

---

# DPQUANT: Efficient and Differentially-Private Model Training via Dynamic Quantization Scheduling

---

**Yubo Gao**

University of Toronto  
Toronto, Ontario  
ybgao@cs.toronto.edu

**Renbo Tu**

University of Toronto  
Toronto, Ontario  
renbo.tu@mail.utoronto.ca

**Gennady Pekhimenko**

University of Toronto  
Toronto, Ontario  
pekhimenko@cs.toronto.edu

**Nandita Vijaykumar**

University of Toronto  
Toronto, Ontario  
nandita@cs.toronto.edu

## Abstract

Differentially-Private SGD (DP-SGD) is a powerful technique to protect user privacy when using sensitive data to train neural networks. During training, converting model weights and activations into low-precision formats, i.e., *quantization*, can drastically reduce training times, energy consumption, and cost, and is thus a widely used technique. In this work, we demonstrate that quantization causes significantly higher accuracy degradation in DP-SGD compared to regular SGD. We observe that this is caused by noise injection in DP-SGD, which amplifies quantization variance, leading to disproportionately large accuracy degradation. To address this challenge, we present DPQUANT, a dynamic quantization framework that adaptively selects a changing subset of layers to quantize at each epoch. Our method combines two key ideas that effectively reduce quantization variance: (i) *probabilistic sampling* of the layers that rotates which layers are quantized every epoch, and (ii) *loss-aware layer prioritization*, which uses a differentially private loss sensitivity estimator to identify layers that can be quantized with minimal impact on model quality. This estimator consumes a negligible fraction of the overall privacy budget, preserving DP guarantees. Empirical evaluations on ResNet18, ResNet50, and DenseNet121 across a range of datasets demonstrate that DPQUANT consistently outperforms static quantization baselines, achieving near Pareto-optimal accuracy-compute trade-offs and up to  $2.21\times$  theoretical throughput improvements on low-precision hardware, with less than 2% drop in validation accuracy.

## 1 Introduction

Differentially Private Stochastic Gradient Descent (DP-SGD) [1] enables training neural networks on sensitive data while providing formal privacy guarantees. To improve the efficiency of such training on modern hardware, the use of *low-precision* arithmetic and data formats, i.e., *quantization*, has gained widespread interest [18, 23]. Quantization can significantly reduce the amount of computation and memory required, thus reducing the latencies, cost, and energy consumption during training and inference, often with little to no loss in model accuracy [32]. These benefits are especially important in resource-constrained settings such as federated learning with edge devices, where support for full-precision arithmetic is limited and compute budgets are constrained.

Modern accelerators, ranging from datacenter GPUs to mobile NPUs, are rapidly adopting *ultra-low precision formats* such as FP8, INT4, or FP4. NVIDIA’s Blackwell architecture [37] is reported

to provide  $4\times$  throughput for FP4 matrix multiplications compared to FP16; AMD Instinct GPUs supports FP8 [3, 4], and Qualcomm Hexagon supports INT4/INT8 [40]. Leveraging these compute capabilities in model training would enable significant performance and scalability improvements.

In this work, we observe that applying low-precision quantization directly to DP-SGD training often leads to *significant accuracy degradation*, as severe as a 40% drop. While non-DP training is typically robust to quantized training, the gradient clipping and noise addition steps in DP-SGD interact poorly with low-precision arithmetic leading to poor convergence as explained in Section 4.

Our goal is to develop an automatic mechanism to effectively quantize DP-SGD, while minimizing the impact on model accuracy and the differential privacy budget. We observe that *quantizing only a subset of layers* and selectively varying this subset every epoch can preserve most of the efficiency gains from quantization while maintaining model accuracy. DPQUANT uses two core techniques that are implemented in a *differentially private* framework for dynamic quantization scheduling:

1. **Probabilistic layer sampling**, which rotates which layers are quantized every epoch to distribute quantization variance across the network, decreasing overall quantization variance;
2. **Loss-aware prioritization**, which uses a loss sensitivity estimator to selectively quantize layers that have minimal impact on model accuracy.

We make the following contributions:

1. To our knowledge, this work is the first to demonstrate and explain the significant accuracy degradation when employing existing quantization techniques with DP-SGD compared to non-private SGD during model training.
2. We introduce DPQUANT, a differentially private lightweight mechanism that minimizes quantization-induced loss by (a) probabilistically sampling which layers to quantize every epoch and (b) prioritizing layers with lower sensitivity—while incurring only a negligible cost to the overall privacy budget.
3. We demonstrate that DPQUANT achieves near Pareto-optimal accuracy-speed tradeoffs across a range of compute and privacy budgets, outperforming static (fixed-layer) quantization baselines.

## 2 Related Works

**Post-training quantization (PTQ)** [7, 23, 36] aims to accelerate inference through low-precision computations. A neural network is first trained in full-precision and then its weights are quantized. The conversion to quantized formats typically involves a small calibration dataset [21, 35] to allocate quantization bit-widths in different parts of the model and to perform bias correction. PTQ methods are orthogonal to this work since they do not optimize training.

**Quantization-aware training (QAT)** [27] ameliorates the aforementioned accuracy loss by training in lower precision. However, this technique requires extra analysis, such as hardware simulation [47] and sensitivity estimation [14, 38]. QAT also involves quantizer updates, such as step-size tuning [17, 13] and quantizer scaling [41], to select bit-widths [49]. The incurred overhead during training typically cancels out any raw bit-width speedups and often increases wall-clock training time, making them ideal for accelerating inference but not training [8].

**Gradient compression** [29, 6, 5, 48, 42] reduces communication costs by compressing gradients in distributed settings, either through sparsification [44, 52] or low-rank approximation [46, 22]. Notably, [50] combines quantization and the noising mechanism to achieve differential privacy while reducing communication. However, compression does not lower the arithmetic cost of training. In addition, these methods often rely on assumptions such as full gradient availability or error feedback accumulation [26], which are difficult to satisfy under DP constraints.

**Mixed-precision training** [10, 53, 45, 9, 32] aims to reduce training cost by operating on lower-precision data types, e.g., FP16, BF16, FP8, and FP4). While effective for standard SGD, mixed-precision training degrades significantly under DP-SGD. To our knowledge, no prior work has explored mixed-precision training when differential privacy mechanisms are employed.

### 3 Preliminaries

#### 3.1 Differentially-Private DNN Training

We first recall the standard definition of differential privacy:

**Definition 1** (Differential Privacy, [16]). A randomized algorithm  $\mathcal{A}$  satisfies  $(\epsilon, \delta)$ -differential privacy if for all adjacent datasets  $D, D'$  differing on at most one example, and for all measurable sets  $S$  in the output space,

$$\Pr[\mathcal{A}(D) \in S] \leq e^\epsilon \Pr[\mathcal{A}(D') \in S] + \delta.$$

**Definition 2** (DP-SGD, [1]). Differentially Private Stochastic Gradient Descent (DP-SGD) is a variant of SGD that satisfies  $(\epsilon, \delta)$ -differential privacy by clipping and perturbing per-example gradients. At each iteration  $t$ , the update rule is:

$$\theta_{t+1} \leftarrow \theta_t - \eta \left( \frac{1}{|B|} \sum_{i \in B} \text{clip}(\nabla \mathcal{L}(\theta_t, x_i)) + \mathcal{N}(0, \sigma^2 C^2 \mathbf{1}) \right),$$

where  $B$  is a minibatch of training examples,  $\text{clip}(\cdot)$  scales the gradient to have  $\ell_2$  norm at most  $C$ , and  $\mathcal{N}(0, \sigma^2 C^2 \mathbf{1})$  is Gaussian noise added to ensure privacy.

#### 3.2 Quantization and Mixed Precision Training

Modern hardware accelerators such as NVIDIA GPUs, Google TPUs, and Qualcomm Hexagon NPUs provide dedicated support for low-precision arithmetic, including fp16, bfloat16, and increasingly lower bitwidth formats like fp8, fp6, and fp4. These formats enable faster matrix multiplications and convolutions by reducing arithmetic complexity, memory usage, and data transfer costs. Lower precision reduces both the number of transistors required per operation and the bandwidth needed for memory and interconnects, resulting in substantial speedups and energy savings.

While prior work has demonstrated that full training in low-bit formats (e.g., fp4) can retain accuracy under standard SGD, extending these techniques to differentially private training remains challenging. The clipping and noise injection steps in DP-SGD amplify quantization errors and increase gradient variance, making DP training more sensitive to precision loss. Fully quantized DP-SGD thus often results in severe degradation unless carefully tuned.

### 4 Degradation of Quantized DP-SGD

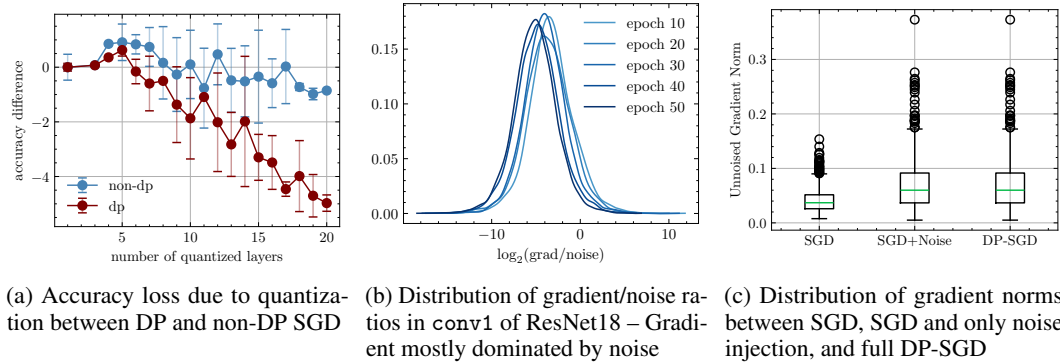


Figure 1: Comparing quantized SGD vs DP-SGD ResNet18 training on the GTSRB dataset

Figure 1 presents a case study of ResNet18 trained on GTSRB, where the forward and backward convolution operators are quantized to evaluate the effects of quantized training. Figure 1a shows the accuracy loss compared to the unquantized baseline for different degrees of quantization (in terms of the number of layers quantized), and the error bars represent the results when different subsets of layers are chosen for quantization, both for DP and non-DP training. For the non-DP SGD baseline, fully-quantized training results in only a modest accuracy drop of around 1%. In contrast, DP-SGD

experiences a much greater degradation, up to 5%. Furthermore, the variance in performance due to different layers being quantized is substantially higher under DP-SGD. We observe similar trends in other neural networks and datasets (included in Appendix A.2).

We hypothesize that the increased sensitivity to quantization can be attributed to noise injection in DP-SGD as follows. In iteration  $t$  in the DP-SGD training, the gradients  $\mathbf{g}_t$  is first clipped to obtain  $\bar{\mathbf{g}}_t$  where  $\|\bar{\mathbf{g}}_t\|_2 \leq C$  (section 3.1). Next, noise  $\mathbf{n}_t \sim \mathcal{N}(0, \sigma^2 C^2 \mathbf{1})$  is sampled, and finally the weights are updated using the sum of the noise and clipped gradient:

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta (\bar{\mathbf{g}}_t + \mathbf{n}_t) \quad (1)$$

We assume the noise scale is  $\sigma \in (0.5, 10)$  [1, 12], in line with common configurations reported in the DP-SGD literature. Since the standard deviation of the injected noise  $\mathbf{n}_t$  is equal to the 2-norm of the clipped gradients, the  $\infty$ -norm of the noise  $\mathbf{n}_t$  (i.e. its largest component) is roughly on the same order as  $\|\bar{\mathbf{g}}_t\|_2$ . Since in higher dimensions  $\|\bar{\mathbf{g}}_t\|_2 \gg \|\bar{\mathbf{g}}_t\|_\infty$  due to the 2-norm growing much faster than the  $\infty$ -norm, combining we have:

$$\|\mathbf{n}_t\|_\infty \approx \|\bar{\mathbf{g}}_t\|_2 \gg \|\bar{\mathbf{g}}_t\|_\infty \quad (2)$$

This relation is also demonstrated empirically in Figure 1b, where on average the magnitude of the clipped gradient elements of  $\bar{\mathbf{g}}$  is  $2^5$  times smaller than that of the injected noise  $\mathbf{n}$ .

The weight update (Equation 1) with noisy gradient updates amplify the norms of raw gradients in subsequent iterations. To show this, we first write the weight update as:

$$\Delta \mathbf{w}_t = \mathbf{w}_{t+1} - \mathbf{w}_t = -\eta (\bar{\mathbf{g}}_t + \mathbf{n}_t) \quad (3)$$

The weight update  $\Delta \mathbf{w}_t \approx \eta \mathbf{n}_t$  due to  $\mathbf{n}_t$  being much larger in norm than  $\bar{\mathbf{g}}_t$ , and thus:

$$\|\Delta \mathbf{w}_t\|_\infty \approx \eta \|\mathbf{n}_t\|_\infty = \mathcal{O}(\|\bar{\mathbf{g}}_t\|_2) \quad (4)$$

where the asymptotic equality holds due to Equation 2. Assuming  $L$ -Lipschitz-smoothness of the loss with respect to the gradients:

$$\|\mathbf{g}_{t+1} - \mathbf{g}_t\|_\infty \leq L \|\mathbf{w}_{t+1} - \mathbf{w}_t\|_\infty = L \|\Delta \mathbf{w}_t\|_\infty \quad (5)$$

We now show that the  $\infty$ -norm of the raw gradients (i.e. before clipping and noising) of the next iteration is bounded by  $\|\bar{\mathbf{g}}_t\|_2$  using the inverse triangle inequality with Equation 4 and 5:

$$\|\mathbf{g}_{t+1}\|_\infty \geq \|\mathbf{g}_{t+1} - \mathbf{g}_t\|_\infty - \|\mathbf{g}_t\|_\infty = \mathcal{O}(\|\bar{\mathbf{g}}_t\|_2) \quad (6)$$

Equation 6 shows that elements of the raw gradients in the next iteration are bound by the much larger  $\mathcal{O}(\|\bar{\mathbf{g}}_t\|_2)$  rather than the usual  $\mathcal{O}(\|\bar{\mathbf{g}}_t\|_\infty)$  in normal SGD. As a result, we expect elements of the raw gradients in DP-SGD to be larger in magnitude than in non-DP training.

To show this empirically, we plot the norms of intermediate gradients under both SGD and DP-SGD using the same hyperparameters in Figure 1c, where the DP-SGD intermediate raw gradients are  $2\times$  larger in the average and worst case. This phenomenon has also been observed in prior work [15], showing a even larger gap than what we observe in gradient norms later in training.

We now demonstrate that the much larger raw gradients under DP-SGD result in much higher quantization variance.

**Proposition 1.** *Let  $q : \mathbb{R}^n \rightarrow \mathbb{R}^n$  be an unbiased (i.e.  $\mathbb{E}[q(\mathbf{x})] = \mathbf{x}$ ) and scale invariant (i.e.  $q(\lambda \mathbf{x}) = \lambda q(\mathbf{x})$ ) quantizer. Assume  $q(\mathbf{x})$  quantizes values onto some finite grid. Let  $\mathbf{x}$  be sampled from a absolutely continuous distribution. Then the quantizer variance  $\text{Var}(q(\mathbf{x})) = \Theta(\|\mathbf{x}\|_\infty^2)$ . Proof: See Appendix A.3.*

Using Prop. 1, we can more precisely express the variance of the quantization as follows:

$$\begin{aligned} \text{(under DP-SGD)} \quad \text{Var}(q(\mathbf{g}_{t+1}) | \mathbf{g}_{t+1}) &= \mathcal{O}(\|\mathbf{g}_{t+1}\|_\infty^2) = \mathcal{O}(\|\bar{\mathbf{g}}_t\|_2^2) \\ \text{(under SGD)} \quad \text{Var}(q(\mathbf{g}_{t+1}) | \mathbf{g}_{t+1}) &= \mathcal{O}(\|\mathbf{g}_t\|_\infty^2) \end{aligned}$$

The quantization variance above is *in addition* to the existing variance of the stochastic gradients, as well as noise injected by DP-SGD. In higher dimensions,  $\|\bar{\mathbf{g}}_t\|_2 \gg \|\bar{\mathbf{g}}_t\|_\infty$ , quantization contributes much more variance to the gradients, hence leads to slower and less reliable convergence [25] and accuracy degradation. To address this challenge, we aim to reduce the quantized-induced variance.

## 5 DPQUANT: Our Proposed Solution

### 5.1 Part I: Probabilistic Layer Sampling

Each time we perform randomized and unbiased quantization on a layer, we introduce additional variance to its gradient updates. While this added variance might be acceptable in standard training, it results in a significant performance degradation under DP-SGD, where gradient stability is already challenged by injected noise.

Suppose a layer is quantized with probability  $p$ , we let  $\mathbf{g}_{\text{fp}}$  to denote its full precision gradients and  $\mathbf{g}_{\text{quant}}$  to be its gradients computed under quantization. By Section 4, quantization incurs additional variance, hence  $\text{Var}(\mathbf{g}_{\text{fp}}) \leq \text{Var}(\mathbf{g}_{\text{quant}})$ . We can write the *expected* gradient variance as:

$$\mathbb{E}(\text{Var}(\mathbf{g})) = (1-p) \text{Var}(\mathbf{g}_{\text{fp}}) + p \text{Var}(\mathbf{g}_{\text{quant}}) \leq \text{Var}(\mathbf{g}_{\text{quant}})$$

From this it follows that whenever  $p < 1$  – that is, when only a subset of layers is quantized at each epoch—the average quantization-induced variance is strictly lower than in full quantization. Furthermore, by rotating which layers are quantized every epoch, no single layer repeatedly incurs the full quantization variance, and hence their expected variance remains smaller than  $\text{Var}(\mathbf{g}_{\text{quant}})$ .

### 5.2 Part II: Loss-Aware Layer Prioritization

Not all layers contribute equally to model performance. Intuitively, we prefer to retain higher precision in layers that have a greater impact on the loss or accuracy. Given a constrained quantization budget, our goal is to prioritize quantization in lower-impact layers, thereby minimizing the overall loss in model quality.

We define a quantization policy  $p$  to be the set of layers to compute under quantization. We define  $R(p)$  as the expected *loss increase* incurred by applying quantization policy  $p$  instead of full precision:

$$R(p) := \mathbb{E}_D[\mathcal{L}(M_p(D)) - \mathcal{L}(M_{\text{fp32}}(D))].$$

Our goal is to find policies that minimally increases loss (or equivalently, a policy  $p$  with small  $R(p)$ ). We note that a key challenge in evaluating  $R(p)$  for a given policy  $p$  is that the expectation is taken over the full private dataset  $D$ . This makes direct computation both expensive and incompatible with tight privacy guarantees. To address this, we instead estimate  $R(p)$  by subsampling  $D$  and running a limited number of DP-SGD iterations under policy  $p$  to obtain a proxy loss, then the same training iterations is done to obtain the baseline full-precision loss, and the difference is used as an empirical estimate of the loss impact.

Since this quantity is computed on the private training dataset  $D$ , any estimation of  $R(\cdot)$  must be performed in a differentially private manner, and therefore consumes part of the overall privacy budget. Any computation using the private data incurs a privacy cost that must be accounted for to ensure that the privacy budget is not exceeded. We outline how DPQUANT privatizes and accounts for loss measurement in Section 5.4.

### 5.3 Dynamic Layer Selection for Quantized DP-SGD Training

Building on the insights from the previous sections, we design a dynamic layer selection strategy for quantized DP-SGD that combines: (i) probabilistic sampling of quantized layers to reduce variance (Section 5.1), and (ii) loss-aware prioritization to preserve performance by avoiding quantization of high-impact layers (Section 5.2).

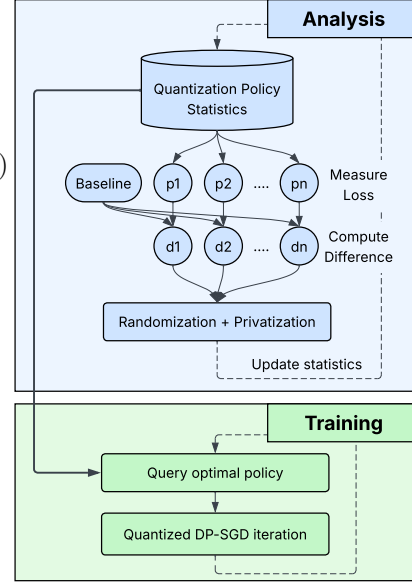


Figure 2: DPQUANT system overview

Table 1: DPQUANT hyperparameters

Parameter	Default	Description
$n$	60	number of epochs to train
$k$	—	layers to execute in low precision.
$n_{\text{sample}}$	1	test samples for loss measurement.
$n_{\text{interval}}$	2	epochs to train before the next measurement.
$R$	2	repetitions during measurement.
$\sigma_{\text{measure}}$	0.5	Noise scale used during loss-difference privatization.
$C_{\text{measure}}$	0.01	Clipping norm used during loss-difference privatization.

Let  $R(l_i)$  denote the estimated quantization loss impact of layer  $i$ . We define the probability of selecting layer  $i$  for quantization as:

$$p_i := \frac{\exp(-\beta R(l_i))}{\sum_{j=1}^n \exp(-\beta R(l_j))}, \quad \text{for } i = 1, \dots, n,$$

where  $\beta > 0$  is a scaling parameter that controls how strongly we prioritize low-impact layers.

As outlined in Figure 2, to quantize  $k$  out of  $n$  layers at each epoch, we sample a subset *without replacement* according to the distribution  $\{p_i\}$ . This allows us to adaptively choose the least sensitive layers for quantization, while still randomly rotating layers with similar loss impact to minimize variance over time. This selection procedure is detailed in Appendix A.8. DPQUANT provides a set of tunable parameters that govern the frequency of the analysis, as well as other privacy parameters.

#### 5.4 Privacy Accounting

Our method begins by measuring loss differences on each user’s private dataset. Specifically, we compute  $\mathcal{L}(M(D))$  which requires inspecting raw data and inherently risks exposing sensitive information if released directly. Without these privacy-preserving measures, simply publishing the loss-difference measurements compromises the privacy guarantee DP-SGD provides.

**Definition 3** (Sampled Gaussian Mechanism (SGM), [33]). *Let  $f$  be a function that maps subsets of a dataset  $S$  to  $\mathbb{R}^d$ . The Sampled Gaussian Mechanism, denoted  $\text{SG}_{q,\sigma}$ , is defined with sampling rate  $0 < q \leq 1$  and noise parameter  $\sigma > 0$  as:*

$$\text{SG}_{q,\sigma}(S) := f(\{x \in S : x \text{ is independently sampled with probability } q\}) + \mathcal{N}(0, \sigma^2 \mathbb{I}^d),$$

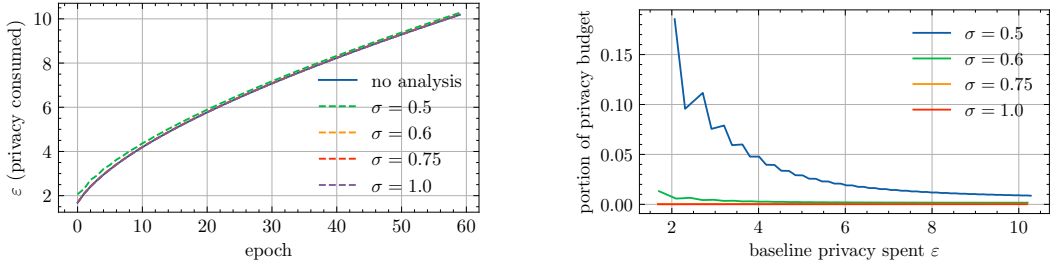
where each element in  $S$  is independently included with probability  $q$ , and  $\mathcal{N}(0, \sigma^2 \mathbb{I}^d)$  denotes  $d$ -dimensional isotropic Gaussian noise with variance  $\sigma^2$  per coordinate.

To protect privacy, we frame this loss computation as a Sampled Gaussian Mechanism (SGM): we draw a random subsample of  $D$ , clip the resulting loss value to bound sensitivity, and then add Gaussian noise of scale  $\sigma$ . These operations correspond to steps 4–5 of Algorithm 1.

**Proposition 2.** *Algorithm 1 is a Sampled Gaussian Mechanism (SGM) with sample rate  $q = n_{\text{sample}}/|B|$  and noise scale  $\sigma = \sigma_{\text{measure}}$ .*

*Proof:* See Appendix A.4.

To account for the privacy cost we incur by performing the analysis in Algorithm 1, we rely on Opacus’s privacy accountants [51]. This is for two reasons: First, these accountants measure the cumulative privacy loss of SGMs [31], where Prop. 2 we can reuse its implementation. Second, by leveraging the advanced composition theorem [1], we obtain a much tighter upper bound on the total privacy expenditure incurred by both the DP-SGD training process and any subsequent analyses performed under the same privacy budget. We explain this in more detail in Appendix A.7.



(a) Total privacy consumption of analysis + training

(b) Fraction of privacy spent on analysis rel. training

Figure 3: Privacy cost of analysis for ResNet18/GTSRB; performing analysis every 2 epochs

In Figure 3, we report the cumulative privacy loss from both the analysis and training components across various configurations. Our results empirically demonstrate that the privacy cost of analysis is negligible compared to training, and does not meaningfully affect the quality of the resulting model.

---

**Algorithm 1** COMPUTELOSSIMPACT

---

```
1: Input:  $P$  (set of quantization policies),  $B$  (batches),  $R$  (iterations),  $\alpha$  (EMA decay),  $\mathcal{C}$  (clipping norm),  $\sigma$  (noise scale),  $L$  (EMA loss differences)
2:  $\ell_{\text{baseline}} \leftarrow 0$ 
3: for  $i = 1$  to  $R$  do
4:   RESTOREMODEL()
5:   for each  $(x, y) \in B$  do
6:     with no quantization do DPSGD-UPDATE( $M$ ,  $\text{loss}(M(x), y)$ )
7:   end for
8:    $\ell_{\text{baseline}} \leftarrow \ell_{\text{baseline}} + \frac{1}{|B|} \sum_{(x,y) \in B} \text{loss}(M(x), y)$  ❶ Compute baseline loss
9: end for
10: for each  $p \in P$  do
11:    $\ell[p] \leftarrow 0$ 
12:   for  $i = 1$  to  $R$  do
13:     RESTOREMODEL()
14:     for each  $(x, y) \in B$  do
15:       with policy  $p$  do DPSGD-UPDATE( $M$ ,  $\text{loss}(M(x), y)$ )
16:     end for
17:      $\ell[p] \leftarrow \ell[p] + \frac{1}{|B|} \sum_{(x,y) \in B} \text{loss}(M(x), y)$  ❷ Compute loss for policy  $p$ 
18:   end for
19: end for
20:  $R[p] \leftarrow \ell[p]/R - \ell_{\text{baseline}}/R$  for all  $p \in P$ 
21:  $\mathbf{R} \leftarrow [R[p_1], \dots, R[p_k]]$ 
22:  $\hat{\mathbf{R}} \leftarrow \mathbf{R} \cdot \min\left(1, \frac{\mathcal{C}}{\|\mathbf{R}\|_2}\right) + \mathcal{N}(0, \sigma^2 \mathbf{C}^2 \mathbf{1})$  ❸ Privatize loss differences
23: UPDATEPRIVACY(rate =  $|B|/|D|$ , steps = 1, noise_scale =  $\sigma$ ) ❹ Update privacy accountant
24: for each  $p \in P$  do
25:    $L[p] \leftarrow (1 - \alpha) \cdot L[p] + \alpha \cdot \hat{\mathbf{R}}[p]$  ❺ Update loss difference EMA
26: end for
27: return  $L$ 
```

---

## 6 Evaluation

**Models and Datasets.** We evaluate our approach on commonly used neural networks for differentially private training [24, 12]: ResNet18 [19], ResNet50 and DenseNet121 [20] from in the torchvision [30] library. These models are trained on the Extended MNIST [11], German Traffic Sign Recognition Benchmark (GTSRB) [43], and CIFAR-10 [28] datasets.

**Hardware.** The experiments were conducted on a server with a Intel(R) Xeon(R) Platinum 8470 CPU, 1TiB of memory, and 8 NVIDIA H100 80GB HBM3 GPUs.

**Implementation.** DPQUANT is implemented on top of Opacus [51], a DP training framework. Opacus provides Poisson sampling, gradient clipping, noising, and gradient accumulation.

**Low Precision Format.** For low precision computations, we used the LUQ-FP4[9] format, the highest-performing 4-bit quantization format. LUQ-FP4 uses a 4-bit representation of floating point numbers. The 4-bit representation consists of 1 sign and 3 exponent bits. LUQ-FP4 first normalizes the values in a tensor, then stochastically quantizes it to one of the 16 levels.

**Measurements.** We use the default parameters in Table 1 unless otherwise noted. We list more detailed training parameters in Appendix A.1.

### 6.1 Quantization-quality trade-off

Quantizing more layers proportionally increases the speed of training. However, it also increases the accuracy degradation in DP-SGD training. Thus, there is a *speed-accuracy* trade-off depending on the number of layers quantized. For a given number of quantized layers, the resulting model accuracy can significantly vary depending on which layers are quantized at any given epoch. DPQUANT aims to automatically identify the subset of layers for each epoch that provides the best accuracy, assuming

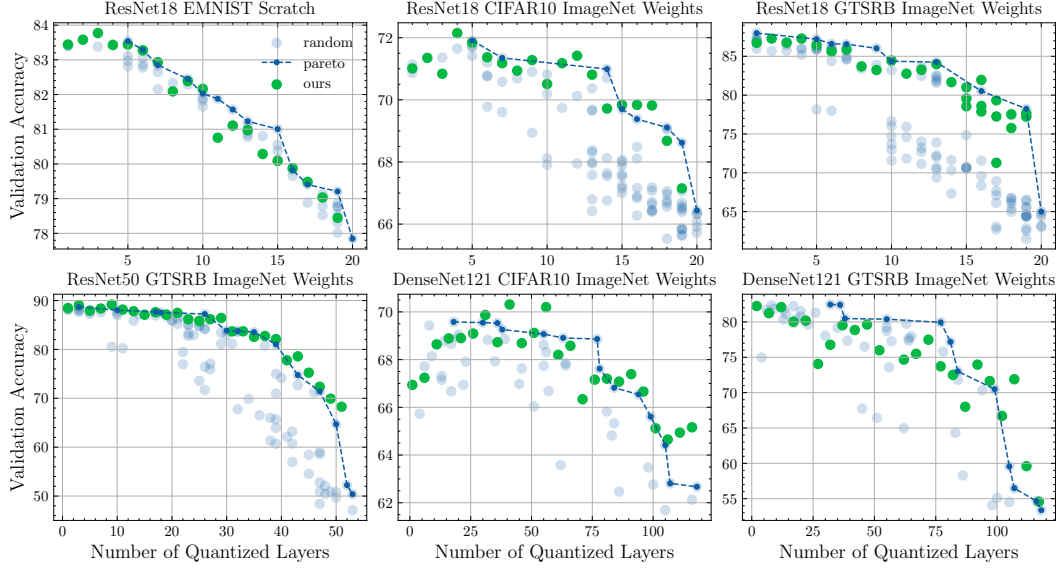


Figure 4: Comparing policies generated by DPQUANT to the speed-accuracy Pareto front

a certain number of layers are quantized. We refer to the desired number of quantized layers as “computational budget” because it determines the speed and compute resources needed.

In Figure 4, we sampled  $\approx 50$  random subsets of layers to execute in fp4. We plotted the empirical Pareto front using these sampled measurements, in addition to the resulting validation accuracy when using DPQUANT’s scheduling technique for a given computational budget.

We make two observations. First, we note that randomly selecting the quantized layers can lead to significant loss in accuracy, as much as 40%. Second, DPQUANT generates scheduling configurations that provide validation accuracy close to the Pareto-front for all evaluated networks and datasets.

## 6.2 Sensitivity to Privacy Budget

Model	Dataset	Percent Quantized	$\epsilon = 4$				$\epsilon = 8$			
			Baseline	$\epsilon$	Ours	$\epsilon$	Baseline	$\epsilon$	Ours	$\epsilon$
ResNet18	EMNIST	0.5	81.27 $\pm$ 1.29	3.14	<b>82.16</b>	3.04	—	—	—	—
		0.75	80.51 $\pm$ 0.37	3.01	80.09	3.04	—	—	—	—
		0.9	78.82 $\pm$ 0.30	3.01	<b>79.03</b>	3.04	—	—	—	—
	GTSRB	0.5	42.34 $\pm$ 5.53	4.01	<b>49.09</b>	3.99	69.06 $\pm$ 5.63	8.01	<b>76.75</b>	7.99
		0.75	39.98 $\pm$ 3.99	4.01	<b>42.62</b>	3.96	63.62 $\pm$ 5.59	8.01	<b>70.07</b>	7.99
		0.9	37.94 $\pm$ 2.23	4.01	<b>39.48</b>	3.99	57.49 $\pm$ 4.46	8.01	<b>67.67</b>	7.99
	CIFAR-10	0.5	64.37 $\pm$ 1.42	4.06	<b>65.39</b>	3.94	69.26 $\pm$ 1.46	7.12	<b>70.51</b>	7.17
		0.75	62.17 $\pm$ 0.61	4.06	<b>63.57</b>	3.94	67.80 $\pm$ 0.81	7.12	<b>69.84</b>	7.17
		0.9	61.09 $\pm$ 1.66	4.06	61.22	3.94	67.21 $\pm$ 1.24	7.12	<b>68.68</b>	7.17
ResNet50	GTSRB	0.5	38.76 $\pm$ 8.16	4.01	<b>42.11</b>	3.99	75.99 $\pm$ 7.33	8.01	<b>80.23</b>	7.99
		0.75	29.48 $\pm$ 4.72	4.01	<b>33.67</b>	3.99	58.13 $\pm$ 8.50	8.01	<b>69.03</b>	7.99
		0.9	24.78 $\pm$ 2.67	4.01	<b>29.00</b>	3.99	47.40 $\pm$ 7.23	8.01	<b>59.87</b>	7.99
DenseNet121	GTSRB <sup>1</sup>	0.5	54.10 $\pm$ 5.58	4.06	55.38	3.97	65.47 $\pm$ 5.42	8.01	<b>71.05</b>	7.93
		0.75	44.60 $\pm$ 5.06	4.06	<b>47.36</b>	3.97	56.14 $\pm$ 7.57	8.01	<b>63.30</b>	7.93
		0.9	40.52 $\pm$ 2.83	4.06	<b>44.15</b>	3.97	51.06 $\pm$ 5.41	8.01	<b>52.60</b>	7.93
	CIFAR-10 <sup>1</sup>	0.5	59.22 $\pm$ 1.15	4.03	<b>61.08</b>	3.97	67.96 $\pm$ 0.93	7.12	<b>68.96</b>	7.28
		0.75	56.43 $\pm$ 1.72	4.03	<b>60.31</b>	3.97	64.81 $\pm$ 1.71	7.12	<b>66.48</b>	7.28
		0.9	55.18 $\pm$ 1.38	4.03	<b>58.89</b>	3.97	63.03 $\pm$ 1.69	7.12	<b>65.13</b>	7.28

Table 2: Model quality across datasets and privacy levels.

<sup>1</sup>Batch size decreased to improve convergence under  $\epsilon = 4$ .



We compared our method to the baseline for two privacy budgets  $\varepsilon = 4$  and  $\varepsilon = 8$ . In Table 2 we plotted the validation accuracy for different privacy budgets. For ResNet18/50, we obtained these values by truncating the training at the respective privacy budgets (i.e. without additional hypermeter tuning), and selected baseline data point with larger  $\varepsilon$  than ours wherever possible.

In most cases, DPQUANT outperforms the baseline performance by at least 1 standard deviation whilst not exceeding the privacy budget. In particular, despite the privacy cost of analysis being more dominating during the  $\varepsilon = 4$  case, DPQUANT produces near-optimal quantization schedules, demonstrating its robustness with respect to  $\varepsilon$ .

### 6.3 Ablation Study

In order to better understand the contributions of the two approaches, we compared our approach (probabilistic layer sampling + loss-aware layer prioritization) with probabilistic layer sampling (PLS) alone. In Figure 5, we observe that PLS consistently performs better than the baseline where the quantized layers are selected statically. However, there is still a large gap between PLS and the *best-performing* layer selections, suggesting that some crucial layers are consistently being subjected to quantization which significantly degrades the quality of trained models.

When PLS is combined with loss-aware layer prioritization, the layers crucial to model training are left in full precision, even when most of the layers are quantized. The benefits of prioritization begins to surface as the proportion of quantized layers increase, as the critical layers have a larger probability of being quantized in the randomized baselines.

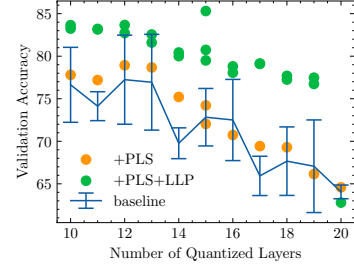


Figure 5: Ablation study, PLS: probabilistic layer selection, LLP: loss-aware layer prioritization

### 6.4 Theoretical Speedup

As hardware with support for FP4 MatMuls and Conv2D (e.g., NVIDIA Blackwell) are not yet widely available, we are unable to evaluate the speed benefits of quantization with DPQUANT. Instead, we use estimates from prior work, along with performance statistics published by NVIDIA [37] to estimate speedups. We estimate that FP4 can provide a  $4\times$  speedup over the FP16 baseline by emulating FP4 computation on existing hardware. Separately, prior works [45, 10, 2] report a  $4 - 7.3\times$  speedup when using FP4 on supported hardware. To remain conservative, we use the lower bound ( $4\times$ ) in our estimates. We assume matrix multiplications, convolutions, and element-wise operations can be accelerated  $4\times$ , and characterize the total runtime as a linear compute cost model:

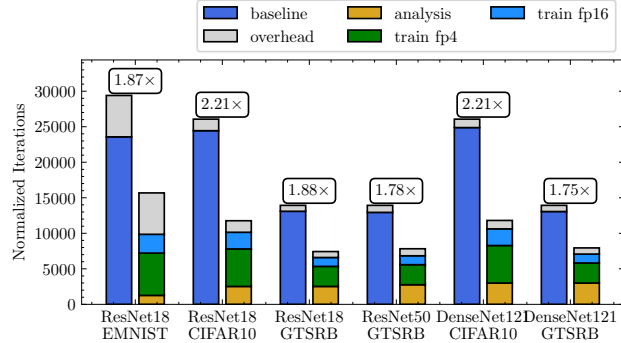


Figure 6: Theoretical speedups for DPQUANT assuming 90% of the layers are quantized.

$$T_{\text{ours}} = T_{\text{analysis}} + (1 - p + p/4)(T_{\text{train baseline}} - T_{\text{overhead}}) + T_{\text{analysis}} + T_{\text{overhead}}$$

where  $T_{\text{analysis}}$  is the time taken by algorithm 1, and  $T_{\text{overhead}}$  captures the time taken by operations that do not have performance benefits from low precision (details in appendix A.6).

We show our speedups in Figure 6. Quantized training with DPQUANT is  $1.75\times$  to  $2.21\times$  faster than the fp16 baseline. In particular, the loss-aware prioritization mechanism in DPQUANT incurs minimal runtime overhead, which is crucial to preserve the performance gains of fp4 computation.

## 7 Conclusion

In this paper, we introduce DPQUANT, a mechanism for efficient quantized DP-SGD training. We make the observation that existing quantized training techniques can significantly degrade the accuracy of models trained with DP-SGD and provided justification which demonstrated the amplified quantization error. To address this challenge, DPQUANT employs techniques to dynamically select layers to quantize such that impact of quantization on model accuracy is minimized. DPQUANT itself is a differentially private mechanism that incurs only small privacy cost. We empirically demonstrate that DPQUANT achieves near-optimal compute-to-accuracy tradeoffs during quantized training, generalizes to different models, datasets and privacy budgets, and can provide up to  $2.21\times$  speedup while minimally impacting accuracy. DPQUANT enables efficient and practical differentially-private training for both centralized and distributed training deployments.

## References

- [1] Martin Abadi, Andy Chu, Ian Goodfellow, H. Brendan McMahan, Ilya Mironov, Kunal Talwar, and Li Zhang. Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security, CCS'16*. ACM, October 2016. doi: 10.1145/2976749.2978318. URL <http://dx.doi.org/10.1145/2976749.2978318>.
- [2] AmirAli Abdolrashidi, Lisa Wang, Shivani Agrawal, Jonathan Malmaud, Oleg Rybakov, Chas Leichner, and Lukasz Lew. Pareto-optimal quantized resnet is mostly 4-bit. In *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, page 3085–3093. IEEE, June 2021. doi: 10.1109/cvprw53098.2021.00345. URL <http://dx.doi.org/10.1109/CVPRW53098.2021.00345>.
- [3] Advanced Micro Devices, Inc. AMD Instinct™ MI300X Accelerator Data Sheet: Leading-Edge Accelerator Module for Generative AI, Training, and High-Performance Computing. Technical report, Advanced Micro Devices, Inc., 2023. URL <https://www.amd.com/content/dam/amd/en/documents/instinct-tech-docs/data-sheets/amd-instinct-mi300x-data-sheet.pdf>. Accessed: 2025-05-13.
- [4] Advanced Micro Devices, Inc. Data types and precision support. <https://rocm.docs.amd.com/en/latest/reference/precision-support.html>, March 2025. ROCm Documentation; Accessed: 2025-05-13.
- [5] Mohammadreza Alimohammadi, Ilia Markov, Elias Frantar, and Dan Alistarh. L-greco: Layerwise-adaptive gradient compression for efficient and accurate deep learning, 2023. URL <https://arxiv.org/abs/2210.17357>.
- [6] Dan Alistarh, Demjan Grubic, Jerry Li, Ryota Tomioka, and Milan Vojnovic. Qsgd: Communication-efficient sgd via gradient quantization and encoding, 2017. URL <https://arxiv.org/abs/1610.02132>.
- [7] Ron Banner, Yury Nahshan, Elad Hoffer, and Daniel Soudry. Post-training 4-bit quantization of convolution networks for rapid-deployment, 2019. URL <https://arxiv.org/abs/1810.05723>.
- [8] Mengzhao Chen, Wenqi Shao, Peng Xu, Jiahao Wang, Peng Gao, Kaipeng Zhang, and Ping Luo. Efficientqat: Efficient quantization-aware training for large language models. *arXiv preprint arXiv:2407.11062*, 2024.
- [9] Brian Chmiel, Ron Banner, Elad Hoffer, Hilla Ben Yaacov, and Daniel Soudry. Accurate neural training with 4-bit matrix multiplications at standard formats, 2024. URL <https://arxiv.org/abs/2112.10769>.
- [10] Jungwook Choi, Zhuo Wang, Swagath Venkataramani, Pierce I-Jen Chuang, Vijayalakshmi Srinivasan, and Kailash Gopalakrishnan. Pact: Parameterized clipping activation for quantized neural networks, 2018. URL <https://arxiv.org/abs/1805.06085>.
- [11] Gregory Cohen, Saeed Afshar, Jonathan Tapson, and André van Schaik. Emnist: an extension of mnist to handwritten letters, 2017. URL <https://arxiv.org/abs/1702.05373>.
- [12] Soham De, Leonard Berrada, Jamie Hayes, Samuel L. Smith, and Borja Balle. Unlocking high-accuracy differentially private image classification through scale, 2022. URL <https://arxiv.org/abs/2204.13650>.
- [13] Xin Ding, Xiaoyu Liu, Zhijun Tu, Yun Zhang, Wei Li, Jie Hu, Hanting Chen, Yehui Tang, Zhiwei Xiong, Baoqun Yin, et al. Cbq: Cross-block quantization for large language models. *arXiv preprint arXiv:2312.07950*, 2023.
- [14] Zhen Dong, Zhewei Yao, Amir Gholami, Michael Mahoney, and Kurt Keutzer. Hawq: Hessian aware quantization of neural networks with mixed-precision, 2019. URL <https://arxiv.org/abs/1905.03696>.
- [15] Jian Du, Song Li, Xiangyi Chen, Siheng Chen, and Mingyi Hong. Dynamic differential-privacy preserving sgd, 2022. URL <https://arxiv.org/abs/2111.00173>.

- [16] Cynthia Dwork and Aaron Roth. The algorithmic foundations of differential privacy. *Foundations and Trends in Theoretical Computer Science*, 9(3–4):211–407, 2014. doi: 10.1561/04000000042. URL <https://www.nowpublishers.com/article/Details/TCS-042>.
- [17] Steven K. Esser, Jeffrey L. McKinstry, Deepika Bablani, Rathinakumar Appuswamy, and Dharmendra S. Modha. Learned step size quantization, 2020. URL <https://arxiv.org/abs/1902.08153>.
- [18] Amir Gholami, Sehoon Kim, Zhen Dong, Zhewei Yao, Michael W. Mahoney, and Kurt Keutzer. A survey of quantization methods for efficient neural network inference, 2021. URL <https://arxiv.org/abs/2103.13630>.
- [19] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition, 2015. URL <https://arxiv.org/abs/1512.03385>.
- [20] Gao Huang, Zhuang Liu, Laurens van der Maaten, and Kilian Q. Weinberger. Densely connected convolutional networks, 2018. URL <https://arxiv.org/abs/1608.06993>.
- [21] Itay Hubara, Yury Nahshan, Yair Hanani, Ron Banner, and Daniel Soudry. Accurate post training quantization with small calibration sets. In *International Conference on Machine Learning*, pages 4466–4475. PMLR, 2021.
- [22] Yerlan Idelbayev and Miguel A Carreira-Perpinán. Low-rank compression of neural nets: Learning the rank of each layer. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8049–8059, 2020.
- [23] Benoit Jacob, Skirmantas Kligys, Bo Chen, Menglong Zhu, Matthew Tang, Andrew Howard, Hartwig Adam, and Dmitry Kalenichenko. Quantization and training of neural networks for efficient integer-arithmetic-only inference, 2017. URL <https://arxiv.org/abs/1712.05877>.
- [24] Matthew Jagielski, Jonathan Ullman, and Alina Oprea. Auditing differentially private machine learning: How private is private sgd?, 2020. URL <https://arxiv.org/abs/2006.07709>.
- [25] Rie Johnson and Tong Zhang. Accelerating stochastic gradient descent using predictive variance reduction. In C.J. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 26. Curran Associates, Inc., 2013. URL [https://proceedings.neurips.cc/paper\\_files/paper/2013/file/ac1dd209cbcc5e5d1c6e28598e8cbb8-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2013/file/ac1dd209cbcc5e5d1c6e28598e8cbb8-Paper.pdf).
- [26] Sai Praneeth Karimireddy, Quentin Rebjock, Sebastian Stich, and Martin Jaggi. Error feedback fixes signsgd and other gradient compression schemes. In *International Conference on Machine Learning*, pages 3252–3261. PMLR, 2019.
- [27] Raghuraman Krishnamoorthi. Quantizing deep convolutional networks for efficient inference: A whitepaper. *arXiv preprint arXiv:1806.08342*, 2018.
- [28] Alex Krizhevsky. Learning Multiple Layers of Features from Tiny Images. Technical report, University of Toronto, April 2009. URL <https://www.cs.toronto.edu/~kriz/learning-features-2009-TR.pdf>.
- [29] Yujun Lin, Song Han, Huizi Mao, Yu Wang, and William J. Dally. Deep gradient compression: Reducing the communication bandwidth for distributed training, 2020. URL <https://arxiv.org/abs/1712.01887>.
- [30] TorchVision maintainers and contributors. Torchvision: Pytorch’s computer vision library. <https://github.com/pytorch/vision>, 2016.
- [31] Mehdi Makni, Kayhan Behdin, Gabriel Afriat, Zheng Xu, Sergei Vassilvitskii, Natalia Ponomareva, Hussein Hazimeh, and Rahul Mazumder. An optimization framework for differentially private sparse fine-tuning, 2025. URL <https://arxiv.org/abs/2503.12822>.

- [32] Paulius Micikevicius, Dusan Stolic, Neil Burgess, Marius Cornea, Pradeep Dubey, Richard Grisenthwaite, Sangwon Ha, Alexander Heinecke, Patrick Judd, John Kamalu, Naveen Mellem-pudi, Stuart Oberman, Mohammad Shoeybi, Michael Siu, and Hao Wu. Fp8 formats for deep learning, 2022. URL <https://arxiv.org/abs/2209.05433>.
- [33] Ilya Mironov, Kunal Talwar, and Li Zhang. Rényi differential privacy of the sampled gaussian mechanism, 2019. URL <https://arxiv.org/abs/1908.10530>.
- [34] Felix Morsbach, Jan Reubold, and Thorsten Strufe. R+r: understanding hyperparameter effects in dp-sgd, 2024. URL <https://arxiv.org/abs/2411.02051>.
- [35] Markus Nagel, Mart van Baalen, Tijmen Blankevoort, and Max Welling. Data-free quantization through weight equalization and bias correction. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 1325–1334, 2019.
- [36] Markus Nagel, Marios Fournarakis, Rana Ali Amjad, Yelysei Bondarenko, Mart Van Baalen, and Tijmen Blankevoort. A white paper on neural network quantization. *arXiv preprint arXiv:2106.08295*, 2021.
- [37] NVIDIA Corporation. Nvidia blackwell architecture technical overview, 2024. URL <https://resources.nvidia.com/en-us-blackwell-architecture>. Accessed: 2025-05-05.
- [38] Eunhyeok Park, Sungjoo Yoo, and Peter Vajda. Value-aware quantization for training and inference of neural networks. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 580–595, 2018.
- [39] Natalia Ponomareva, Hussein Hazimeh, Alex Kurakin, Zheng Xu, Carson Denison, H. Brendan McMahan, Sergei Vassilvitskii, Steve Chien, and Abhradeep Guha Thakurta. How to dp-fy ml: A practical guide to machine learning with differential privacy. *Journal of Artificial Intelligence Research*, 77:1113–1201, July 2023. ISSN 1076-9757. doi: 10.1613/jair.1.14649. URL <http://dx.doi.org/10.1613/jair.1.14649>.
- [40] Qualcomm Technologies, Inc. Ai hardware cores/accelerators, 2024. URL <https://docs.qualcomm.com/bundle/publicresource/topics/80-63195-1/AI-hardware-cores-accelerators.html>. Accessed: 2025-05-05.
- [41] Charbel Sakr, Steve Dai, Rangha Venkatesan, Brian Zimmer, William Dally, and Brucek Khailany. Optimal clipping and magnitude-aware differentiation for improved quantization-aware training. In *International Conference on Machine Learning*, pages 19123–19138. PMLR, 2022.
- [42] Shaohuai Shi, Xianhao Zhou, Shutao Song, Xingyao Wang, Zilin Zhu, Xue Huang, Xinan Jiang, Feihu Zhou, Zhenyu Guo, Liqiang Xie, Rui Lan, Xianbin Ouyang, Yan Zhang, Jieqian Wei, Jing Gong, Weiliang Lin, Ping Gao, Peng Meng, Xiaomin Xu, Chenyang Guo, Bo Yang, Zhibo Chen, Yongjian Wu, and Xiaowen Chu. Towards scalable distributed training of deep learning on public cloud clusters, 2020. URL <https://arxiv.org/abs/2010.10458>.
- [43] Johannes Stalldkamp, Marc Schlipsing, Jan Salmen, and Christian Igel. The german traffic sign recognition benchmark: A multi-class classification competition. In *The 2011 International Joint Conference on Neural Networks*, pages 1453–1460, 2011. doi: 10.1109/IJCNN.2011.6033395.
- [44] Sebastian U. Stich, Jean-Baptiste Cordonnier, and Martin Jaggi. Sparsified sgd with memory, 2018. URL <https://arxiv.org/abs/1809.07599>.
- [45] Xiao Sun, Naigang Wang, Chia-Yu Chen, Jiamin Ni, Ankur Agrawal, Xiaodong Cui, Swagath Venkataramani, Kaoutar El Maghraoui, Vijayalakshmi (Viji) Srinivasan, and Kailash Gopalakrishnan. Ultra-low precision 4-bit training of deep neural networks. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 1796–1807. Curran Associates, Inc., 2020. URL [https://proceedings.neurips.cc/paper\\_files/paper/2020/file/13b919438259814cd5be8cb45877d577-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/13b919438259814cd5be8cb45877d577-Paper.pdf).

- [46] Thijs Vogels, Sai Praneeth Karimireddy, and Martin Jaggi. Powersgd: Practical low-rank gradient compression for distributed optimization. *Advances in Neural Information Processing Systems*, 32, 2019.
- [47] Kuan Wang, Zhijian Liu, Yujun Lin, Ji Lin, and Song Han. Haq: Hardware-aware automated quantization with mixed precision, 2019. URL <https://arxiv.org/abs/1811.08886>.
- [48] Wei Wen, Cong Xu, Feng Yan, Chunpeng Wu, Yandan Wang, Yiran Chen, and Hai Li. Terngrad: Ternary gradients to reduce communication in distributed deep learning, 2017. URL <https://arxiv.org/abs/1705.07878>.
- [49] Jiseok Youn, Jaehun Song, Hyung-Sin Kim, and Saewoong Bahk. Bitwidth-adaptive quantization-aware neural network training: a meta-learning approach. In *European Conference on Computer Vision*, pages 208–224. Springer, 2022.
- [50] Yeojoon Youn, Zihao Hu, Juba Ziani, and Jacob Abernethy. Randomized quantization is all you need for differential privacy in federated learning, 2023. URL <https://arxiv.org/abs/2306.11913>.
- [51] Ashkan Yousefpour, Igor Shilov, Alexandre Sablayrolles, Davide Testuggine, Karthik Prasad, Mani Malek, John Nguyen, Sayan Ghosh, Akash Bharadwaj, Jessica Zhao, Graham Cormode, and Ilya Mironov. Opacus: User-friendly differential privacy library in pytorch, 2022. URL <https://arxiv.org/abs/2109.12298>.
- [52] Xiyu Yu, Tongliang Liu, Xinchao Wang, and Dacheng Tao. On compressing deep models by low rank and sparse decomposition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7370–7379, 2017.
- [53] Shuchang Zhou, Yuxin Wu, Zekun Ni, Xinyu Zhou, He Wen, and Yuheng Zou. Dorefa-net: Training low bitwidth convolutional neural networks with low bitwidth gradients, 2018. URL <https://arxiv.org/abs/1606.06160>.

## A Appendix / supplemental material

### A.1 Training Hyperparameters

While the learning rate might seem too high for regular SGD training, previous results [34, 39] have shown that large learning rates are more beneficial for DP-SGD training.

Table 3: Experimental configurations (6 runs)

	1	2	3	4	5	6
Model	ResNet18	ResNet18	ResNet18	ResNet50	DenseNet121	DenseNet121
Dataset	EMNIST	CIFAR10	GTSRB	GTSRB	CIFAR10	GTSRB
$\sigma$	1	1	1	1	1	1
$\delta$	$10^{-5}$	$10^{-5}$	$10^{-5}$	$10^{-5}$	$10^{-5}$	$10^{-5}$
Clipping norm	1	1	1	1	1	1
Batch size	1024	1024	1024	1024	512	512
Physical batch size	128	128	128	128	128	128
Weights	None	ImageNet	ImageNet	ImageNet	ImageNet	ImageNet
Optimizer	SGD	SGD	SGD	SGD	SGD	SGD
Learning rate (lr)	0.5	0.5	0.5	0.5	0.5	0.5
Epochs	30	60	60	60	60	60

### A.2 Accuracy Degradation for DP-SGD under Naive Quantization

Prior works [45, 9, 32] have demonstrated minimal degradation during quantized fp4/8 training compared to the full precision counterpart. We tabulate their results below:

Table 4: Ultra-Low and LUQ vs. baseline accuracy

Model	Baseline	Ultra-Low [45]	LUQ [9]
ResNet-18	69.7%	68.27% (-1.43%)	69.09% (-0.61%)
ResNet-50	76.5%	74.01% (-2.49%)	75.42% (-1.08%)
MobileNet-V2	71.9%	68.85% (-3.05%)	69.55% (-2.35%)
ResNext50	77.6%	N/A	76.02% (-1.58%)
Transformer-base	27.5 (BLEU)	25.4 (-2.10)	27.17 (-0.33)
BERT fine-tune	87.03 (F1)	N/A	85.75 (-1.28)

As demonstrated in the Figure 4, the performance degradation of DP-SGD under quantization is much larger.

Table 5: Validation accuracy for DP-SGD training: baseline vs. LUQ-FP4 (all layers quantized)

Model	Dataset	Pretraining	Baseline	LUQ-FP4	$\Delta$
ResNet-18	EMNIST	None	83.4%	77.8%	-5.6%
ResNet-18	CIFAR-10	ImageNet	71.0%	65.8%	-5.2%
ResNet-18	GTSRB	ImageNet	85.6%	64.0%	-21.6%
ResNet-50	GTSRB	ImageNet	89.8%	49.0%	-40.8%
DenseNet-121	CIFAR-10	ImageNet	67.0%	62.9%	-4.1%
DenseNet-121	GTSRB	ImageNet	82.0%	53.0%	-29.0%

### A.3 Proof of Proposition 1

**Proposition 1.** *Let  $q : \mathbb{R}^n \rightarrow \mathbb{R}^n$  be an unbiased ( $\mathbb{E}[q(\mathbf{x})] = \mathbf{x}$ ) and scale-invariant ( $q(\lambda\mathbf{x}) = \lambda q(\mathbf{x})$ ) quantizer whose outputs lie on a fixed finite grid. If  $\mathbf{x}$  is drawn from an absolutely continuous distribution, then*

$$\text{Var}(q(\mathbf{x})) = \Theta(\|\mathbf{x}\|_\infty^2).$$

*Proof.* We begin by first showing the upper-bound:  $\text{Var}(q(\mathbf{x})) = \mathcal{O}(\|\mathbf{x}\|_\infty^2)$ . We define  $M = \|\mathbf{x}\|_\infty$  and  $\mathbf{v} = \mathbf{x}/M$ , so  $\|\mathbf{v}\|_\infty = 1$ . By scale-invariance of  $q$ ,

$$\text{Var}(q(\mathbf{x})) = \text{Var}(q(M\mathbf{v})) = \text{Var}(M q(\mathbf{v})) = M^2 \text{Var}(q(\mathbf{v})).$$

Since  $q(\mathbf{v}) \in [-1, 1]^n$ , there exists a finite  $C$  such that  $\text{Var}(q(\mathbf{v})) \leq C$ , giving

$$\text{Var}(q(\mathbf{x})) = M^2 \text{Var}(q(\mathbf{v})) \leq C M^2 = C \|\mathbf{x}\|_\infty^2.$$

Next, we show the lower-bound:  $\text{Var}(q(\mathbf{x})) = \Omega(\|\mathbf{x}\|_\infty^2)$ . On the compact set  $\{\mathbf{v} : \|\mathbf{v}\|_\infty = 1\}$ , the continuous function  $\mathbf{v} \mapsto \text{Var}(q(\mathbf{v}))$  attains a minimum  $m \geq 0$ . Because the finite quantizer grid has measure-zero, absolute continuity of  $\mathbf{x}$  ensures the probability of  $\mathbf{v}$  landing on the grid is 0, so  $\text{Var}(q(\mathbf{v})) > 0$  and hence  $m > 0$ . Therefore

$$m M^2 \leq \text{Var}(q(\mathbf{x})) \leq C M^2,$$

From this we conclude  $\text{Var}(q(\mathbf{x})) = \Theta(\|\mathbf{x}\|_\infty^2)$ , as desired.  $\square$

#### A.4 Justification of Privacy Guarantees (Theorem 2)

We issue a correction of the Proposition 2, listed below (with the difference **highlighted in red**).

**Proposition 2.** *Algorithm 1 is a Sampled Gaussian Mechanism (SGM) with sample rate  $q = |B|/|D|$  and noise scale  $\sigma = \sigma_{\text{measure}}$ .*

*Proof.* We first characterize Algorithm 1 as an analysis on the user’s private dataset. The function accepts a subsampled batch of size  $|B|$  from a dataset with  $|D|$  samples.

Using this batch of user data, as well as some non-private sources of data such as the model weights, we compute the loss differences which is vectorized in  $\mathbf{R} \in \mathbb{R}^p$ , where  $p$  is the number of available quantization policies.

In step (3) of Algorithm 1, we clip the vector  $\mathbf{R}$  to norm  $\mathcal{C}$ , to which independent Gaussian noise proportional to  $\sigma^2 \mathcal{C}^2$  is added to obtain  $\hat{\mathbf{R}}$ . This is equivalent to adding noise proportional to  $\sigma^2$  when the sensitivity of  $\hat{\mathbf{R}}$  is 1 through a scaling argument.

Algorithm 1 ceases to access private data in  $B$  after the computation of  $\hat{\mathbf{R}}$ , which results in all the following steps (i.e. updating privacy accountant and EMA) post-processing [16] which does not impact the privacy consumption.

Furthermore, the privacy accounting step makes use of Opacus’ [51] privacy accountant, which assumes<sup>2</sup> the the noise scale  $\sigma$  is proportional to the clipping constant (i.e. equivalent to adding a noise proportional to  $\sigma^2 \mathcal{C}^2$ ).  $\square$

#### A.5 Low Precision Simulation Setup

As FP4 hardware support is forthcoming, we employ the following simulation setup to emulate the effect of training under FP4. Notably, we quantize both inputs to the conv2d forward, wgrad, and dgrad operators as well as its output.

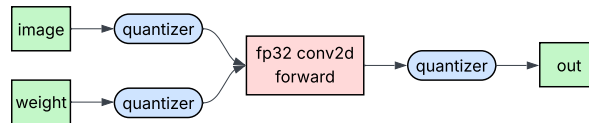


Figure 7: Quantization simulation setup

<sup>2</sup>This assumption is stated in <https://github.com/pytorch/opacus/blob/main/opacus/accountants/analysis/rdp.py>



## A.6 Theoretical Speedup Calculation

Due to the unavailability of accelerators and reliable software support for FP4, we instead rely on a performance model to estimate the theoretical throughputs of FP4 computation.

We first decompose the DP-SGD training computation into the following operations, listed in table 6.

We performed profiling on the models (on their respective datasets) stated in the paper, and we plot the runtime decomposition in Figure 8. Using this data, we can compute the amount of “overhead” (i.e. the time spent on operators which will not benefit from lower precision) for each model/dataset. This is tabulated in Table 7.

Table 6: Decomposition of DP-SGD training

Computation Stage	Description	Benefits from FP4
<i>Total Forward</i>	Time spent on the forward pass through the model, where input data is processed layer by layer to produce the output.	✓
<i>Total Backward</i>	Time for backpropagation, where gradients are calculated for model parameter updates.	✓
<i>Optimizer Clip</i>	Time for clipping gradients to a predefined threshold to ensure stability and prevent large updates during training.	✓
<i>Optimizer Noise</i>	Time for adding random noise to the gradients to ensure differential privacy by masking individual data point contributions.	
<i>Optimizer Scale</i>	Time for scaling the gradients after clipping to adjust the magnitude of the updates.	✓
<i>Other Optimizer</i>	Time spent on other optimizer-related operations, such as learning rate management.	
<i>Other Time</i>	Time for all other operations during the training iteration, including data loading, synchronization, and auxiliary tasks.	

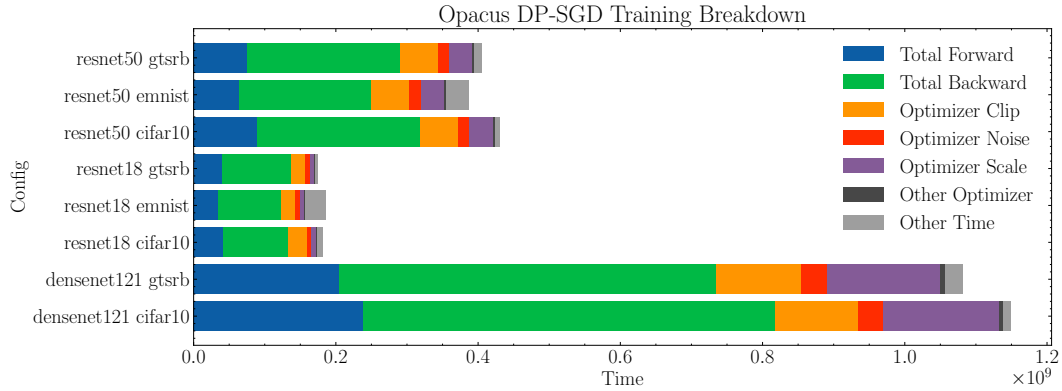


Figure 8: Runtime decomposition of DP-SGD training

Table 7: Breakdown of total time, good ops, bad ops, and overhead percentage for different model configurations.

Config	Total Time	Ops with Speedup	Overhead Ops	Overhead %
DenseNet121 CIFAR10	$1.15 \times 10^9$	$1.10 \times 10^9$	$5.23 \times 10^7$	4.55
DenseNet121 GTSRB	$1.08 \times 10^9$	$1.01 \times 10^9$	$6.74 \times 10^7$	6.23
ResNet18 CIFAR10	$1.82 \times 10^8$	$1.66 \times 10^8$	$1.68 \times 10^7$	9.20
ResNet18 EMNIST	$1.86 \times 10^8$	$1.49 \times 10^8$	$3.68 \times 10^7$	19.81
ResNet18 GTSRB	$1.74 \times 10^8$	$1.63 \times 10^8$	$1.04 \times 10^7$	5.99
ResNet50 CIFAR10	$4.31 \times 10^8$	$4.05 \times 10^8$	$2.55 \times 10^7$	5.92
ResNet50 EMNIST	$3.88 \times 10^8$	$3.36 \times 10^8$	$5.13 \times 10^7$	13.22
ResNet50 GTSRB	$4.05 \times 10^8$	$3.76 \times 10^8$	$2.87 \times 10^7$	7.10

### A.7 Opacus Privacy Accounting

Opacus maintains a tuple of the form (sample rate, noise scale, number of steps), which is incrementally updated during training. At any point, we can query the current privacy cost in terms of  $(\epsilon, \delta)$  by specifying a target  $\delta$  and using either the `rdp` or `prv` accountant. This mechanism enables flexible and precise tracking of privacy usage, allowing us to assess how much additional privacy is consumed by our analysis relative to standard training.

### A.8 Sampling of Quantized Layers

We outline the algorithm DPQUANT uses to select layers to compute under quantization.

---

#### Algorithm 2 SELECTTARGETS

---

```

1: Input:  $L$  (EMA scores),  $P$  (set of policies),  $s$  (temperature),  $m$  (number to sample), layers
   (set of layers to quantize under policy  $p$ )
2:  $v \leftarrow [L[p]]$  for  $p \in P$ 
3:  $v \leftarrow (v - \min(v)) / (\max(v) - \min(v))$  ▷ Normalize
4:  $\pi \leftarrow \text{softmax}(-s \cdot v)$ 
5:  $Q \leftarrow \text{Multinomial}(\pi, m, \text{without replacement})$  ▷ Sample  $m$  policies
6:  $S \leftarrow \emptyset$ 
7: for each  $p \in Q$  do
8:    $S \leftarrow S \cup \text{layers}[p]$ 
9: end for
10: return  $S$ 

```

---