{mk}+d) \log\frac{1}{\epsilon}}{\delta}\right)\right)$	✓	✓

Table 3: Comparison of results with compression and periodic averaging in the heterogeneous setting. UG and PP stand for Unbounded Gradient and Privacy Property respectively.

Reference	non-convex	General Convex	UG	PP
Basu et al. (Basu et al., 2019) (with $\gamma = m/d$)	$R = O\left(\frac{d}{m\epsilon^{1-\gamma}}\right)$ $\tau = O\left(\frac{m}{pd\sqrt{\epsilon}}\right)$ $B = O\left(\frac{d}{d}\right)$ $RB = O\left(\frac{d^2}{m\epsilon^{1-\gamma}}\right)$	–	✗	✗
Li et al. (Li et al., 2019)	–	$R = O\left(\frac{d}{m\epsilon^2}\right)$ $\tau = 1$ $B = O\left(m \log\left(\frac{d^2}{m\epsilon^2\delta}\right)\right)$	✗	✓
Rothchild et al. (Rothchild et al., 2020)	$R = O\left(\max\left(\frac{1}{\epsilon^2}, \frac{d^2-m\delta}{m^2\epsilon}\right)\right)$ $\tau = 1$ $B = O\left(m \log\left(\frac{d}{\delta} \max\left(\frac{1}{\epsilon^2}, \frac{d^2-m\delta}{m^2\epsilon}\right)\right)\right)$ $RB = O\left(m \max\left(\frac{1}{\epsilon^2}, \frac{d^2-m\delta}{m^2\epsilon}\right) \log\left(\frac{d}{\delta} \max\left(\frac{1}{\epsilon^2}, \frac{d^2-m\delta}{m^2\epsilon}\right)\right)\right)$	–	✗	✗
Rothchild et al. (Rothchild et al., 2020)	$R = O\left(\frac{\max(I^{2/3}, 2-\alpha)}{\epsilon^3}\right)$ $\tau = 1$ $B = O\left(\frac{m}{\alpha} \log\left(\frac{d \max(I^{2/3}, 2-\alpha)}{\epsilon^3\delta}\right)\right)$ $RB = O\left(\frac{m \max(I^{2/3}, 2-\alpha)}{\epsilon^3\alpha} \log\left(\frac{d \max(I^{2/3}, 2-\alpha)}{\epsilon^3\delta}\right)\right)$	–	✗	✗
Theorem 2	$R = O\left(\frac{d}{m\epsilon}\right)$ $\tau = O\left(\frac{1}{p\epsilon}\right)$ $B = O\left(m \log\left(\frac{d^2}{m\epsilon\delta}\right)\right)$ $RB = O\left(\frac{d}{\epsilon} \log\left(\frac{d^2}{m\epsilon\delta} \log\left(\frac{1}{\epsilon}\right)\right)\right)$	$R = O\left(\frac{d}{m\epsilon} \log\left(\frac{1}{\epsilon}\right)\right)$ $\tau = O\left(\frac{1}{p\epsilon^2}\right)$ $B = O\left(m \log\left(\frac{d^2}{m\epsilon\delta}\right)\right)$	✓	✓

Comparison with Haddadpour et al. (2020) and Reisizadeh et al. (2020) Convergence analysis of algorithms in Haddadpour et al. (2020) relies on unbiased compression, while in this paper our FL algorithm based on HEAPRIX enjoys from unbiased compression with equivalent biased compression variance. Moreover, we highlight that the convergence analysis of FedCOMGATE is based on the extra assumption of boundedness of the difference between the average of compressed vectors and compressed averages of vectors. However, we do not need this extra assumption as it is satisfied naturally due to linearity of sketching. Finally, as pointed out in Remark 2, our algorithms enjoy from a bidirectional compression property, unlike FedCOMGATE in general. Furthermore, since results in Haddadpour et al. (2020) improve the communication complexity of FedPAQ algorithm, developed in Reisizadeh et al. (2020), hence FedSKETCH and FedSKETCHGATE improves the communication complexity obtained in Reisizadeh et al. (2020).

Comparison with Basu et al. (2019). We note that the algorithm in Basu et al. (2019) uses a composed compression and quantization while our algorithm is solely based on compression. So, in order

to compare with algorithms in Basu et al. (2019) we only consider Qsparse-local-SGD with compression and we let compression factor $\gamma = \frac{m}{d}$ (to compare with the same compression ratio induced with sketch size of mt). For strongly convex objective in Qsparse-local-SGD to achieve convergence error of ϵ they require $R = O\left(\kappa \frac{d}{m\sqrt{\epsilon}}\right)$ and $\tau = O\left(\frac{m}{pd\sqrt{\epsilon}}\right)$, which is improved to $R = O\left(\frac{\kappa d}{m} \log(1/\epsilon)\right)$ and $\tau = O\left(\frac{1}{p\epsilon}\right)$ for PL objectives. Similarly, for non-convex objective Basu et al. (2019) requires $R = O\left(\frac{d}{m\epsilon^{1.5}}\right)$ and $\tau = O\left(\frac{m}{pd\sqrt{\epsilon}}\right)$, which is improved to $R = O\left(\frac{d}{m\epsilon}\right)$ and $\tau = O\left(\frac{1}{p\epsilon}\right)$. We note that we reduce communication rounds at the cost of increasing number of local updates (which scales down with number of devices, p). Additionally, we highlight that our FedSKETCHGATE exploits the gradient tracking idea to deal with data heterogeneity, while algorithms in Basu et al. (2019) does not develop such mechanism and may suffer from poor convergence in heterogeneous setting. We also note that setting $\tau = 1$ and using top_m compressor, the QSPARSE-local-SGD algorithm becomes similar to distributed SGD with sketching as they both use the error feedback framework to improve the compression variance. Finally, since the average of sparse vectors may not be sparse in general the number of transmitted bits from server to devices in QSPARSE-Local-SGD in Basu et al. (2019) may not be sparse in general ($B = O(d)$), however our algorithms enjoy from bidirectional compression properly due to lower dimension and linearity properties of sketching ($B = O(m \log(\frac{Rd}{\delta}))$). Therefore, the total number of bits per device for strongly convex and non-convex objective is improved respectively from $RB = O\left(\kappa \frac{d^2}{m\sqrt{\epsilon}}\right)$ and $RB = O\left(\frac{d^2}{m\epsilon^{1.5}}\right)$ in Basu et al. (2019) to $RB = O\left(\kappa d \log(\frac{\kappa d^2}{m\delta} \log(\frac{1}{\epsilon})) \log(1/\epsilon)\right) = O\left(\kappa d \max\left(\log(\frac{\kappa d^2}{m\delta}), \log^2(1/\epsilon)\right)\right)$ and $RB = O\left(\log(\frac{d^2}{m\epsilon\delta}) \frac{d}{\epsilon}\right)$.

Additionally, as we noted using sketching for transmission implies two way communication from master to devices and vice versa. Therefore, in order to show efficacy of our algorithm we compare our convergence analysis with the obtained rates in the following related work:

Comparison with Philippenko and Dieuleveut (2020). The reference (Philippenko and Dieuleveut, 2020) considers two-way compression from parameter server to devices and vice versa. They provide the convergence rate of $R = O\left(\frac{\omega^{\text{Up}} \omega^{\text{Down}}}{\epsilon^2}\right)$ for strongly-objective functions where ω^{Up} and ω^{Down} are uplink and downlink's compression noise (specializing to our case for the sake of comparison $\omega^{\text{Up}} = \omega^{\text{Down}} = \theta(d)$) for general heterogeneous data distribution. In contrast, while our algorithms are using bidirectional compression due to use of sketching for communication, our convergence rate for strongly-convex objective is $R = O(\kappa \mu^2 d \log(\frac{1}{\epsilon}))$ with probability $1 - \delta$.

We would like to also mention that there prior studies such as Tang et al. (2019) and Zheng et al. (2019) that analyze the two-way compression, but since Philippenko and Dieuleveut (2020) is the state-of-the-art on this topic we only compared our results with these papers.

C THEORETICAL PROOFS

We will use the following fact (which is also used in Li et al. (2020c); Haddadpour and Mahdavi (2019)) in proving results.

Fact 3 (Li et al. (2020c); Haddadpour and Mahdavi (2019)). *Let $\{x_i\}_{i=1}^p$ denote any fixed deterministic sequence. We sample a multiset \mathcal{P} (with size K) uniformly at random where x_j is sampled with probability q_j for $1 \leq j \leq p$ with replacement. Let $\mathcal{P} = \{i_1, \dots, i_K\} \subset [p]$ (some i_j s may have the same value). Then*

$$\mathbb{E}_{\mathcal{P}} \left[\sum_{i \in \mathcal{P}} x_i \right] = \mathbb{E}_{\mathcal{P}} \left[\sum_{k=1}^K x_{i_k} \right] = K \mathbb{E}_{\mathcal{P}} [x_{i_k}] = K \left[\sum_{j=1}^p q_j x_j \right] \quad (2)$$

For the sake of the simplicity, we review an assumption for the quantization/compression, that naturally holds for PRIVIX and HEAPRIX.

Assumption 4 (Haddadpour et al. (2020)). *The output of the compression operator $Q(\mathbf{x})$ is an unbiased estimator of its input \mathbf{x} , and its variance grows with the squared of the squared of ℓ_2 -norm of its argument, i.e., $\mathbb{E}[Q(\mathbf{x})] = \mathbf{x}$ and $\mathbb{E}[\|Q(\mathbf{x}) - \mathbf{x}\|^2] \leq \omega \|\mathbf{x}\|^2$.*

We note that the sketching PRIVIX and HEAPRIX, satisfy Assumption 4 with $\omega = c \frac{d}{m}$ and $\omega = c \frac{d}{m} - 1$ respectively with probability $1 - \frac{\delta}{R}$ per communication round. Therefore, all the results in Theorem 1, by taking union over the all probabilities of each communication rounds, are concluded with probability $1 - \delta$ by plugging $\omega = c \frac{d}{m}$ and $\omega = c \frac{d}{m} - 1$ respectively into the corresponding convergence bounds.

C.1 PROOF OF THEOREM 1

In this section, we study the convergence properties of our FedSKETCH method presented in Algorithm 3. Before developing the proofs for FedSKETCH in the homogeneous setting, we first mention the following intermediate lemmas.

Lemma 1. *Using unbiased compression and under Assumption 2, we have the following bound:*

$$\mathbb{E}_{\mathcal{K}} \left[\mathbb{E}_{\mathbf{S}, \xi^{(r)}} \left[\|\tilde{\mathbf{g}}_{\mathbf{S}}^{(r)}\|^2 \right] \right] = \mathbb{E}_{\xi^{(r)}} \mathbb{E}_{\mathbf{S}} \left[\|\tilde{\mathbf{g}}_{\mathbf{S}}^{(r)}\|^2 \right] \leq \tau \left(\frac{\omega}{k} + 1 \right) \sum_{j=1}^m q_j \left[\sum_{c=0}^{\tau-1} \|\mathbf{g}_j^{(c,r)}\|^2 + \sigma^2 \right] \quad (3)$$

Proof.

$$\begin{aligned} & \mathbb{E}_{\xi^{(r)} | \mathbf{w}^{(r)}} \mathbb{E}_{\mathcal{K}} \left[\mathbb{E}_{\mathbf{S}} \left[\left\| \frac{1}{k} \sum_{j \in \mathcal{K}} \mathbf{S} \left(\sum_{c=0}^{\tau-1} \tilde{\mathbf{g}}_j^{(c,r)} \right) \right\|^2 \right] \right] \\ &= \mathbb{E}_{\xi^{(r)}} \left[\mathbb{E}_{\mathcal{K}} \left[\mathbb{E}_{\mathbf{S}} \left[\left\| \frac{1}{k} \sum_{j \in \mathcal{K}} \underbrace{\mathbf{S} \left(\sum_{c=0}^{\tau-1} \tilde{\mathbf{g}}_j^{(c,r)} \right)}_{\tilde{\mathbf{g}}_{\mathbf{S}j}^{(r)}} \right\|^2 \right] \right] \right] \\ &\stackrel{\textcircled{1}}{=} \mathbb{E}_{\xi^{(r)}} \left[\mathbb{E}_{\mathcal{K}} \left[\left\| \frac{1}{k} \sum_{j \in \mathcal{K}} \tilde{\mathbf{g}}_{\mathbf{S}j}^{(r)} - \frac{1}{k} \sum_{j \in \mathcal{K}} \mathbb{E}_{\mathbf{S}} [\tilde{\mathbf{g}}_{\mathbf{S}j}^{(r)}] \right\|^2 + \left\| \mathbb{E}_{\mathbf{S}} \left[\frac{1}{k} \sum_{j \in \mathcal{K}} \tilde{\mathbf{g}}_{\mathbf{S}j}^{(r)} \right] \right\|^2 \right] \right] \\ &\stackrel{\textcircled{2}}{=} \mathbb{E}_{\xi^{(r)}} \left[\mathbb{E}_{\mathcal{K}} \left[\mathbb{E}_{\mathbf{S}} \left[\left\| \frac{1}{k} \left[\sum_{j \in \mathcal{K}} \tilde{\mathbf{g}}_{\mathbf{S}j}^{(r)} - \sum_{j \in \mathcal{K}} \tilde{\mathbf{g}}_j^{(r)} \right] \right\|^2 + \left\| \frac{1}{k} \sum_{j \in \mathcal{K}} \tilde{\mathbf{g}}_j^{(r)} \right\|^2 \right] \right] \right] \end{aligned}$$

$$\begin{aligned}
&= \mathbb{E}_{\xi^{(r)}} \left[\mathbb{E}_{\mathcal{K}} \left[\left[\text{Var}_{\mathbf{S}} \left[\frac{1}{k} \sum_{j \in \mathcal{K}} \tilde{\mathbf{g}}_{\mathbf{S}j}^{(r)} \right] + \left\| \frac{1}{k} \sum_{j \in \mathcal{K}} \tilde{\mathbf{g}}_j^{(r)} \right\|^2 \right] \right] \right] \\
&= \mathbb{E}_{\xi^{(r)}} \left[\mathbb{E}_{\mathcal{K}} \left[\frac{1}{k^2} \sum_{j \in \mathcal{K}} \text{Var}_{\mathbf{S}j} \left[\tilde{\mathbf{g}}_{\mathbf{S}j}^{(r)} \right] + \left\| \frac{1}{k} \sum_{j \in \mathcal{K}} \tilde{\mathbf{g}}_j^{(r)} \right\|^2 \right] \right] \\
&\leq \mathbb{E}_{\xi^{(r)}} \left[\mathbb{E}_{\mathcal{K}} \left[\frac{1}{k^2} \sum_{j \in \mathcal{K}} \omega \left\| \tilde{\mathbf{g}}_j^{(r)} \right\|^2 + \left\| \frac{1}{k} \sum_{j \in \mathcal{K}} \tilde{\mathbf{g}}_j^{(r)} \right\|^2 \right] \right] \\
&= \left[\mathbb{E}_{\xi} \left[\frac{1}{k} \sum_{j \in \mathcal{K}} \omega \left\| \tilde{\mathbf{g}}_j^{(r)} \right\|^2 + \mathbb{E}_{\mathcal{K}} \mathbb{E}_{\xi^{(r)}} \left\| \frac{1}{k} \sum_{j \in \mathcal{K}} \tilde{\mathbf{g}}_j^{(r)} \right\|^2 \right] \right] \\
&= \left[\mathbb{E}_{\xi} \left[\frac{\omega}{k} \sum_{j=1}^p q_j \left\| \tilde{\mathbf{g}}_j^{(r)} \right\|^2 + \mathbb{E}_{\mathcal{K}} \left[\text{Var} \left(\frac{1}{k} \sum_{j \in \mathcal{K}} \tilde{\mathbf{g}}_j^{(r)} \right) + \left\| \frac{1}{k} \sum_{j \in \mathcal{K}} \mathbf{g}_j^{(r)} \right\|^2 \right] \right] \right] \\
&= \frac{\omega}{k} \sum_{j=1}^p q_j \mathbb{E}_{\xi} \left\| \tilde{\mathbf{g}}_j^{(r)} \right\|^2 + \mathbb{E}_{\mathcal{K}} \left[\frac{1}{k^2} \sum_{j \in \mathcal{K}} \text{Var} \left(\tilde{\mathbf{g}}_j^{(r)} \right) + \left\| \frac{1}{k} \sum_{j \in \mathcal{K}} \mathbf{g}_j^{(r)} \right\|^2 \right] \\
&\leq \frac{\omega}{k} \sum_{j=1}^p q_j \mathbb{E}_{\xi} \left\| \tilde{\mathbf{g}}_j^{(r)} \right\|^2 + \mathbb{E}_{\mathcal{K}} \left[\frac{1}{k^2} \sum_{j \in \mathcal{K}} \tau \sigma^2 + \frac{1}{k} \sum_{j \in \mathcal{K}} \left\| \mathbf{g}_j^{(r)} \right\|^2 \right] \\
&= \frac{\omega}{k} \sum_{j=1}^p q_j \left[\text{Var} \left(\tilde{\mathbf{g}}_j^{(r)} \right) + \left\| \mathbf{g}_j^{(r)} \right\|^2 \right] + \left[\frac{\tau \sigma^2}{k} + \sum_{j=1}^p q_j \left\| \mathbf{g}_j^{(r)} \right\|^2 \right] \\
&\leq \frac{\omega}{k} \sum_{j=1}^p q_j \left[\tau \sigma^2 + \left\| \mathbf{g}_j^{(r)} \right\|^2 \right] + \left[\frac{\tau \sigma^2}{k} + \sum_{j=1}^p q_j \left\| \mathbf{g}_j^{(r)} \right\|^2 \right] \\
&= (\omega + 1) \frac{\tau \sigma^2}{k} + \left(\frac{\omega}{k} + 1 \right) \left[\sum_{j=1}^p q_j \left\| \mathbf{g}_j^{(r)} \right\|^2 \right] \tag{4}
\end{aligned}$$

where ① holds due to $\mathbb{E} \left[\left\| \mathbf{x} \right\|^2 \right] = \text{Var}[\mathbf{x}] + \left\| \mathbb{E}[\mathbf{x}] \right\|^2$, ② is due to $\mathbb{E}_{\mathbf{S}} \left[\frac{1}{p} \sum_{j=1}^p \tilde{\mathbf{g}}_{\mathbf{S}j}^{(r)} \right] = \frac{1}{p} \sum_{j=1}^m \tilde{\mathbf{g}}_j^{(r)}$.

Next we show that from Assumptions 3, we have

$$\mathbb{E}_{\xi^{(r)}} \left[\left\| \tilde{\mathbf{g}}_j^{(r)} - \mathbf{g}_j^{(r)} \right\|^2 \right] \leq \tau \sigma^2 \tag{5}$$

To do so, note that

$$\begin{aligned}
\text{Var} \left(\tilde{\mathbf{g}}_j^{(r)} \right) &= \mathbb{E}_{\xi^{(r)}} \left[\left\| \tilde{\mathbf{g}}_j^{(r)} - \mathbf{g}_j^{(r)} \right\|^2 \right] \stackrel{\text{①}}{=} \mathbb{E}_{\xi^{(r)}} \left[\left\| \sum_{c=0}^{\tau-1} \left[\tilde{\mathbf{g}}_j^{(c,r)} - \mathbf{g}_j^{(c,r)} \right] \right\|^2 \right] = \text{Var} \left(\sum_{c=0}^{\tau-1} \tilde{\mathbf{g}}_j^{(c,r)} \right) \\
&\stackrel{\text{②}}{=} \sum_{c=0}^{\tau-1} \text{Var} \left(\tilde{\mathbf{g}}_j^{(c,r)} \right) \\
&= \sum_{c=0}^{\tau-1} \mathbb{E} \left[\left\| \tilde{\mathbf{g}}_j^{(c,r)} - \mathbf{g}_j^{(c,r)} \right\|^2 \right] \\
&\stackrel{\text{③}}{\leq} \tau \sigma^2 \tag{6}
\end{aligned}$$

where in ① we use the definition of $\tilde{\mathbf{g}}_j^{(r)}$ and $\mathbf{g}_j^{(r)}$, in ② we use the fact that mini-batches are chosen in i.i.d. manner at each local machine, and ③ immediately follows from Assumptions 2.

Replacing $\mathbb{E}_{\xi^{(r)}} [\|\tilde{\mathbf{g}}_j^{(r)} - \mathbf{g}_j^{(r)}\|^2]$ in (4) by its upper bound in (5) implies that

$$\mathbb{E}_{\xi^{(r)}|\mathbf{w}^{(r)}} \mathbb{E}_{\mathbf{S}, \mathcal{K}} \left[\left\| \frac{1}{k} \sum_{j \in \mathcal{K}} \mathbf{S} \left(\sum_{c=0}^{\tau-1} \tilde{\mathbf{g}}_j^{(c,r)} \right) \right\|^2 \right] \leq (\omega + 1) \frac{\tau \sigma^2}{k} + \left(\frac{\omega}{k} + 1 \right) \sum_{j=1}^p q_j \|\mathbf{g}_j^{(r)}\|^2 \quad (7)$$

Further note that we have

$$\left\| \mathbf{g}_j^{(r)} \right\|^2 = \left\| \sum_{c=0}^{\tau-1} \mathbf{g}_j^{(c,r)} \right\|^2 \leq \tau \sum_{c=0}^{\tau-1} \|\mathbf{g}_j^{(c,r)}\|^2 \quad (8)$$

where the last inequality is due to $\left\| \sum_{i=1}^n \mathbf{a}_i \right\|^2 \leq n \sum_{i=1}^n \|\mathbf{a}_i\|^2$, which together with (7) leads to the following bound:

$$\mathbb{E}_{\xi^{(r)}|\mathbf{w}^{(r)}} \mathbb{E}_{\mathbf{S}} \left[\left\| \frac{1}{k} \sum_{j \in \mathcal{K}} \mathbf{S} \left(\sum_{c=0}^{\tau-1} \tilde{\mathbf{g}}_j^{(c,r)} \right) \right\|^2 \right] \leq (\omega + 1) \frac{\tau \sigma^2}{k} + \tau \left(\frac{\omega}{k} + 1 \right) \sum_{j=1}^p q_j \|\mathbf{g}_j^{(c,r)}\|^2, \quad (9)$$

and the proof is complete. \square

Lemma 2. *Under Assumption 1, and according to the FedCOM algorithm the expected inner product between stochastic gradient and full batch gradient can be bounded with:*

$$-\mathbb{E}_{\xi, \mathbf{S}, \mathcal{K}} \left[\left\langle \nabla f(\mathbf{w}^{(r)}), \tilde{\mathbf{g}}^{(r)} \right\rangle \right] \leq \frac{1}{2} \eta \frac{1}{m} \sum_{j=1}^m \sum_{c=0}^{\tau-1} \left[-\|\nabla f(\mathbf{w}^{(r)})\|_2^2 - \|\nabla f(\mathbf{w}_j^{(c,r)})\|_2^2 + L^2 \|\mathbf{w}^{(r)} - \mathbf{w}_j^{(c,r)}\|_2^2 \right] \quad (10)$$

Proof. We have:

$$\begin{aligned} & -\mathbb{E}_{\{\xi_1^{(t)}, \dots, \xi_m^{(t)} | \mathbf{w}_1^{(t)}, \dots, \mathbf{w}_m^{(t)}\}} \mathbb{E}_{\mathbf{S}, \mathcal{K}} \left[\left\langle \nabla f(\mathbf{w}^{(r)}), \tilde{\mathbf{g}}_{\mathbf{S}, \mathcal{K}}^{(r)} \right\rangle \right] \\ &= -\mathbb{E}_{\{\xi_1^{(t)}, \dots, \xi_m^{(t)} | \mathbf{w}_1^{(t)}, \dots, \mathbf{w}_m^{(t)}\}} \left[\left\langle \nabla f(\mathbf{w}^{(r)}), \eta \sum_{j \in \mathcal{K}} q_j \sum_{c=0}^{\tau-1} \tilde{\mathbf{g}}_j^{(c,r)} \right\rangle \right] \\ &= -\left\langle \nabla f(\mathbf{w}^{(r)}), \eta \sum_{j=1}^m q_j \sum_{c=0}^{\tau-1} \mathbb{E}_{\xi, \mathbf{S}} [\tilde{\mathbf{g}}_{j, \mathbf{S}}^{(c,r)}] \right\rangle \\ &= -\eta \sum_{c=0}^{\tau-1} \sum_{j=1}^m q_j \left\langle \nabla f(\mathbf{w}^{(r)}), \mathbf{g}_j^{(c,r)} \right\rangle \\ &\stackrel{\textcircled{1}}{=} \frac{1}{2} \eta \sum_{c=0}^{\tau-1} \sum_{j=1}^m q_j \left[-\|\nabla f(\mathbf{w}^{(r)})\|_2^2 - \|\nabla f(\mathbf{w}_j^{(c,r)})\|_2^2 + \|\nabla f(\mathbf{w}^{(r)}) - \nabla f(\mathbf{w}_j^{(c,r)})\|_2^2 \right] \\ &\stackrel{\textcircled{2}}{\leq} \frac{1}{2} \eta \sum_{c=0}^{\tau-1} \sum_{j=1}^m q_j \left[-\|\nabla f(\mathbf{w}^{(r)})\|_2^2 - \|\nabla f(\mathbf{w}_j^{(c,r)})\|_2^2 + L^2 \|\mathbf{w}^{(r)} - \mathbf{w}_j^{(c,r)}\|_2^2 \right] \end{aligned} \quad (11)$$

where $\textcircled{1}$ is due to $2\langle \mathbf{a}, \mathbf{b} \rangle = \|\mathbf{a}\|^2 + \|\mathbf{b}\|^2 - \|\mathbf{a} - \mathbf{b}\|^2$, and $\textcircled{2}$ follows from Assumption 1. \square

The following lemma bounds the distance of local solutions from global solution at r th communication round.

Lemma 3. *Under Assumptions 2 we have:*

$$\mathbb{E} \left[\|\mathbf{w}^{(r)} - \mathbf{w}_j^{(c,r)}\|_2^2 \right] \leq \eta^2 \tau \sum_{c=0}^{\tau-1} \left\| \mathbf{g}_j^{(c,r)} \right\|_2^2 + \eta^2 \tau \sigma^2$$

Proof. Note that

$$\begin{aligned}
\mathbb{E} \left[\left\| \mathbf{w}^{(r)} - \mathbf{w}_j^{(c,r)} \right\|_2^2 \right] &= \mathbb{E} \left[\left\| \mathbf{w}^{(r)} - \left(\mathbf{w}^{(r)} - \eta \sum_{k=0}^c \tilde{\mathbf{g}}_j^{(k,r)} \right) \right\|_2^2 \right] \\
&= \mathbb{E} \left[\left\| \eta \sum_{k=0}^c \tilde{\mathbf{g}}_j^{(k,r)} \right\|_2^2 \right] \\
&\stackrel{\textcircled{1}}{=} \mathbb{E} \left[\left\| \eta \sum_{k=0}^c \left(\tilde{\mathbf{g}}_j^{(k,r)} - \mathbf{g}_j^{(k,r)} \right) \right\|_2^2 \right] + \mathbb{E} \left[\left\| \eta \sum_{k=0}^c \mathbf{g}_j^{(k,r)} \right\|_2^2 \right] \\
&\stackrel{\textcircled{2}}{=} \eta^2 \sum_{k=0}^c \mathbb{E} \left[\left\| \left(\tilde{\mathbf{g}}_j^{(k,r)} - \mathbf{g}_j^{(k,r)} \right) \right\|_2^2 \right] + (c+1) \eta^2 \sum_{k=0}^c \mathbb{E} \left[\left\| \mathbf{g}_j^{(k,r)} \right\|_2^2 \right] \\
&\leq \eta^2 \sum_{k=0}^{\tau-1} \mathbb{E} \left[\left\| \left(\tilde{\mathbf{g}}_j^{(k,r)} - \mathbf{g}_j^{(k,r)} \right) \right\|_2^2 \right] + \tau \eta^2 \sum_{k=0}^{\tau-1} \mathbb{E} \left[\left\| \mathbf{g}_j^{(k,r)} \right\|_2^2 \right] \\
&\stackrel{\textcircled{3}}{\leq} \eta^2 \sum_{k=0}^{\tau-1} \sigma^2 + \tau \eta^2 \sum_{k=0}^{\tau-1} \mathbb{E} \left[\left\| \mathbf{g}_j^{(k,r)} \right\|_2^2 \right] \\
&= \eta^2 \tau \sigma^2 + \eta^2 \sum_{k=0}^{\tau-1} \tau \left\| \mathbf{g}_j^{(k,r)} \right\|_2^2
\end{aligned} \tag{12}$$

where $\textcircled{1}$ comes from $\mathbb{E}[\mathbf{x}^2] = \text{Var}[\mathbf{x}] + [\mathbb{E}[\mathbf{x}]]^2$ and $\textcircled{2}$ holds because $\text{Var}\left(\sum_{j=1}^n \mathbf{x}_j\right) = \sum_{j=1}^n \text{Var}(\mathbf{x}_j)$ for i.i.d. vectors \mathbf{x}_i (and i.i.d. assumption comes from i.i.d. sampling), and finally $\textcircled{3}$ follows from Assumption 2. \square

C.1.1 MAIN RESULT FOR THE NON-CONVEX SETTING

Now we are ready to present our result for the homogeneous setting. We first state and prove the result for the general non-convex objectives.

Theorem 4 (non-convex). *For FedSKETCH(τ, η, γ), for all $0 \leq t \leq R\tau - 1$, under Assumptions 1 to 2, if the learning rate satisfies*

$$1 \geq \tau^2 L^2 \eta^2 + \left(\frac{\omega}{k} + 1 \right) \eta \gamma L \tau \tag{13}$$

and all local model parameters are initialized at the same point $\mathbf{w}^{(0)}$, then the average-squared gradient after τ iterations is bounded as follows:

$$\frac{1}{R} \sum_{r=0}^{R-1} \left\| \nabla f(\mathbf{w}^{(r)}) \right\|_2^2 \leq \frac{2(f(\mathbf{w}^{(0)}) - f(\mathbf{w}^{(*)}))}{\eta \gamma \tau R} + \frac{L \eta \gamma (\omega + 1)}{k} \sigma^2 + L^2 \eta^2 \tau \sigma^2, \tag{14}$$

where $\mathbf{w}^{(*)}$ is the global optimal solution with function value $f(\mathbf{w}^{(*)})$.

Proof. Before proceeding with the proof of Theorem 4, we would like to highlight that

$$\mathbf{w}^{(r)} - \mathbf{w}_j^{(\tau,r)} = \eta \sum_{c=0}^{\tau-1} \tilde{\mathbf{g}}_j^{(c,r)}. \tag{15}$$

From the updating rule of Algorithm 3 we have

$$\mathbf{w}^{(r+1)} = \mathbf{w}^{(r)} - \gamma \eta \left(\frac{1}{k} \sum_{j \in \mathcal{K}} \mathbf{s} \left(\sum_{c=0, r}^{\tau-1} \tilde{\mathbf{g}}_j^{(c,r)} \right) \right) = \mathbf{w}^{(r)} - \gamma \left[\frac{\eta}{k} \sum_{j \in \mathcal{K}} \mathbf{s} \left(\sum_{c=0}^{\tau-1} \tilde{\mathbf{g}}_j^{(c,r)} \right) \right].$$

In what follows, we use the following notation to denote the stochastic gradient used to update the global model at r th communication round

$$\tilde{\mathbf{g}}_{\mathbf{S},\mathcal{K}}^{(r)} \triangleq \frac{\eta}{p} \sum_{j=1}^p \mathbf{S} \left(\frac{\mathbf{w}^{(r)} - \mathbf{w}_j^{(\tau,r)}}{\eta} \right) = \frac{1}{k} \sum_{j \in \mathcal{K}} \mathbf{S} \left(\sum_{c=0}^{\tau-1} \tilde{\mathbf{g}}_j^{(c,r)} \right).$$

and notice that $\mathbf{w}^{(r)} = \mathbf{w}^{(r-1)} - \gamma \tilde{\mathbf{g}}^{(r)}$.

Then using the unbiased estimation property of sketching we have:

$$\mathbb{E}_{\mathbf{S}} [\tilde{\mathbf{g}}_{\mathbf{S}}^{(r)}] = \frac{1}{k} \sum_{j \in \mathcal{K}} \left[-\eta \mathbb{E}_{\mathbf{S}} \left[\mathbf{S} \left(\sum_{c=0}^{\tau-1} \tilde{\mathbf{g}}_j^{(c,r)} \right) \right] \right] = \frac{1}{k} \sum_{j \in \mathcal{K}} \left[-\eta \left(\sum_{c=0}^{\tau-1} \tilde{\mathbf{g}}_j^{(c,r)} \right) \right] \triangleq \tilde{\mathbf{g}}_{\mathbf{S},\mathcal{K}}^{(r)}.$$

From the L -smoothness gradient assumption on global objective, by using $\tilde{\mathbf{g}}^{(r)}$ in inequality (15) we have:

$$f(\mathbf{w}^{(r+1)}) - f(\mathbf{w}^{(r)}) \leq -\gamma \langle \nabla f(\mathbf{w}^{(r)}), \tilde{\mathbf{g}}^{(r)} \rangle + \frac{\gamma^2 L}{2} \|\tilde{\mathbf{g}}^{(r)}\|^2 \quad (16)$$

By taking expectation on both sides of above inequality over sampling, we get:

$$\begin{aligned} \mathbb{E} \left[\mathbb{E}_{\mathbf{S}} [f(\mathbf{w}^{(r+1)}) - f(\mathbf{w}^{(r)})] \right] &\leq -\gamma \mathbb{E} \left[\mathbb{E}_{\mathbf{S}} [\langle \nabla f(\mathbf{w}^{(r)}), \tilde{\mathbf{g}}_{\mathbf{S}}^{(r)} \rangle] \right] + \frac{\gamma^2 L}{2} \mathbb{E} \left[\mathbb{E}_{\mathbf{S}} \|\tilde{\mathbf{g}}_{\mathbf{S}}^{(r)}\|^2 \right] \\ &\stackrel{(a)}{=} \underbrace{-\gamma \mathbb{E} \left[\langle \nabla f(\mathbf{w}^{(r)}), \tilde{\mathbf{g}}^{(r)} \rangle \right]}_{(I)} + \frac{\gamma^2 L}{2} \underbrace{\mathbb{E} \left[\mathbb{E}_{\mathbf{S}} \|\tilde{\mathbf{g}}_{\mathbf{S}}^{(r)}\|^2 \right]}_{(II)}. \end{aligned} \quad (17)$$

We proceed to use Lemma 1, Lemma 2, and Lemma 3, to bound terms (I) and (II) in right hand side of (17), which gives

$$\begin{aligned} &\mathbb{E} \left[\mathbb{E}_{\mathbf{S}} [f(\mathbf{w}^{(r+1)}) - f(\mathbf{w}^{(r)})] \right] \\ &\leq \frac{\gamma}{2} \eta \sum_{j=1}^p q_j \sum_{c=0}^{\tau-1} \left[-\|\nabla f(\mathbf{w}^{(r)})\|_2^2 - \|\mathbf{g}_j^{(c,r)}\|_2^2 + L^2 \eta^2 \sum_{c=0}^{\tau-1} \left[\tau \|\mathbf{g}_j^{(c,r)}\|_2^2 + \sigma^2 \right] \right] \\ &\quad + \frac{\gamma^2 L (\frac{\omega}{k} + 1)}{2} \left[\eta^2 \tau \sum_{j=1}^p q_j \sum_{c=0}^{\tau-1} \|\mathbf{g}_j^{(c,r)}\|_2^2 \right] + \frac{\gamma^2 \eta^2 L (\omega + 1)}{2} \frac{\tau \sigma^2}{k} \\ &\stackrel{\textcircled{1}}{\leq} \frac{\gamma \eta}{2} \sum_{j=1}^p q_j \sum_{c=0}^{\tau-1} \left[-\|\nabla f(\mathbf{w}^{(r)})\|_2^2 - \|\mathbf{g}_j^{(c,r)}\|_2^2 + \tau L^2 \eta^2 \left[\tau \|\mathbf{g}_j^{(c,r)}\|_2^2 + \sigma^2 \right] \right] \\ &\quad + \frac{\gamma^2 L (\frac{\omega}{k} + 1)}{2} \left[\eta^2 \tau \sum_{j=1}^p q_j \sum_{c=0}^{\tau-1} \|\mathbf{g}_j^{(c,r)}\|_2^2 \right] + \frac{\gamma^2 \eta^2 L (\omega + 1)}{2} \frac{\tau \sigma^2}{k} \\ &= -\eta \gamma \frac{\tau}{2} \|\nabla f(\mathbf{w}^{(r)})\|_2^2 \\ &\quad - \left(1 - \tau L^2 \eta^2 \tau - (\frac{\omega}{k} + 1) \eta \gamma L \tau \right) \frac{\eta \gamma}{2} \sum_{j=1}^p q_j \sum_{c=0}^{\tau-1} \|\mathbf{g}_j^{(c,r)}\|_2^2 + \frac{L \tau \gamma \eta^2}{2k} (k L \tau \eta + \gamma (\omega + 1)) \sigma^2 \\ &\stackrel{\textcircled{2}}{\leq} -\eta \gamma \frac{\tau}{2} \|\nabla f(\mathbf{w}^{(r)})\|_2^2 + \frac{L \tau \gamma \eta^2}{2k} (k L \tau \eta + \gamma (\omega + 1)) \sigma^2, \end{aligned} \quad (18)$$

where in $\textcircled{1}$ we incorporate outer summation $\sum_{c=0}^{\tau-1}$, and $\textcircled{2}$ follows from condition

$$1 \geq \tau L^2 \eta^2 \tau + (\frac{\omega}{k} + 1) \eta \gamma L \tau.$$

Summing up for all R communication rounds and rearranging the terms gives:

$$\frac{1}{R} \sum_{r=0}^{R-1} \|\nabla f(\mathbf{w}^{(r)})\|_2^2 \leq \frac{2(f(\mathbf{w}^{(0)}) - f(\mathbf{w}^{(*)}))}{\eta \gamma \tau R} + \frac{L \eta \gamma (\omega + 1)}{k} \sigma^2 + L^2 \eta^2 \tau \sigma^2.$$

From the above inequality, it is easy to see that in order to achieve a linear speed up, we need to have $\eta\gamma = O\left(\frac{\sqrt{k}}{\sqrt{R\tau}}\right)$. \square

Corollary 3 (Linear speed up). *In (14) for the choice of $\eta\gamma = O\left(\frac{1}{L}\sqrt{\frac{k}{R\tau(\omega+1)}}\right)$, and $\gamma \geq k$ the convergence rate reduces to:*

$$\frac{1}{R} \sum_{r=0}^{R-1} \left\| \nabla f(\mathbf{w}^{(r)}) \right\|_2^2 \leq O \left(\frac{L\sqrt{(\omega+1)} (f(\mathbf{w}^{(0)}) - f(\mathbf{w}^*))}{\sqrt{kR\tau}} + \frac{(\sqrt{(\omega+1)})^2 \sigma^2}{\sqrt{kR\tau}} + \frac{k\sigma^2}{R\gamma^2} \right). \quad (19)$$

Note that according to (19), if we pick a fixed constant value for γ , in order to achieve an ϵ -accurate solution, $R = O\left(\frac{1}{\epsilon}\right)$ communication rounds and $\tau = O\left(\frac{\omega+1}{k\epsilon}\right)$ local updates are necessary. We also highlight that (19) also allows us to choose $R = O\left(\frac{\omega+1}{\epsilon}\right)$ and $\tau = O\left(\frac{1}{k\epsilon}\right)$ to get the same convergence rate.

Remark 3. Condition in (13) can be rewritten as

$$\begin{aligned} \eta &\leq \frac{-\gamma L\tau \left(\frac{\omega}{k} + 1\right) + \sqrt{\gamma^2 (L\tau \left(\frac{\omega}{k} + 1\right))^2 + 4L^2\tau^2}}{2L^2\tau^2} \\ &= \frac{-\gamma L\tau \left(\frac{\omega}{k} + 1\right) + L\tau \sqrt{\left(\frac{\omega}{k} + 1\right)^2 \gamma^2 + 4}}{2L^2\tau^2} \\ &= \frac{\sqrt{\left(\frac{\omega}{k} + 1\right)^2 \gamma^2 + 4} - \left(\frac{\omega}{k} + 1\right) \gamma}{2L\tau}. \end{aligned} \quad (20)$$

So based on (20), if we set $\eta = O\left(\frac{1}{L\gamma} \sqrt{\frac{k}{R\tau(\omega+1)}}\right)$, it implies that:

$$R \geq \frac{\tau k}{(\omega+1) \gamma^2 \left(\sqrt{\left(\frac{\omega}{k} + 1\right)^2 \gamma^2 + 4} - \left(\frac{\omega}{k} + 1\right) \gamma \right)^2}. \quad (21)$$

We note that $\gamma^2 \left(\sqrt{\left(\frac{\omega}{k} + 1\right)^2 \gamma^2 + 4} - \left(\frac{\omega}{k} + 1\right) \gamma \right)^2 = \Theta(1) \leq 5$ therefore even for $\gamma \geq m$ we need to have

$$R \geq \frac{\tau k}{5(\omega+1)} = O\left(\frac{\tau k}{(\omega+1)}\right). \quad (22)$$

Therefore, for the choice of $\tau = O\left(\frac{\omega+1}{k\epsilon}\right)$, due to condition in (22), we need to have $R = O\left(\frac{1}{\epsilon}\right)$. Similarly, we can have $R = O\left(\frac{\omega+1}{\epsilon}\right)$ and $\tau = O\left(\frac{1}{k\epsilon}\right)$.

Corollary 4 (Special case, $\gamma = 1$). *By letting $\gamma = 1$, $\omega = 0$ and $k = p$ the convergence rate in (14) reduces to*

$$\frac{1}{R} \sum_{r=0}^{R-1} \left\| \nabla f(\mathbf{w}^{(r)}) \right\|_2^2 \leq \frac{2(f(\mathbf{w}^{(0)}) - f(\mathbf{w}^{(*)}))}{\eta R\tau} + \frac{L\eta}{p} \sigma^2 + L^2 \eta^2 \tau \sigma^2,$$

which matches the rate obtained in Wang and Joshi (2018). In this case the communication complexity and the number of local updates become

$$R = O\left(\frac{p}{\epsilon}\right), \quad \tau = O\left(\frac{1}{\epsilon}\right),$$

which simply implies that in this special case the convergence rate of our algorithm reduces to the rate obtained in Wang and Joshi (2018), which indicates the tightness of our analysis.

C.1.2 MAIN RESULT FOR THE PL/STRONGLY CONVEX SETTING

We now turn to stating the convergence rate for the homogeneous setting under PL condition which naturally leads to the same rate for strongly convex functions.

Theorem 5 (PL or strongly convex). *For $\text{FedSKETCH}(\tau, \eta, \gamma)$, for all $0 \leq t \leq R\tau - 1$, under Assumptions 1 to 2 and 3, if the learning rate satisfies*

$$1 \geq \tau^2 L^2 \eta^2 + \left(\frac{\omega}{k} + 1\right) \eta \gamma L \tau$$

and if the all the models are initialized with $\mathbf{w}^{(0)}$ we obtain:

$$\mathbb{E}[f(\mathbf{w}^{(R)}) - f(\mathbf{w}^{(*)})] \leq (1 - \eta \gamma \mu \tau)^R \left(f(\mathbf{w}^{(0)}) - f(\mathbf{w}^{(*)})\right) + \frac{1}{\mu} \left[\frac{1}{2} L^2 \tau \eta^2 \sigma^2 + (1 + \omega) \frac{\gamma \eta L \sigma^2}{2k}\right]$$

Proof. From (18) under condition:

$$1 \geq \tau L^2 \eta^2 \tau + \left(\frac{\omega}{k} + 1\right) \eta \gamma L \tau$$

we obtain:

$$\begin{aligned} \mathbb{E}[f(\mathbf{w}^{(r+1)}) - f(\mathbf{w}^{(r)})] &\leq -\eta \gamma \frac{\tau}{2} \left\| \nabla f(\mathbf{w}^{(r)}) \right\|_2^2 + \frac{L \tau \gamma \eta^2}{2k} (k L \tau \eta + \gamma(\omega + 1)) \sigma^2 \\ &\leq -\eta \mu \gamma \tau \left(f(\mathbf{w}^{(r)}) - f(\mathbf{w}^{(*)})\right) + \frac{L \tau \gamma \eta^2}{2k} (k L \tau \eta + \gamma(\omega + 1)) \sigma^2 \end{aligned} \quad (23)$$

which leads to the following bound:

$$\mathbb{E}[f(\mathbf{w}^{(r+1)}) - f(\mathbf{w}^{(*)})] \leq (1 - \eta \mu \gamma \tau) \left[f(\mathbf{w}^{(r)}) - f(\mathbf{w}^{(*)})\right] + \frac{L \tau \gamma \eta^2}{2k} (k L \tau \eta + (\omega + 1) \gamma) \sigma^2$$

By setting $\Delta = 1 - \eta \mu \gamma \tau$ we obtain the following bound:

$$\begin{aligned} &\mathbb{E}[f(\mathbf{w}^{(R)}) - f(\mathbf{w}^{(*)})] \\ &\leq \Delta^R \left[f(\mathbf{w}^{(0)}) - f(\mathbf{w}^{(*)})\right] + \frac{1 - \Delta^R}{1 - \Delta} \frac{L \tau \gamma \eta^2}{2k} (k L \tau \eta + (\omega + 1) \gamma) \sigma^2 \\ &\leq \Delta^R \left[f(\mathbf{w}^{(0)}) - f(\mathbf{w}^{(*)})\right] + \frac{1}{1 - \Delta} \frac{L \tau \gamma \eta^2}{2k} (k L \tau \eta + (\omega + 1) \gamma) \sigma^2 \\ &= (1 - \eta \mu \gamma \tau)^R \left[f(\mathbf{w}^{(0)}) - f(\mathbf{w}^{(*)})\right] + \frac{1}{\eta \mu \gamma \tau} \frac{L \tau \gamma \eta^2}{2k} (k L \tau \eta + (\omega + 1) \gamma) \sigma^2 \end{aligned} \quad (24)$$

□

Corollary 5. *If we let $\eta \gamma \mu \tau \leq \frac{1}{2}$, $\eta = \frac{1}{2L(\frac{\omega}{k} + 1)\tau \gamma}$ and $\kappa = \frac{L}{\mu}$ the convergence error in Theorem 5, with $\gamma \geq k$ results in:*

$$\begin{aligned} &\mathbb{E}[f(\mathbf{w}^{(R)}) - f(\mathbf{w}^{(*)})] \\ &\leq e^{-\eta \gamma \mu \tau R} \left(f(\mathbf{w}^{(0)}) - f(\mathbf{w}^{(*)})\right) + \frac{1}{\mu} \left[\frac{1}{2} L^2 \tau \eta^2 \sigma^2 + (1 + \omega) \frac{\gamma \eta L \sigma^2}{2k}\right] \\ &\leq e^{-\frac{R}{2(\frac{\omega}{k} + 1)\kappa}} \left(f(\mathbf{w}^{(0)}) - f(\mathbf{w}^{(*)})\right) + \frac{1}{\mu} \left[\frac{1}{2} L^2 \frac{\tau \sigma^2}{L^2 (\frac{\omega}{k} + 1)^2 \gamma^2 \tau^2} + \frac{(1 + \omega) L \sigma^2}{2 (\frac{\omega}{k} + 1) L \tau k}\right] \\ &= O \left(e^{-\frac{R}{2(\frac{\omega}{k} + 1)\kappa}} \left(f(\mathbf{w}^{(0)}) - f(\mathbf{w}^{(*)})\right) + \frac{\sigma^2}{(\frac{\omega}{k} + 1)^2 \gamma^2 \mu \tau} + \frac{(\omega + 1) \sigma^2}{\mu (\frac{\omega}{k} + 1) \tau k} \right) \end{aligned}$$

$$= O \left(e^{-\frac{R}{2(\frac{\omega}{k}+1)\kappa}} \left(f(\mathbf{w}^{(0)}) - f(\mathbf{w}^{(*)}) \right) + \frac{\sigma^2}{\gamma^2 \mu \tau} + \frac{(\omega+1)\sigma^2}{\mu \left(\frac{\omega}{k}+1\right) \tau k} \right) \quad (25)$$

which indicates that to achieve an error of ϵ , we need to have $R = O \left(\left(\frac{\omega}{k} + 1 \right) \kappa \log \left(\frac{1}{\epsilon} \right) \right)$ and $\tau = \frac{(\omega+1)}{k \left(\frac{\omega}{k} + 1 \right) \epsilon}$. Additionally, we note that if $\gamma \rightarrow \infty$, yet $R = O \left(\left(\frac{\omega}{k} + 1 \right) \kappa \log \left(\frac{1}{\epsilon} \right) \right)$ and $\tau = \frac{(\omega+1)}{k \left(\frac{\omega}{k} + 1 \right) \epsilon}$ will be necessary.

C.1.3 MAIN RESULT FOR THE GENERAL CONVEX SETTING

Theorem 6 (Convex). For a general convex function $f(\mathbf{w})$ with optimal solution $\mathbf{w}^{(*)}$, using $\text{FedSKETCH}(\tau, \eta, \gamma)$ to optimize $\tilde{f}(\mathbf{w}, \phi) = f(\mathbf{w}) + \frac{\phi}{2} \|\mathbf{w}\|^2$, for all $0 \leq t \leq R\tau - 1$, under Assumptions 1 to 2, if the learning rate satisfies

$$1 \geq \tau^2 L^2 \eta^2 + \left(\frac{\omega}{k} + 1 \right) \eta \gamma L \tau$$

and if the all the models initiate with $\mathbf{w}^{(0)}$, with $\phi = \frac{1}{\sqrt{k\tau}}$ and $\eta = \frac{1}{2L\gamma\tau(1+\frac{\omega}{k})}$ we obtain:

$$\begin{aligned} \mathbb{E} \left[f(\mathbf{w}^{(R)}) - f(\mathbf{w}^{(*)}) \right] &\leq e^{-\frac{R}{2L(1+\frac{\omega}{k})\sqrt{m\tau}}} \left(f(\mathbf{w}^{(0)}) - f(\mathbf{w}^{(*)}) \right) \\ &\quad + \left[\frac{\sqrt{k}\sigma^2}{8\sqrt{\tau}\gamma^2 \left(1 + \frac{\omega}{k} \right)^2} + \frac{(\omega+1)\sigma^2}{4 \left(\frac{\omega}{k} + 1 \right) \sqrt{k\tau}} \right] + \frac{1}{2\sqrt{k\tau}} \left\| \mathbf{w}^{(*)} \right\|^2 \end{aligned} \quad (26)$$

We note that above theorem implies that to achieve a convergence error of ϵ we need to have $R = O \left(L \left(1 + \frac{\omega}{k} \right) \frac{1}{\epsilon} \log \left(\frac{1}{\epsilon} \right) \right)$ and $\tau = O \left(\frac{(\omega+1)^2}{k \left(\frac{\omega}{k} + 1 \right)^2 \epsilon} \right)$.

Proof. Since $\tilde{f}(\mathbf{w}^{(r)}, \phi) = f(\mathbf{w}^{(r)}) + \frac{\phi}{2} \|\mathbf{w}^{(r)}\|^2$ is ϕ -PL, according to Theorem 5, we have:

$$\begin{aligned} &\tilde{f}(\mathbf{w}^{(R)}, \phi) - \tilde{f}(\mathbf{w}^{(*)}, \phi) \\ &= f(\mathbf{w}^{(r)}) + \frac{\phi}{2} \|\mathbf{w}^{(r)}\|^2 - \left(f(\mathbf{w}^{(*)}) + \frac{\phi}{2} \|\mathbf{w}^{(*)}\|^2 \right) \\ &\leq (1 - \eta\gamma\phi\tau)^R \left(f(\mathbf{w}^{(0)}) - f(\mathbf{w}^{(*)}) \right) + \frac{1}{\phi} \left[\frac{1}{2} L^2 \tau \eta^2 \sigma^2 + (1 + \omega) \frac{\gamma\eta L \sigma^2}{2k} \right] \end{aligned} \quad (27)$$

Next rearranging (27) and replacing μ with ϕ leads to the following error bound:

$$\begin{aligned} &f(\mathbf{w}^{(R)}) - f^* \\ &\leq (1 - \eta\gamma\phi\tau)^R \left(f(\mathbf{w}^{(0)}) - f(\mathbf{w}^{(*)}) \right) + \frac{1}{\phi} \left[\frac{1}{2} L^2 \tau \eta^2 \sigma^2 + (1 + \omega) \frac{\gamma\eta L \sigma^2}{2k} \right] \\ &\quad + \frac{\phi}{2} \left(\|\mathbf{w}^*\|^2 - \|\mathbf{w}^{(r)}\|^2 \right) \\ &\leq e^{-(\eta\gamma\phi\tau)R} \left(f(\mathbf{w}^{(0)}) - f(\mathbf{w}^{(*)}) \right) + \frac{1}{\phi} \left[\frac{1}{2} L^2 \tau \eta^2 \sigma^2 + (1 + \omega) \frac{\gamma\eta L \sigma^2}{2k} \right] + \frac{\phi}{2} \|\mathbf{w}^{(*)}\|^2 \end{aligned}$$

Next, if we set $\phi = \frac{1}{\sqrt{k\tau}}$ and $\eta = \frac{1}{2(1+\frac{\omega}{k})L\gamma\tau}$, we obtain that

$$\begin{aligned} &f(\mathbf{w}^{(R)}) - f^* \\ &\leq e^{-\frac{R}{2(1+\frac{\omega}{k})L\sqrt{m\tau}}} \left(f(\mathbf{w}^{(0)}) - f(\mathbf{w}^{(*)}) \right) + \sqrt{k\tau} \left[\frac{\sigma^2}{8\tau\gamma^2 \left(1 + \frac{\omega}{k} \right)^2} + \frac{(\omega+1)\sigma^2}{4 \left(\frac{\omega}{k} + 1 \right) \tau k} \right] + \frac{1}{2\sqrt{k\tau}} \left\| \mathbf{w}^{(*)} \right\|^2, \end{aligned}$$

thus the proof is complete. \square

C.2 PROOF OF THEOREM 2

The proof of Theorem 2 follows directly from the results in Haddadpour et al. (2020). We first mention the general Theorem 7 from (Haddadpour et al., 2020) for general compression noise ω . Next, since the sketching PRIVIX and HEAPRIX, satisfy Assumption 4 with $\omega = c \frac{d}{m}$ and $\omega = c \frac{d}{m} - 1$ respectively with probability $1 - \frac{\delta}{R}$ per communication round, all the results in Theorem 2, conclude from Theorem 7 with probability $1 - \delta$ (by taking union over the all probabilities of each communication rounds with probability $1 - \delta/R$) and plugging $\omega = c \frac{d}{m}$ and $\omega = c \frac{d}{m} - 1$ respectively into the corresponding convergence bounds. For the heterogeneous setting, the results in Haddadpour et al. (2020) requires the following extra assumption that naturally holds for the sketching:

Assumption 5 (Haddadpour et al. (2020)). *The compression scheme Q for the heterogeneous data distribution setting satisfies the following condition $\mathbb{E}_Q[\|\frac{1}{m} \sum_{j=1}^m Q(\mathbf{x}_j)\|^2 - \|Q(\frac{1}{m} \sum_{j=1}^m \mathbf{x}_j)\|^2] \leq G_q$.*

We note that since sketching is a linear compressor, in the case of our algorithms for heterogeneous setting we have $G_q = 0$.

Next, we restate the Theorem in Haddadpour et al. (2020) here as follows:

Theorem 7. *Consider FedCOMGATE in Haddadpour et al. (2020). If Assumptions 1, 3, 4 and 5 hold, then even for the case the local data distribution of users are different (heterogeneous setting) we have*

- **non-convex:** By choosing stepsizes as $\eta = \frac{1}{L\gamma} \sqrt{\frac{p}{R\tau(\omega+1)}}$ and $\gamma \geq p$, we obtain that the iterates satisfy $\frac{1}{R} \sum_{r=0}^{R-1} \|\nabla f(\mathbf{w}^{(r)})\|_2^2 \leq \epsilon$ if we set $R = O\left(\frac{\omega+1}{\epsilon}\right)$ and $\tau = O\left(\frac{1}{p\epsilon}\right)$.
- **Strongly convex or PL:** By choosing stepsizes as $\eta = \frac{1}{2L(\frac{\omega}{p}+1)\tau\gamma}$ and $\gamma \geq \sqrt{p\tau}$, we obtain that the iterates satisfy $\mathbb{E}[f(\mathbf{w}^{(R)}) - f(\mathbf{w}^{(*)})] \leq \epsilon$ if we set $R = O\left((\omega+1)\kappa \log\left(\frac{1}{\epsilon}\right)\right)$ and $\tau = O\left(\frac{1}{p\epsilon}\right)$.
- **Convex:** By choosing stepsizes as $\eta = \frac{1}{2L(\omega+1)\tau\gamma}$ and $\gamma \geq \sqrt{p\tau}$, we obtain that the iterates satisfy $\mathbb{E}[f(\mathbf{w}^{(R)}) - f(\mathbf{w}^{(*)})] \leq \epsilon$ if we set $R = O\left(\frac{L(1+\omega)}{\epsilon} \log\left(\frac{1}{\epsilon}\right)\right)$ and $\tau = O\left(\frac{1}{p\epsilon^2}\right)$.

Proof. Since the sketching methods PRIVIX and HEAPRIX, satisfy the Assumption 4 with $\omega = c \frac{d}{m}$ and $\omega = c \frac{d}{m} - 1$ respectively with probability $1 - \frac{\delta}{R}$ per communication round, we conclude the proofs of Theorem 2 using Theorem 7 with probability $1 - \delta$ (by taking union over all communication rounds) and plugging $\omega = c \frac{d}{m}$ and $\omega = c \frac{d}{m} - 1$ respectively into the convergence bounds. \square

D NUMERICAL EXPERIMENTS AND ADDITIONAL RESULTS

D.1 IMPLEMENTATION OF FETCHSGD

Our implementation of `FetchSGD` basically follows the original paper (Algorithm 1 in Rothchild et al. (2020)). The only difference is that, in the original algorithm, the local workers compress the gradient (in every local step) and transmit it to the central server. In our setting, we extend to the case with multiple local updates, where the difference in local weights are transmitted (same as the standard FL framework). Also, TopK compression is used to decode the sketches at the central server. We apply the same implementation trick that when accumulating the errors, we only count the non-zero coordinates and leave other coordinates zero for the accumulator. This greatly improves the empirical performance.

D.2 ADDITIONAL PLOTS FOR THE MNIST EXPERIMENTS

D.2.1 HOMOGENEOUS SETTING

In the homogeneous case, each node has same data distribution. To achieve this setting, we randomly choose samples uniformly from 10 classes of hand-written digits. The train loss and test accuracy are provided in Figure 3, where we report local epochs $\tau = 2$ in addition to the main context (single local update). The number of users is set to 50, and in each round of training we randomly pick half of the nodes to be active (i.e., receiving data and performing local updates). We can draw similar conclusion: FS-HEAPRIX consistently performs better than other competing methods. The test accuracy increases with larger τ in homogeneous setting.

D.2.2 HETEROGENEOUS SETTING

Analogously, we present experiments on MNIST dataset under heterogeneous data distribution, including $\tau = 2$. We simulate the setting by only sending samples from one digit to each local worker (very few nodes get two classes). We see from Figure 4 that FS-HEAPRIX shows consistent advantage over competing methods. SketchedSGD performs poorly in this case.

D.3 ADDITIONAL EXPERIMENTS: CIFAR-10

We conduct similar sets of experiments on CIFAR10 dataset. We also use the simple LeNet CNN structure, as in practice small models are more favorable in federated learning, due to the limitation of mobile devices. The test accuracy is presented in Figure 5 and Figure 6, for respectively homogeneous and heterogeneous data distribution. In general, we retrieve similar information as from MNIST experiments: our proposed FS-HEAPRIX improves FS-PRIVIX and SketchedSGD in all cases. We note that although the test accuracy provided by LeNet cannot reach the state-of-the-art accuracy given by some huge models, it is also informative in terms of comparing the relative performance of different sketching methods.

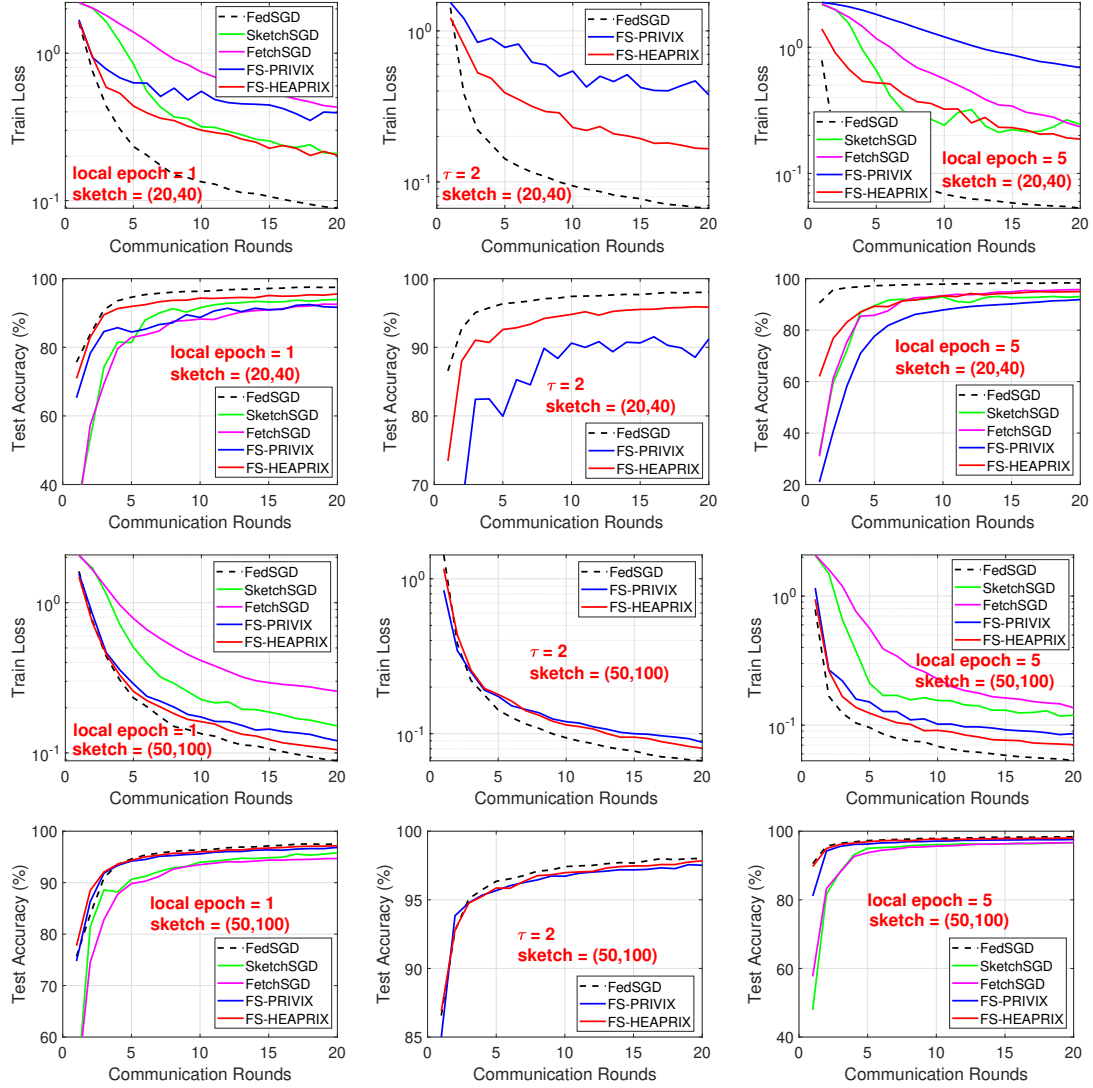


Figure 3: MNIST Homogeneous case: Comparison of compressed optimization methods on LeNet CNN architecture.

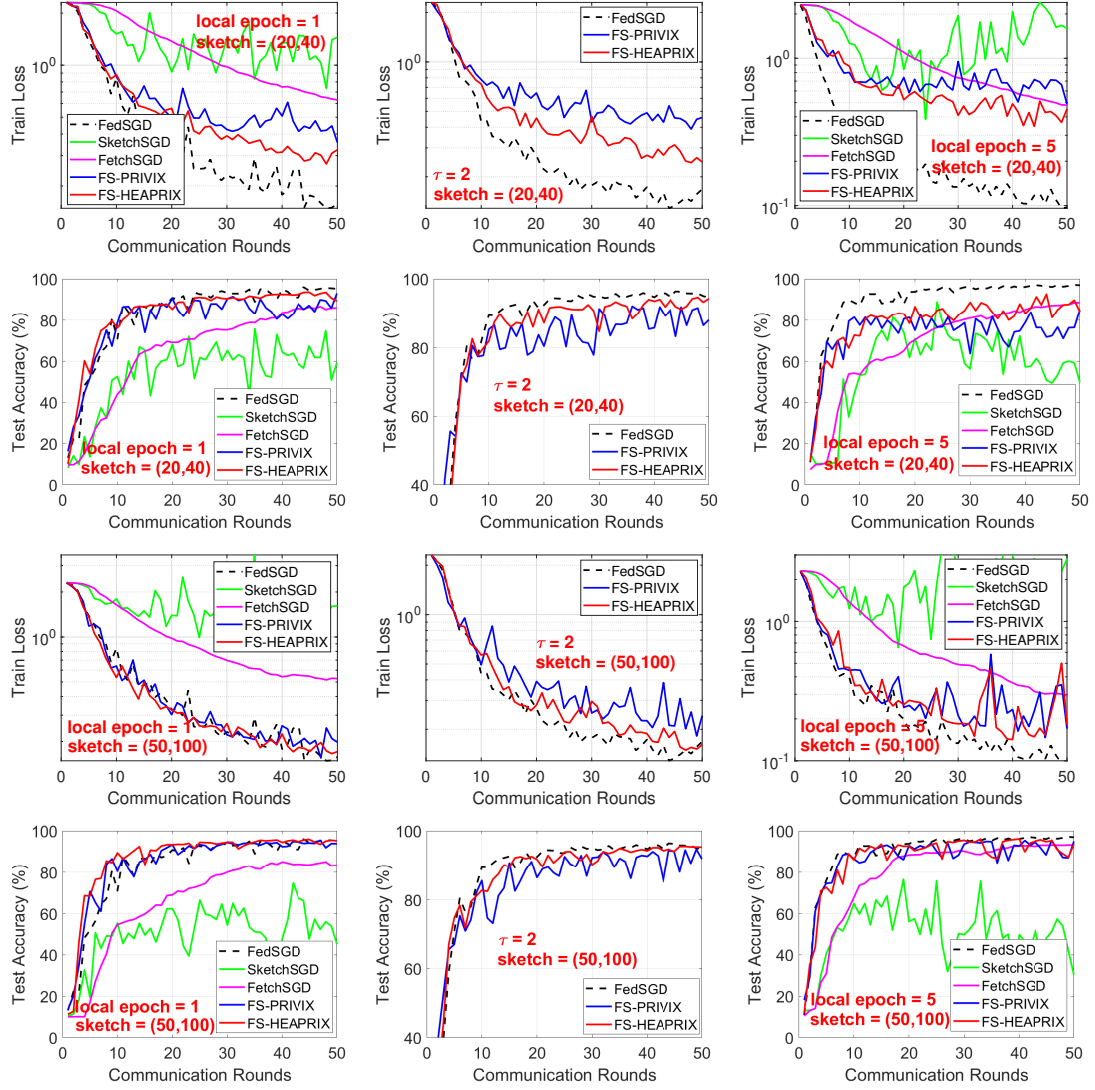


Figure 4: MNIST Heterogeneous case: Comparison of compressed optimization algorithms on LeNet CNN architecture.

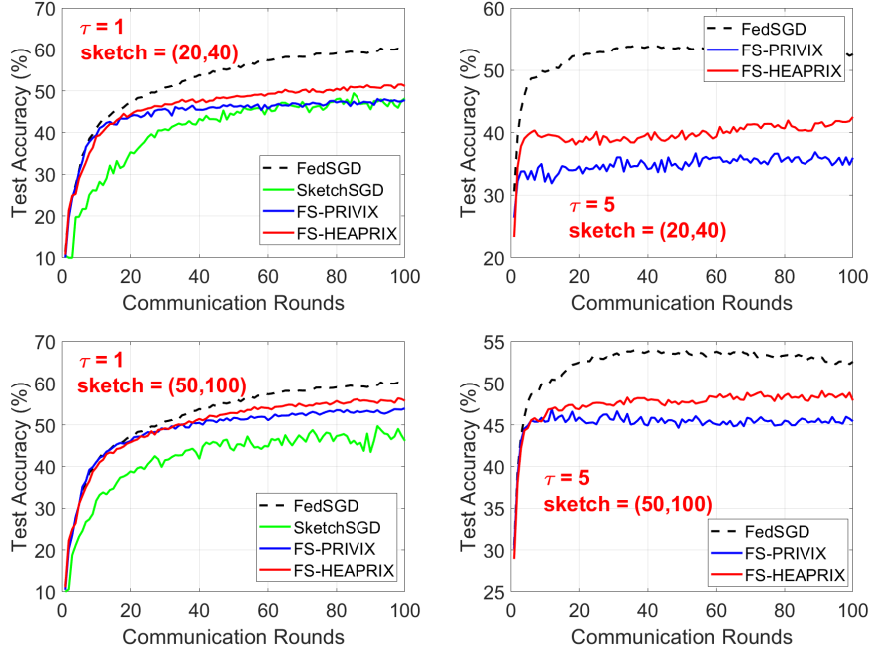


Figure 5: Homogeneous case: CIFAR10: Comparison of compressed optimization methods on LeNet CNN.

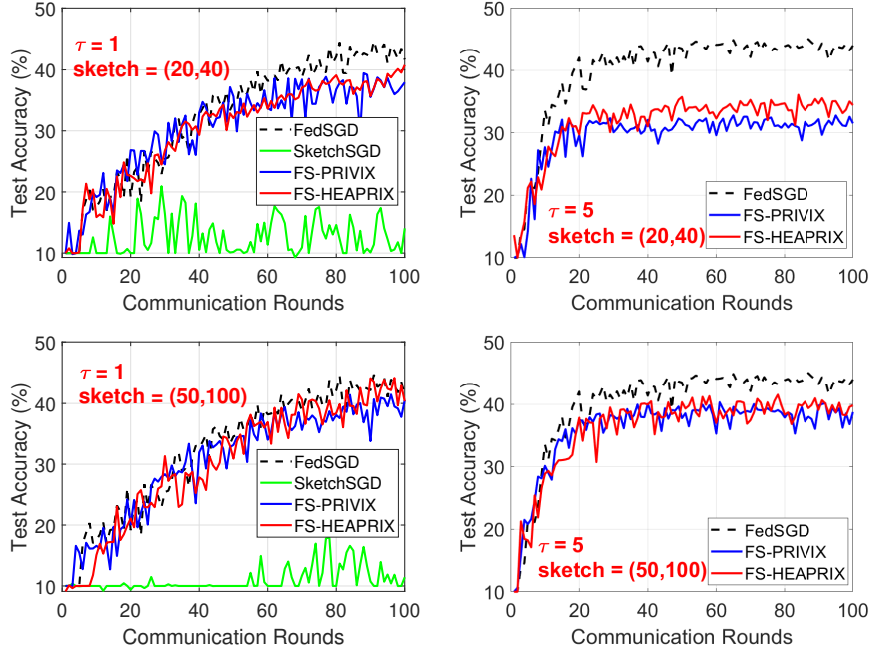


Figure 6: Heterogeneous case: CIFAR10: Comparison of compressed optimization methods on LeNet CNN.