

# REVISITING PRIVACY, UTILITY, AND EFFICIENCY TRADE-OFFS WHEN FINE-TUNING LARGE LANGUAGE MODELS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

We study the inherent trade-offs in *minimizing privacy risks and maximizing utility, while maintaining high computational efficiency*, when fine-tuning large language models (LLMs). A number of recent works in privacy research have attempted to mitigate privacy risks posed by memorizing fine-tuning data by using differentially private training methods (e.g., DP-SGD), albeit at a significantly higher computational cost (inefficiency). In parallel, several works in systems research have focused on developing (parameter) efficient fine-tuning methods (e.g., LoRA). However, few works, if any, investigated whether such efficient methods, in isolation, enhance or diminish privacy risks.

In this paper, we investigate this gap and arrive at a surprising conclusion: efficient fine-tuning methods like LoRA mitigate privacy-risks similar to private fine-tuning methods like DP-SGD. Our empirical finding contradicts the prevailing wisdom that privacy and efficiency objectives are at odds during fine-tuning. Our finding is established by (a) carefully defining measures of privacy and utility that distinguish between recollecting *sensitive and non-sensitive* tokens in training and test datasets used in fine-tuning and (b) extensive evaluations using multiple open-source language models from Pythia, Gemma, Llama, and Qwen families and different domain-specific datasets.

## 1 INTRODUCTION

Large language models (LLMs) have shown proficiency across diverse natural language tasks (Naveed et al., 2023), finding applications in education (Wang et al., 2024), medical chatbots (Thirunavukarasu et al., 2023), and AI assistants (Dong et al., 2023). Their capabilities stem from a two-phase process: (a) pre-training on extensive web data (Kaplan et al., 2020) to develop general language understanding (Brown, 2020), (b) fine-tuning on domain-specific data for specific tasks (Zhang et al., 2023). In both phases, the key challenges involve enhancing *privacy* and *efficiency* while maintaining the models’ *utility*. Privacy is related to reducing the risk of LLMs leaking sensitive user information contained in the training data, efficiency is related to reducing the computational cost of training, while utility is related to generative performance over test data.

A long line of recent research in the privacy community has focussed on methods to mitigate privacy risks when training LLMs (Shi et al., 2022; Zhao et al., 2022; Abadi et al., 2016; Yu et al., 2022). A notable example of such methods is differential privacy based stochastic gradient descent (DP-SGD)(Abadi et al., 2016). Simultaneously, a flurry of recent research in the systems community has explored parameter-efficient fine-tuning methods in LLMs (Han et al., 2024). A notable example of this class of methods is low-rank adaptation (LoRA) (Hu et al., 2022). However, no prior works, to the best of our knowledge, have investigated the privacy risks associated with the efficient training methods.

The central question driving our research here is: *do efficient fine-tuning methods enhance or mitigate privacy risks during training?* As DP-SGD incurs significant additional computational overhead (Dupuy et al., 2022), while LoRA significantly reduces the computational costs, the answer to the above question can have significant consequences for achieving good privacy-efficiency-utility tradeoffs when fine-tuning. For instance, if LoRA mitigates privacy risks of training that would

054 suggest that it can simultaneously achieve both privacy and efficiency objectives, contradicting the  
 055 conventional wisdom drawn from DP literature that privacy comes at a computational cost. A key  
 056 (surprising) finding of our work lies in establishing that *LoRA does indeed mitigate privacy risks*.  
 057

058 A conjecture that might explain our finding is rooted in the following high-level observations about  
 059 DP-SGD and LoRA: methodologically, both DP-SGD and LoRA restrict the impact that training  
 060 examples can have on model parameters – DP-SGD deliberately through its noisy gradient update,  
 061 and LoRA through low-rank adaptation. We formalize this intuition in this work.

062 When attempting to answer the above question, we encountered a more foundational question: *how*  
 063 *should one quantify such privacy risks associated with a fine-tuning method, so that it allows for*  
 064 *a performance comparison across different methods?* Numerous studies have highlighted privacy  
 065 risks in LLMs due to their tendency to memorize and regurgitate training data containing sensitive  
 066 personally identifiable information (PII) such as names, emails, and credentials (Mattern et al.,  
 067 2023; Fu et al., 2023; Kaneko et al., 2024; Carlini et al., 2021; Mireshghallah et al., 2022b; Panda  
 068 et al., 2024). A natural way to account for privacy risks from memorization might be to measure  
 069 loss on recollecting tokens in training data sequences.

070 However, we find that LLMs exhibit very different losses in recollecting sensitive vs. non-sensitive  
 071 tokens in training data (see Figure 1). This difference is due to inherent randomness and unpre-  
 072 dictability of sensitive data (e.g., phone numbers, SSNs) compared to non-sensitive data (e.g., “The  
 073 dog chases the \_”), which is often more structured and predictable. Consequently, we propose a new  
 074 privacy measure that explicitly account for this difference: aiming for high loss on sensitive tokens  
 075 from training data.

076 **Contributions.** We summarize our main find-  
 077 ings:

078 1. *Quantifying privacy and utility, when train-*  
 079 *ing LLMs:* We conceptually argue and empir-  
 080 ically demonstrate that LLM’s ability to recol-  
 081 lect (predict) sensitive and non-sensitive data in  
 082 training (test) datasets are so starkly different  
 083 that we need to account for them when quanti-  
 084 fying privacy and utility. Privacy is best cap-  
 085 tured by a model’s ability to *recollect sensi-*  
 086 *tive tokens in training data*, while utility is best  
 087 captured by the model’s ability to *predict non-*  
 088 *sensitive tokens in test data*.

089 2. *Comparing privacy-utility-efficiency trade-*  
 090 *offs for three different fine-tuning methods:* Our  
 091 measures allow us to conduct a systematic  
 092 and extensive empirical study of three different  
 093 fine-tuning methods: full fine-tuning (FFT), DP-SGD, and LoRA, using models from four different  
 094 LLM families, Pythia (Biderman et al., 2023), Gemma (Team et al., 2024), Llama (Touvron et al.,  
 095 2023) and Qwen (Yang et al., 2024) over two datasets. Our comparative study yields several in-  
 096 teresting insights: we find that – FFT results in poor utility-privacy tradeoffs; DP offers reasonable  
 097 utility-privacy tradeoffs, but is computationally very expensive; LoRA is almost on par with DP in  
 098 terms of privacy with good privacy-utility tradeoffs, but is more computationally efficient. To re-  
 099 confirm our results, we evaluate the three methods using existing privacy measures—*privacy loss*  
 100 (Abadi et al., 2016) and *canary exposure* (Carlini et al., 2019). We also show formally that the effect  
 101 of fine-tuning a model using LoRA and DP on a datapoint is analogous, leading to privacy benefits.

102 3. *Feasibility of achieving all the three privacy, utility and efficiency objectives simultaneously:*  
 103 Our findings about LoRA performance challenge prevailing wisdom that enhancing privacy during  
 104 training is more computationally expensive. This calls for investigating privacy benefits of existing  
 105 and new parameter efficient fine-tuning methods.

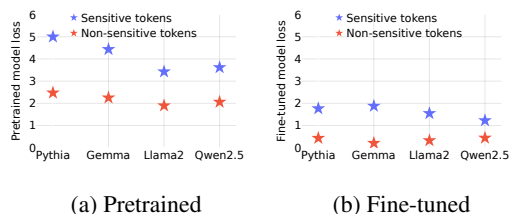


Figure 1: Sensitive and non-sensitive tokens have different predictability, measured as the recollection loss by pre-trained models (Figure 1a) and fine-tuned models (Figure 1b). This distinction motivates to quantify privacy using sensitive tokens from the training data.

## 2 BACKGROUND & RELATED WORK

Privacy in large language models (LLMs) is threatened by membership inference attacks (MIAs) and data extraction attacks (DEAs) (Kibriya et al., 2024; Yao et al., 2024; Mireshghallah et al., 2022a; Mattern et al., 2023; Fu et al., 2023; Kaneko et al., 2024; Carlini et al., 2021; Mireshghallah et al., 2022b; Panda et al., 2024), with DEAs posing a more practical risk (Zhang et al., 2024). Differential privacy (DP-SGD) mitigates these risks (Abadi et al., 2016) but is computationally expensive and sometimes vulnerable (Hayes et al., 2023). Parameter-efficient fine-tuning (PEFT) methods (Han et al., 2024) like LoRA (Hu et al., 2022) provide efficiency and utility, though their privacy impact is less explored. We find that LoRA achieves privacy comparable to DP. Unlike prior approaches combining DP-SGD with LoRA (Liu et al., 2023; Yu et al., 2021; Ma et al., 2024), we show all three objectives—privacy, utility, and efficiency—can be achieved with LoRA alone. An extended version of the related work is in Appendix C.

## 3 EXPERIMENTAL SETUP AND METHODOLOGY

In this work, we want to answer the following questions for three different fine-tuning methods – full fine-tuning (FFT), Differential Privacy (DP-SGD), and Low-Rank Adaptation (LoRA):

- *Privacy: How prone is each method to recollecting the sensitive parts of the training data?*
- *Utility : How effective is each method at predicting non-sensitive parts of test data?*
- *Efficiency: What is the computational cost associated with each method?*

We begin by outlining our experimental setup and motivating our definitions of privacy and utility. In Section 4, we study the privacy–utility trade-off in three fine-tuning methods through an exhaustive hyperparameter search. Using the best configurations identified, we then compare the methods via Pareto-optimal trade-offs in Section 5 under both our proposed and existing privacy measures, ensuring a fair evaluation.

We train models from four different families, Pythia-1B (Biderman et al., 2023), Gemma-2B (Team et al., 2024), Llama2-7B (Touvron et al., 2023), and Qwen2.5-7B (Yang et al., 2024) for 50 epochs, on two publicly available datasets, CustomerSim and SynBio. CustomerSim (Shi et al., 2022) is a simulated dialog dataset for conversation with 10k samples, and SynBio (originally called PII) (Holmes et al., 2024) is an LLM generated dataset with student biographies containing PII with  $\sim 5k$  samples. We use 80% for training and 20% for evaluation. Table 1 in Appendix D shows excerpts from the datasets. Details on libraries and hyperparameters are provided in Appendix E.

In order to measure privacy and utility, we would need to distinguish between the sensitive and non-sensitive tokens. We define sensitive data as per GDPR Article 4(1) (EU, 2016) which says that it can be any personal data or identifier like name, ID number, or location. Owing to the large scale of data, we leverage two tools for annotating the sensitive information: Presidio (Mendels et al., 2018), which helps in a fast identification of private entities in text, and GPT-4 (Achiam et al., 2023), which is provided with a particular prompt for returning the annotated portions. An example of such a prompt for annotation is provided in Appendix F.

To confirm the reliability of the annotations at scale, we run a survey among 40 Prolific<sup>1</sup> crowdworkers with a random subset of 100 samples. We provide the details of the survey in Appendix G that depicts that  $\sim 75\%$  participants found the GPT-4 annotations to be accurate while Presidio annotations were mostly mixed or under-annotated. Hence, the rest of the paper has our results using GPT-4 annotations, and those using Presidio annotations are shown in Appendix I.

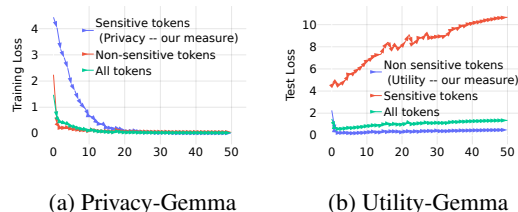


Figure 2: Our measures offer a precise assessment of privacy and utility by distinguishing between sensitive and non-sensitive tokens.

<sup>1</sup><https://www.prolific.com>

Before diving into the tradeoffs from different fine-tuning methods, we want to demonstrate the usefulness of our privacy and utility measures. We quantify **privacy** as the *recollection of sensitive entities in the training data* and **utility** as the *prediction of non-sensitive entities in the test data*.

We present evidence supporting the importance of distinguishing between sensitive and non-sensitive entities in natural language.

**Results.** Figure 2 shows privacy (left) and utility (right) during fine-tuning of the Gemma model on the Customersim dataset (Shi et al., 2022), with training loss on the left column and test loss on the right (results for other models in Appendix H).

In Figure 2a, when no distinction is made between sensitive and non-sensitive tokens, training loss (green curve) starts lower, falsely suggesting strong memorization and higher privacy leakage. Separating losses reveals that *sensitive tokens (blue curve) show much higher loss than non-sensitive ones (red curve)*, especially early in training, since they are less predictable. Similarly, in Figure 2b, test loss without distinction (green curve) appears higher, giving the impression of poor generalization. However, separating losses shows non-sensitive tokens (blue curve) achieve lower loss, confirming better utility. Using these measures, we observe the tradeoffs among different fine-tuning methods.

#### 4 PRIVACY-UTILITY-EFFICIENCY INTERPLAY

In this section, we study the impact of privacy, utility, and efficiency on three different fine-tuning methods - full fine-tuning, differential privacy and LoRA. We conduct an extensive hyperparameter search for both differential privacy and LoRA, systematically tuning their respective parameters to explore the privacy-utility trade-off in depth.

**Update rules:** For each fine-tuning method, we describe how updated weights  $W_{t+1}$  are computed from the previous weights  $W_t$  in each step, where  $W_0$  are the weights of the pre-trained base model before fine-tuning. We use  $X$  to refer to a batch of  $|X|$  datapoints and  $x_i$  to refer to individual datapoint,  $\mathcal{M}_W$  to refer to the model parameterized by weights  $W$ , and  $\mathcal{L}(\mathcal{M}_W(X), X)$  to refer to the autoregressive cross-entropy loss of the model on data  $X$ . We denote by  $\nabla_W \mathcal{L}(\dots)$  as the gradient of the loss wrt. weights  $W$  and  $\eta$  is the learning rate.

**Efficiency:** The efficiency of each method is determined by the amount of computation it requires, and also other factors such as memory requirements. Following (Kaplan et al., 2020), we estimate the amount of training compute ( $C$ ) in floating point operations (FLOPs) for full fine-tuning as  $C_{FFT} = 6DN$ , where  $D$  is the number of training tokens and  $N$ , the number of model parameters. For each method, we report its compute requirements relative to the FFT-baseline based on measurements using the PyTorch profiler<sup>2</sup>.

##### 4.1 FULL FINE-TUNING FOR UTILITY COSTS PRIVACY

**Update rules:** Full fine-tuning (FFT) updates all model parameters at each step:

$$W_{t+1} = W_t - \eta \nabla_{W_t} \mathcal{L}(\mathcal{M}_{W_t}(X), X) \tag{1}$$

**Privacy-Utility trade-off:** Figures 3a and 3b show the *privacy-utility trade-off* for the CustomerSim and SynBio datasets, respectively. *Privacy increases with increasing training loss on sensitive tokens (up  $\uparrow$  on the y-axis)*, while *utility increases with decreasing test loss on non-sensitive tokens (right  $\implies$  on the x-axis)*. Each curve starts with the baseline performance of the pre-trained model.

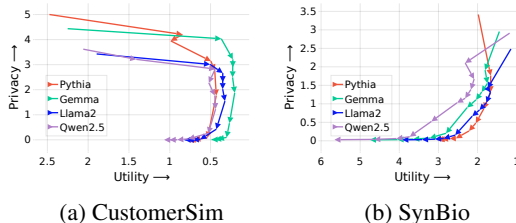


Figure 3: Privacy-utility trade-off shows that privacy increases with higher training loss on sensitive tokens, while utility improves with lower test loss on non-sensitive tokens, enabling desired checkpoint selection to balance both objectives.

<sup>2</sup><https://pytorch.org/docs/stable/profiler>

For CustomerSim (Figure 3a), as training advances (*denoted by an arrow  $\rightarrow$  on the lines*), privacy progressively decreases (*lower on the y-axis*), while utility improves (*rightward on the x-axis*) for approximately the first 5 epochs across all models before stabilizing and eventually declining (*leftward on the x-axis*). However, for SynBio (Figure 3b), the privacy-utility trade-off primarily worsens for Gemma and Llama models. On examining these curves, one can select a desired checkpoint that aligns with specified privacy and utility thresholds.

**Efficiency:** FFT serves as our efficiency baseline. It has moderate compute requirements (discussed above), and relatively high memory requirements, since in addition to the input-dependent activations, we need to keep four numbers per model parameter in GPU memory: the parameter value, its gradient, and two optimizer states (first and second moments for Adam (Kingma, 2014)).

**Takeaway:** FFT offers *poor privacy-utility* trade-offs, since gains in utility in most cases come at the cost of a significant loss in privacy. During FFT, models learn to predict the training distribution better, but also quickly learn to recollect sensitive tokens. FFT is *moderately efficient* and has relatively high memory requirements.

## 4.2 DP-SGD FINE-TUNING FOR PRIVACY AND UTILITY COSTS EFFICIENCY

Differential Privacy (DP-SGD) algorithms (Abadi et al., 2016) aim to safeguard the privacy of individual training data points by limiting their influence on the gradient updates during training. While DP-SGD is designed for MIAs, we show its privacy benefits even in the case of DEAs.

**Update rules:** DP-SGD clips the  $l_2$  norm of each data point’s gradient at a threshold  $T$ , followed by adding noise with magnitude  $\sigma$  to each clipped gradient. Clipping helps to reduce data sensitivity by ensuring that data points with high gradients have limited impact on the parameters. Adding noise further obscures the contribution, making it difficult to infer specific data points from the model.

$$W_{t+1} = W_t - \eta \text{Noise} \left( \frac{1}{B} \sum_i \text{Clip}(\nabla_{W_t} \mathcal{L}(\mathcal{M}_{W_t}(x_i), x_i)) \right); \text{Clip}(y) = y / \max \left( 1, \frac{\|y\|_2}{T} \right)$$

$$\text{Noise}(y) = y + \mathcal{N}(0, \sigma^2 T^2 \mathbf{1}) \quad (2)$$

Between the two hyperparameters,  $\sigma$  and  $T$ , we vary the noise hyperparameter  $\sigma$  for the experiments. We found that fixing  $\sigma$  to 0.1, and varying  $T$  from 0 to 1 was not creating a significant change in privacy. Hence, we resorted to a value of  $10^{-2}$  as used in (Shi et al., 2022).

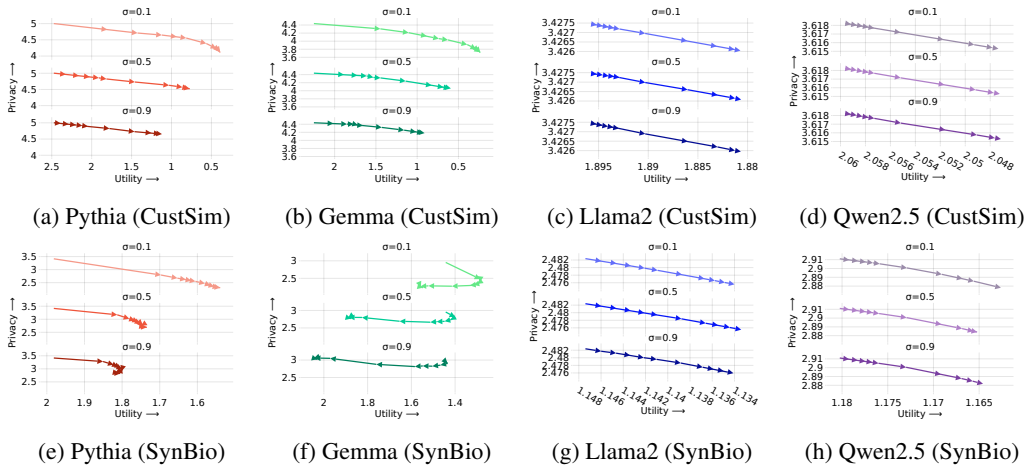


Figure 4: Privacy-utility trade-offs for differential privacy (DP-SGD) across models and datasets with varying noise levels ( $\sigma$ ). Larger  $\sigma$  increases privacy but reduces utility, with the decline more pronounced in larger models. On SynBio, DP-SGD exhibits a similar pattern. Gemma shows a unique utility drop after two epochs, highlighting the complexity of the privacy-utility trade-off.

**Privacy-Utility trade-off:** Figures 4a - 4d illustrate the privacy-utility trade-off when varying the noise  $\sigma$  on the CustomerSim dataset across the models. Similar to previous plots, each curve begins with the performance of the pre-trained model. It is evident that DP-SGD maintains privacy effectively with minimal degradation, but with a trade-off between privacy and utility. For all models, lower noise value of  $\sigma = 0.1$  achieves better utility while higher values of  $\sigma = 0.5, 0.9$  provide better privacy. The total amount of utility achievable with DP-SGD is limited for the larger models like Llama2 and Qwen2.5, possibly due to clipping of larger norms that leads to suboptimal performance (Shen et al., 2021; Tramer & Boneh, 2020). Similar trends can be observed for SynBio in Figures 4e, 4g, 4h, but for Figure 4f (with Gemma) utility declines after approximately two epochs.

**Efficiency:** Differential privacy comes with a high computational cost, since the gradients of each datapoint need norm, clipping computations, and noise additions. Empirically, we observe a relative FLOPs requirement of  $C_{DP}/C_{FFT} = 1.33$ . Additionally, the per-sample operations for clipping mean that we need to keep copies of the gradient for each datapoint in memory, which requires substantially more GPU memory decreasing the feasible batch size and the overall training throughput.

**Takeaway:** DP-SGD offers the *the best privacy-utility trade-off* for small-scale models (Pythia, Gemma), but the privacy gains come at the *cost of efficiency*. DP-SGD finds it difficult to achieve good utility in larger models (Llama2, Qwen2.5) as also observed in (Shen et al., 2021; Tramer & Boneh, 2020).

### 4.3 LORA FINE-TUNING FOR EFFICIENCY AND UTILITY IMPROVES PRIVACY

LoRA (Hu et al., 2022) is a parameter-efficient fine-tuning method developed to reduce the compute and memory requirements of fine-tuning LLMs.

#### Update rules:

$$W_{t+1} = W_0 + \frac{\alpha}{r} \Delta W_{t+1}; \Delta W_{t+1} = \Delta W_t - \eta \nabla_{\Delta W_t} \mathcal{L}(\mathcal{M}_{W_t}(X), X) \quad (3)$$

LoRA freezes the weights of the pretrained base model  $W_0$  and only fine-tunes an adapter matrix  $\Delta W_t$ . During training,  $\Delta W_t$  is stored separately from  $W_0$  as two low-rank matrices with rank  $r$ :  $\Delta W_t = B_t A_t$  with  $B_t \in \mathbb{R}^{d \times r}$ ,  $A_t \in \mathbb{R}^{r \times k}$ ,  $W_0 \in \mathbb{R}^{d \times k}$ .  $r$  is typically very small, often between 4 and 32.  $\alpha$  is a scaling factor for the LoRA adapters affecting the behavior of the base weights  $W_0$ .

**Privacy-Utility trade-off:** We investigate the effects of varying both the rank  $r$  and scaling parameter  $\alpha$  of LoRA. We use common rank values of 16 and 32 with  $\alpha \in \{16, 32, 64, 128\}$ . While LoRA has been explored for privacy in conjunction with DP-SGD (Yu et al., 2022), *there has been no prior work that examines the privacy benefits of LoRA alone*.

Figures 5a- 5d show the trade-off with rank 16 and varying parameters of  $\alpha$  for the CustomerSim dataset across all the models. As seen in Figures 5a and 5b, smaller-scale models (Pythia and Gemma) exhibit a better privacy-utility trade-off when the rank and  $\alpha$  values are equal, compared to the other configurations. For the larger models, the privacy declines at later epochs, while utility is mostly retained. We also analyze the trade-off for the SynBio dataset in Figures 5e- 5h for rank 16 which have similar observations. Similar trends are observed for rank 32 shown in Appendix I. However, relative to rank 16, rank 32 exhibits reduced privacy while maintaining similar utility.

**Efficiency:** LoRA has negligible additional FLOPs required for the low-rank matrices  $\Delta W$ , though they are much smaller than the full base weight  $W_0$  ( $r \ll \min(d, k)$ ). However, during the backward pass, LoRA requires much less compute. We observe a relative FLOPs requirement of  $C_{LoRA}/C_{FFT} \approx 0.65$ . The original paper (Hu et al., 2022) reports a 25% speedup during training. LoRA has needs of less GPU memory than FFT, since no optimizer states and gradients need to be stored for the base weights  $W_0$ , which makes it possible to run it with larger batch-sizes and thus an overall increased training throughput.

**Takeaway:** Ours is the first work to investigate the privacy benefits of LoRA, in isolation. The optimal privacy-utility trade-off is mostly achieved with  $r = \alpha$ . Finally, LoRA is much more *computationally efficient than FFT and especially DP-SGD*. Overall it provides the *best privacy-utility trade-offs* and shows that *privacy can be achieved without additional computational costs*.

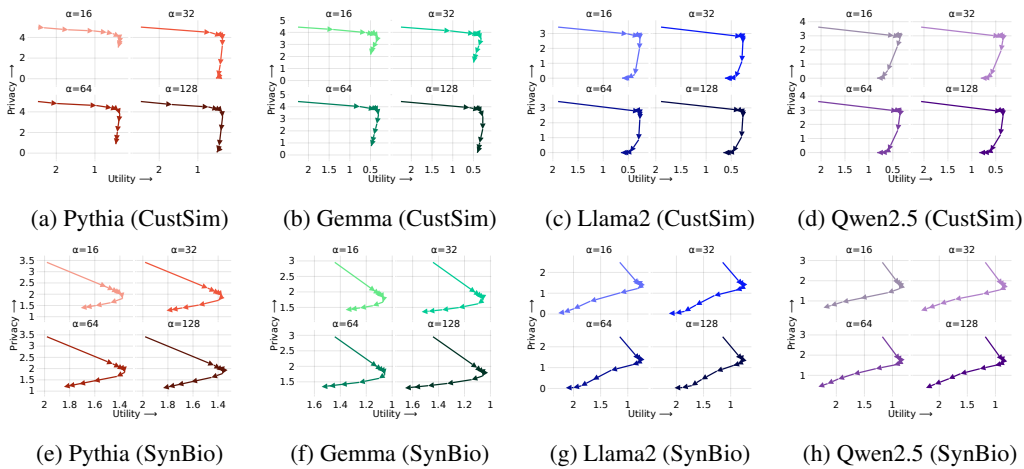


Figure 5: Privacy–utility trade-offs for LoRA fine-tuning with rank 16 and varying scaling factor  $\alpha$ . Increasing  $\alpha$  generally improves or maintains utility but reduces privacy. Smaller models achieve more favorable privacy–utility trade-offs, while larger models retain utility at the cost of reduced privacy during extended training.

## 5 COMPARISON OF FINE-TUNING METHODS

In this section, we compare all fine-tuning methods along three key dimensions—privacy, utility, and efficiency—using the best hyperparameter configurations identified in Section 4. We present Pareto-optimal curves under a fair setting, where each method is evaluated at its best configuration. By examining these curves, one can select the checkpoint most appropriate for the task at hand.

Recall that *utility is measured as test loss on non-sensitive tokens*. Besides measuring privacy as the *training loss on sensitive tokens*, we define two more metrics for privacy: (a) Privacy Loss over sensitive tokens (Abadi et al., 2016) and (b) Canary exposure (Carlini et al., 2019).

**Privacy Loss:** Let  $\mathcal{D}$  be the dataset with the sensitive token  $d$  under consideration, and  $\mathcal{D}'$  be the dataset without it. Considering  $\mathcal{M}(\mathcal{D})$  as the fine-tuned model that has seen the datapoint  $d$ ,  $\mathcal{M}(\mathcal{D}')$  as the model that has not seen  $d$ , and leveraging the definition from (Abadi et al., 2016), we can define the privacy loss (PL) as:

$$PL = \log \frac{P(\mathcal{M}(\mathcal{D}) = d)}{P(\mathcal{M}(\mathcal{D}') = d)} = \log \frac{P_{\mathcal{D}}(d)}{P_{\mathcal{D}'}(d)} \quad (4)$$

A natural approximation of such a model with the data unseen is the pretrained model. Recall that negative log-likelihood for a token  $d$  is  $NLL_{(\mathcal{D})}(d) = -\log P_{\mathcal{D}}(d)$ . We may rewrite the privacy loss as the difference between the negative log-likelihood of the pretrained model and the fine-tuned model over the token  $d$ :

$$PL = \log \frac{P_{\mathcal{D}}(d)}{P_{\mathcal{D}'}(d)} = \log P_{\mathcal{D}}(d) - \log P_{\mathcal{D}'}(d) = NLL_{\mathcal{D}'}(d) - NLL_{\mathcal{D}}(d) \quad (5)$$

Unlike our measure where we observe the absolute loss on sensitive tokens from the training data, this metric observes the relative change in loss w.r.t a model that has never seen the datapoint.

**Canary exposure:** Originally proposed by (Carlini et al., 2019) to measure unintended memorization, canaries are random sequences (e.g.,  $s = \text{“My ID is } \circ \circ \circ \circ \circ \text{”}$ ) inserted into training data. Exposure is computed by enumerating all possible sequences from the randomness space  $\mathcal{R}$  and calculating the negative log-rank. Following the definition of exposure in (Carlini et al., 2019), the exposure of a canary  $s[r]$  in a model  $\mathcal{M}$  over randomness space  $\mathcal{R}$  is defined as:

$$exposure_{\mathcal{M}} = \log_2 |\mathcal{R}| - \log_2 \text{rank}_{\mathcal{M}}(s[r]) \quad (6)$$

where the rank of a canary is its index in the list of all possibly-instantiated canaries, ordered by the model perplexity of all those sequences. In this setting, we inserted the canary “My ID is 34175” into the training data for 10 times and measured its exposure as a metric for privacy.

378  
379  
380  
381  
382  
383  
384  
385  
386  
387  
388  
389  
390  
391  
392  
393  
394  
395  
396  
397  
398  
399  
400  
401  
402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427  
428  
429  
430  
431

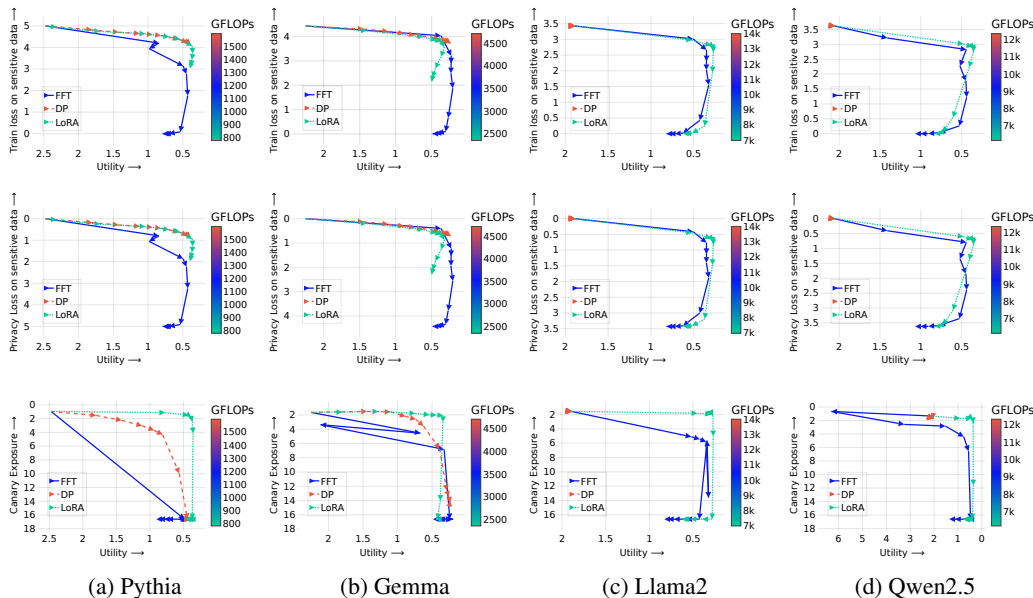


Figure 6: Pareto-optimal tradeoffs (*top-row* : Training loss on sensitive tokens ; *middle row* : Privacy Loss on sensitive tokens ; *bottom-row* : Canary exposure). DP-SGD achieves best privacy-utility tradeoffs on small models with reduced utility on large models, albeit with a significantly high computational cost; FFT offers reasonable efficiency-utility tradeoffs with the worst privacy; and LoRA provides a balanced trade-off among all objectives.

For efficiency, we measure floating point operations (FLOPs) based on the number of operations incurred (e.g., matrix multiplication, addition, etc.) during training.

We obtain the best configuration for DP-SGD with noise ratio of  $\sigma = 0.1$ , and for LoRA with  $r = \alpha = 16, \sigma = 16$ . Figure 6 presents the pareto-optimal curves for CustomerSim, comparing the fine-tuning strategies over all the metrics.

**Privacy:** Regarding *privacy*, FFT shows *poor privacy* over extended training over all 3 metrics, while DP achieves the *highest privacy levels*. LoRA provides *similar privacy as DP* throughout most epochs but declines gradually with extended training, especially on the large-scale model.

**Utility:** FFT maintains *relatively strong utility* on CustomerSim. DP-SGD while *yielding good utility in smaller models* (Pythia,Gemma), *performs poorly in the larger models* (Llama2,Qwen2.5). In contrast, LoRA *consistently preserves higher utility* across the entire training period.

**Efficiency:** The color bar in Figure 6 highlights the FLOPs intensity associated with each fine-tuning strategy. DP-SGD requires the *highest number of FLOPs* (in red) due to the need for per-sample gradient computation, where each sample corresponds to a token within each training sequence. FFT demands a *moderate number of FLOPs* (in blue), proportional to the total number of parameters. Finally, LoRA requires the *fewest FLOPs* (in green), as especially during backpropagation, most operations only operate on the low-rank matrices.

Assessing all the above aspects, this challenges the belief that privacy comes at the cost of efficiency, showing that **LoRA, an efficient fine-tuning method can provide privacy benefits**. Alongside, LoRA-tuned models retain performance on benchmark datasets similar to pre-trained models throughout training as shown in Figure 28. For a full comparison of methods, refer to Appendix J.

## 6 CONJECTURE: LoRA’S PRIVACY BENEFITS RELATIVE TO DP-SGD

In this section, we present a formal conjecture on LoRA’s privacy benefits, drawing an analogy to differential privacy (DP-SGD).

**Analogy between LoRA and DP-SGD:** We hypothesize that LoRA’s low-rank constraint helps in preserving privacy of the sensitive data. By forcing all the parameter updates into a low-rank subspace, LoRA can reduce the influence a training point can have on the model. Essentially, this is the same quantity that DP-SGD also aims to reduce by clipping and adding noise to the gradients. DP-SGD follows the route of adding randomness, and LoRA follows the route of compressing updates. The effect becomes more on the sensitive data as they are rarely occurring in the training set, which may possibly lead to it being absent in the subspace learnt by LoRA. Empirically, Figure 7 shows that low-rank LoRA model results in better privacy than its higher rank counterparts or full fine-tuning. We formally show the analogy in the theorem below:

**Theorem 1.** *Let  $D$  and  $D'$  be two neighbouring datasets differing by a datapoint  $i$ . Both LoRA and DP-SGD reduce the influence of  $i$  in comparison to a fully trained model, i.e.  $\Delta_{DP}^i \leq \Delta_{FFT}^i$  and  $\Delta_{LoRA}^i \leq \Delta_{FFT}^i$ . This indicates privacy preservation of datapoint  $i$  as less the influence, less will be its distinguishability.*

*Proof.* Consider we have  $D$  and  $D'$  as two neighbouring datasets differing by a datapoint  $i$ . Let  $W(D)$  and  $W(D')$  be the respective model parameters where  $W(\cdot) \in \mathbb{R}^{d \times k}$ . Similarly, let  $\Delta W(D)$  and  $\Delta W(D')$  be the change in respective aggregate model gradients. Since  $D$  and  $D'$  differ by a point  $i$ , we can attribute the difference between  $\Delta W(D)$  and  $\Delta W(D')$  to the contribution or influence of datapoint  $i$  on the aggregate update.

$$\therefore \Delta_{FFT}^i = \|\Delta W(D) - \Delta W(D')\|_F = \|g_i\|_F$$

For DP-SGD, it is well established that noise (of scale  $\sigma$ ) gets added to the clipped gradients (gradients clipped by a threshold of  $T$ ) that masks the impact of the datapoint (Abadi et al., 2016). We can thereby frame the influence of the datapoint  $i$  w.r.t DP-SGD model as:

$$\Delta_{DP}^i = \frac{\|g_i\|_F}{\max(1, \frac{\|g_i\|_F}{T})} + \mathcal{N}(0, \sigma^2 \mathcal{I}) \approx \frac{\|g_i\|_F/T}{\sigma} = \frac{\|g_i\|_F}{T\sigma} \leq \Delta_{FFT}^i$$

$$\therefore \Delta_{DP}^i \leq \Delta_{FFT}^i$$

When LoRA with low rank  $r$  is applied, the weight updates get projected to a subspace  $\mathcal{S}$  of rank  $r$ . Let the projection be  $\mathcal{P}_r : \mathbb{R}^{d \times k} \rightarrow \mathcal{S}_r$ . We can frame the influence of the datapoint  $i$  w.r.t LoRA model as:

$$\Delta_{LoRA}^i = \|\mathcal{P}_r(\Delta W(D)) - \mathcal{P}_r(\Delta W(D'))\|_F = \|\mathcal{P}_r(g_i)\|_F ; \mathcal{P}_r \text{ being a linear projection operator.}$$

Now,  $\|\mathcal{P}_r(g_i)\|_F \leq \|g_i\|_F$  as projection  $\mathcal{P}_r$  will drop any point from  $g_i$  lying outside its subspace.

$$\therefore \Delta_{LoRA}^i \leq \Delta_{FFT}^i \quad \square$$

We can see that both LoRA and DP-SGD reduce the influence of the datapoint ; while DP-SGD does it probabilistically using noise  $\sigma$ , LoRA does it deterministically through low-rank projection  $\mathcal{P}_r$ . This aligns with prior work (Chaudhuri et al., 2013; Kapralov & Talwar, 2013) linking principal component analysis (PCA) or low-rank factorization to DP, where PCA limits leakage by discarding components (similar as in LoRA) and DP-SGD by adding noise.

## 7 CONCLUDING DISCUSSION

We investigated the interplay of privacy, utility, and efficiency in fine-tuning LLMs introducing measures that distinguish between sensitive and non-sensitive tokens. The conventional wisdom is that achieving privacy, using methods like DP-SGD, comes at the cost of computational efficiency. In contrast, we demonstrate that efficient fine-tuning methods like LoRA, initially designed for efficiency, achieves privacy benefits without any computational overhead. Simultaneously, LoRA retains the utility of language understanding compared to DP-SGD, or even FFT, highlighting its superiority in balancing all three objectives. Our paper calls for a joint venture of privacy and systems communities in achieving privacy-aware efficient fine-tuning of LLMs while retaining utility.

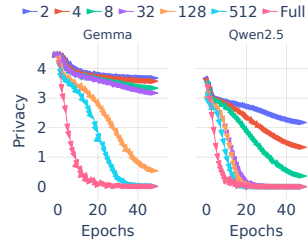


Figure 7: Privacy improves with low-rank matrix.

486 ETHICS STATEMENT  
487

488 This study utilized publicly available datasets (Shi et al., 2022; Holmes et al., 2024), some of which  
489 included identifiable information such as personal details. However, third-party organizations pre-  
490 processed and validated the datasets to ensure that no real individuals’ data were present, thus miti-  
491 gating potential privacy concerns.

492 This project received ethical clearance from the Ethical Review Board of the affiliated institution on  
493 October 21, 2024 (Approval No. 24-09-4), with no ethical concerns raised.  
494

495  
496 REPRODUCIBILITY STATEMENT  
497

498 To ensure the reproducibility of the results in the paper, we will release the code, execution envi-  
499 ronment, and artifacts publicly upon acceptance. The datasets – CustomerSim (Shi et al., 2022)  
500 and SynBio (Holmes et al., 2024) are publicly available. The scripts for fine-tuning are included in  
501 supplementary material. All hyperparameters for reproducing the results in the paper are listed in  
502 Appendix E.  
503

504 REFERENCES  
505

506 Martin Abadi, Andy Chu, Ian Goodfellow, H. Brendan McMahan, Ilya Mironov, Kunal Talwar,  
507 and Li Zhang. Deep learning with differential privacy. In *Proceedings of the 2016 ACM*  
508 *SIGSAC Conference on Computer and Communications Security, CCS ’16*, pp. 308–318, New  
509 York, NY, USA, 2016. Association for Computing Machinery. ISBN 9781450341394. doi:  
510 10.1145/2976749.2978318. URL <https://doi.org/10.1145/2976749.2978318>.

511 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale-  
512 man, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical  
513 report. *arXiv preprint arXiv:2303.08774*, 2023.  
514

515 Stella Biderman, Hailey Schoelkopf, Quentin Gregory Anthony, Herbie Bradley, Kyle O’Brien, Eric  
516 Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, et al.  
517 Pythia: A suite for analyzing large language models across training and scaling. In *International*  
518 *Conference on Machine Learning*, pp. 2397–2430. PMLR, 2023.

519 Hannah Brown, Katherine Lee, Fatemehsadat Mireshghallah, Reza Shokri, and Florian Tramèr.  
520 What does it mean for a language model to preserve privacy? In *Proceedings of the 2022 ACM*  
521 *Conference on Fairness, Accountability, and Transparency, FAccT ’22*, pp. 2280–2292, New  
522 York, NY, USA, 2022. Association for Computing Machinery. ISBN 9781450393522. doi: 10.  
523 1145/3531146.3534642. URL <https://doi.org/10.1145/3531146.3534642>.

524 Tom B Brown. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020.  
525

526 Nicholas Carlini, Chang Liu, Úlfar Erlingsson, Jernej Kos, and Dawn Song. The secret sharer:  
527 Evaluating and testing unintended memorization in neural networks. In *28th USENIX security*  
528 *symposium (USENIX security 19)*, pp. 267–284, 2019.  
529

530 Nicholas Carlini, Florian Tramer, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine  
531 Lee, Adam Roberts, Tom Brown, Dawn Song, Ulfar Erlingsson, et al. Extracting training data  
532 from large language models. In *30th USENIX security symposium (USENIX Security 21)*, pp.  
533 2633–2650, 2021.

534 Kamalika Chaudhuri, Anand D Sarwate, and Kaushik Sinha. A near-optimal algorithm for  
535 differentially-private principal components. *The Journal of Machine Learning Research*, 14(1):  
536 2905–2943, 2013.  
537

538 Xin Luna Dong, Seungwhan Moon, Yifan Ethan Xu, Kshitiz Malik, and Zhou Yu. Towards next-  
539 generation intelligent assistants leveraging llm techniques. In *Proceedings of the 29th ACM*  
*SIGKDD Conference on Knowledge Discovery and Data Mining*, pp. 5792–5793, 2023.

- 540 Christophe Dupuy, Radhika Arava, Rahul Gupta, and Anna Rumshisky. An efficient dp-sgd mecha-  
541 nism for large scale nlu models. In *ICASSP 2022-2022 IEEE International Conference on Acous-  
542 tics, Speech and Signal Processing (ICASSP)*, pp. 4118–4122. IEEE, 2022.
- 543  
544 EU. General data protection regulation, 2016. URL <https://gdpr-info.eu/>.
- 545  
546 Wenjie Fu, Huandong Wang, Chen Gao, Guanghua Liu, Yong Li, and Tao Jiang. Practical member-  
547 ship inference attacks against fine-tuned large language models via self-prompt calibration. *arXiv  
548 preprint arXiv:2311.06062*, 2023.
- 549  
550 Zeyu Han, Chao Gao, Jinyang Liu, Jeff Zhang, and Sai Qian Zhang. Parameter-efficient fine-tuning  
551 for large models: A comprehensive survey, 2024. URL [https://arxiv.org/abs/2403.  
552 14608](https://arxiv.org/abs/2403.14608).
- 553  
554 Jamie Hayes, Borja Balle, and Saeed Mahloujifar. Bounding training data reconstruction in DP-  
555 SGD. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL  
<https://openreview.net/forum?id=7Lz4tZrYlx>.
- 556  
557 Junxian He, Chunting Zhou, Xuezhe Ma, Taylor Berg-Kirkpatrick, and Graham Neubig. Towards a  
558 unified view of parameter-efficient transfer learning. *arXiv preprint arXiv:2110.04366*, 2021.
- 559  
560 Langdon Holmes, Scott Crossley, Perpetual Baffour, Jules King, Lauryn Burleigh, Mag-  
561 gie Demkin, Ryan Holbrook, Walter Reade, and Addison Howard. The learning  
562 agency lab - pii data detection, 2024. URL [https://kaggle.com/competitions/  
563 pii-detection-removal-from-educational-data](https://kaggle.com/competitions/pii-detection-removal-from-educational-data).
- 564  
565 Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, An-  
566 drea Gesmundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning for nlp.  
567 In *International conference on machine learning*, pp. 2790–2799. PMLR, 2019.
- 568  
569 Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,  
570 and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *International Con-  
571 ference on Learning Representations*, 2022. URL [https://openreview.net/forum?  
572 id=nZeVKeeFYf9](https://openreview.net/forum?id=nZeVKeeFYf9).
- 573  
574 Masahiro Kaneko, Youmi Ma, Yuki Wata, and Naoaki Okazaki. Sampling-based pseudo-likelihood  
575 for membership inference attacks. *arXiv preprint arXiv:2404.11262*, 2024.
- 576  
577 Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child,  
578 Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language  
579 models. *arXiv preprint arXiv:2001.08361*, 2020.
- 580  
581 Michael Kapralov and Kunal Talwar. On differentially private low rank approximation. In *Proceed-  
582 ings of the twenty-fourth annual ACM-SIAM symposium on Discrete algorithms*, pp. 1395–1414.  
583 SIAM, 2013.
- 584  
585 Hareem Kibriya, Wazir Zada Khan, Ayesha Siddiq, and Muhammad Khurram Khan. Privacy  
586 issues in large language models: A survey. *Computers and Electrical Engineering*, 120:109698,  
587 2024. ISSN 0045-7906. doi: <https://doi.org/10.1016/j.compeleceng.2024.109698>. URL [https:  
588 //www.sciencedirect.com/science/article/pii/S0045790624006256](https://www.sciencedirect.com/science/article/pii/S0045790624006256).
- 589  
590 Siwon Kim, Sangdoon Yun, Hwaran Lee, Martin Gubri, Sungroh Yoon, and Seong Joon Oh. Propile:  
591 probing privacy leakage in large language models. In *Proceedings of the 37th International Con-  
592 ference on Neural Information Processing Systems, NIPS '23*, Red Hook, NY, USA, 2024. Curran  
593 Associates Inc.
- 594  
595 Diederik P Kingma. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*,  
596 2014.
- 597  
598 Tao Lei, Junwen Bai, Siddhartha Brahma, Joshua Ainslie, Kenton Lee, Yanqi Zhou, Nan Du, Vincent  
599 Zhao, Yuexin Wu, Bo Li, et al. Conditional adapters: Parameter-efficient transfer learning with  
fast inference. *Advances in Neural Information Processing Systems*, 36:8152–8172, 2023.

- 594 Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt  
595 tuning. *arXiv preprint arXiv:2104.08691*, 2021.  
596
- 597 Jonathan Li, Will Aitken, Rohan Bhambhoria, and Xiaodan Zhu. Prefix propagation: Parameter-  
598 efficient tuning for long sequences. *arXiv preprint arXiv:2305.12086*, 2023.  
599
- 600 Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation.  
601 In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli (eds.), *Proceedings of the 59th*  
602 *Annual Meeting of the Association for Computational Linguistics and the 11th International Joint*  
603 *Conference on Natural Language Processing (Volume 1: Long Papers)*, pp. 4582–4597, Online,  
604 August 2021a. Association for Computational Linguistics. doi: 10.18653/v1/2021.acl-long.353.  
605 URL <https://aclanthology.org/2021.acl-long.353/>.
- 606 Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. *arXiv*  
607 *preprint arXiv:2101.00190*, 2021b.  
608
- 609 Shih-Yang Liu, Chien-Yi Wang, Hongxu Yin, Pavlo Molchanov, Yu-Chiang Frank Wang, Kwang-  
610 Ting Cheng, and Min-Hung Chen. Dora: Weight-decomposed low-rank adaptation. *arXiv*  
611 *preprint arXiv:2402.09353*, 2024a.  
612
- 613 Xiao Liu, Kaixuan Ji, Yicheng Fu, Weng Lam Tam, Zhengxiao Du, Zhilin Yang, and Jie Tang. P-  
614 tuning v2: Prompt tuning can be comparable to fine-tuning universally across scales and tasks.  
615 *arXiv preprint arXiv:2110.07602*, 2021.
- 616 Xiao Liu, Kaixuan Ji, Yicheng Fu, Weng Tam, Zhengxiao Du, Zhilin Yang, and Jie Tang. P-tuning:  
617 Prompt tuning can be comparable to fine-tuning across scales and tasks. In Smaranda Muresan,  
618 Preslav Nakov, and Aline Villavicencio (eds.), *Proceedings of the 60th Annual Meeting of the*  
619 *Association for Computational Linguistics (Volume 2: Short Papers)*, pp. 61–68, Dublin, Ireland,  
620 May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.acl-short.8. URL  
621 <https://aclanthology.org/2022.acl-short.8/>.  
622
- 623 Xiao Liu, Yanan Zheng, Zhengxiao Du, Ming Ding, Yujie Qian, Zhilin Yang, and Jie Tang. Gpt  
624 understands, too. *AI Open*, 5:208–215, 2024b.
- 625 Xiao-Yang Liu, Rongyi Zhu, Daochen Zha, Jiechao Gao, Shan Zhong, Matt White, and Meikang  
626 Qiu. Differentially private low-rank adaptation of large language model using federated learning.  
627 *ACM Transactions on Management Information Systems*, 2023.  
628
- 629 Nils Lukas, Ahmed Salem, Robert Sim, Shruti Tople, Lukas Wutschitz, and Santiago Zanella-  
630 Beguelin. Analyzing Leakage of Personally Identifiable Information in Language Models . In  
631 *2023 IEEE Symposium on Security and Privacy (SP)*, pp. 346–363, Los Alamitos, CA, USA,  
632 May 2023. IEEE Computer Society. doi: 10.1109/SP46215.2023.10179300. URL <https://doi.ieeecomputersociety.org/10.1109/SP46215.2023.10179300>.  
633  
634
- 635 Olivia Ma, Jonathan Passerat-Palmbach, and Dmitrii Usynin. Efficient and private: Memorisation  
636 under differentially private parameter-efficient fine-tuning in language models. *arXiv preprint*  
637 *arXiv:2411.15831*, 2024.
- 638 Sourab Mangrulkar, Sylvain Gugger, Lysandre Debut, Younes Belkada, Sayak Paul, and Benjamin  
639 Bossan. Peft: State-of-the-art parameter-efficient fine-tuning methods. <https://github.com/huggingface/peft>, 2022.  
640  
641
- 642 Justus Mattern, Fatemehsadat Mireshghallah, Zhijing Jin, Bernhard Schölkopf, Mrinmaya Sachan,  
643 and Taylor Berg-Kirkpatrick. Membership inference attacks against language models via neigh-  
644 bourhood comparison. *arXiv preprint arXiv:2305.18462*, 2023.  
645
- 646 Omri Mendels, Coby Peled, Nava Vaisman Levy, Sharon Hart, Tomer Rosenthal, Limor Lahiani,  
647 et al. Microsoft Presidio: Context aware, pluggable and customizable pii anonymization service  
for text and images, 2018. URL <https://microsoft.github.io/presidio>.

- 648 Fatemehsadat Mireshghallah, Kartik Goyal, Archit Uniyal, Taylor Berg-Kirkpatrick, and Reza  
649 Shokri. Quantifying privacy risks of masked language models using membership inference  
650 attacks. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (eds.), *Proceedings of the*  
651 *2022 Conference on Empirical Methods in Natural Language Processing*, pp. 8332–8347,  
652 Abu Dhabi, United Arab Emirates, December 2022a. Association for Computational Linguistics. doi: 10.18653/v1/2022.emnlp-main.570. URL [https://aclanthology.org/2022.](https://aclanthology.org/2022.emnlp-main.570)  
653 [emnlp-main.570](https://aclanthology.org/2022.emnlp-main.570).
- 655 Fatemehsadat Mireshghallah, Archit Uniyal, Tianhao Wang, David Evans, and Taylor Berg-  
656 Kirkpatrick. An empirical analysis of memorization in fine-tuned autoregressive language models.  
657 In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (eds.), *Proceedings of the 2022 Conference*  
658 *on Empirical Methods in Natural Language Processing*, pp. 1816–1826, Abu Dhabi, United Arab  
659 Emirates, December 2022b. Association for Computational Linguistics. doi: 10.18653/v1/2022.  
660 [emnlp-main.119](https://aclanthology.org/2022.emnlp-main.119). URL <https://aclanthology.org/2022.emnlp-main.119>.
- 661 Humza Naveed, Asad Ullah Khan, Shi Qiu, Muhammad Saqib, Saeed Anwar, Muhammad Usman,  
662 Naveed Akhtar, Nick Barnes, and Ajmal Mian. A comprehensive overview of large language  
663 models. *arXiv preprint arXiv:2307.06435*, 2023.
- 664 Ashwinee Panda, Christopher A Choquette-Choo, Zhengming Zhang, Yaoqing Yang, and Prateek  
665 Mittal. Teach llms to phish: Stealing private information from language models. *arXiv preprint*  
666 *arXiv:2403.00871*, 2024.
- 667 Yinchen Shen, Zhiguo Wang, Ruoyu Sun, and Xiaojing Shen. Towards understanding the impact of  
668 model size on differential private classification. *arXiv preprint arXiv:2111.13895*, 2021.
- 670 Weiyang Shi, Aiqi Cui, Evan Li, Ruoxi Jia, and Zhou Yu. Selective differential privacy for language  
671 modeling. In Marine Carpuat, Marie-Catherine de Marneffe, and Ivan Vladimir Meza Ruiz (eds.),  
672 *Proceedings of the 2022 Conference of the North American Chapter of the Association for Com-*  
673 *putational Linguistics: Human Language Technologies*, pp. 2848–2859, Seattle, United States,  
674 July 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.naacl-main.205.  
675 URL <https://aclanthology.org/2022.naacl-main.205>.
- 676 Yan Shvartzshnaider and Vasisht Duddu. Position: Contextual integrity is inadequately applied to  
677 language models. In *Forty-second International Conference on Machine Learning Position Paper*  
678 *Track*, 2025. URL <https://openreview.net/forum?id=YmTxir1HUX>.
- 679 Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya  
680 Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. Gemma: Open  
681 models based on gemini research and technology. *arXiv preprint arXiv:2403.08295*, 2024.
- 682 Arun James Thirunavukarasu, Darren Shu Jeng Ting, Kabilan Elangovan, Laura Gutierrez,  
683 Ting Fang Tan, and Daniel Shu Wei Ting. Large language models in medicine. *Nature medicine*,  
684 29(8):1930–1940, 2023.
- 686 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko-  
687 lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open founda-  
688 tion and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- 689 Florian Tramer and Dan Boneh. Differentially private learning needs better features (or much more  
690 data). *arXiv preprint arXiv:2011.11660*, 2020.
- 692 Florian Tramèr, Gautam Kamath, and Nicholas Carlini. Position: considerations for differentially  
693 private learning with large-scale public pretraining. In *Proceedings of the 41st International*  
694 *Conference on Machine Learning, ICML’24*. JMLR.org, 2024.
- 695 Shen Wang, Tianlong Xu, Hang Li, Chaoli Zhang, Joleen Liang, Jiliang Tang, Philip S Yu, and  
696 Qingsong Wen. Large language models for education: A survey and outlook. *arXiv preprint*  
697 *arXiv:2403.18105*, 2024.
- 699 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi,  
700 Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. Transformers: State-of-the-art  
701 natural language processing. In *Proceedings of the 2020 conference on empirical methods in*  
*natural language processing: system demonstrations*, pp. 38–45, 2020.

- 702 An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li,  
703 Dayiheng Liu, Fei Huang, Haoran Wei, et al. Qwen2. 5 technical report. *arXiv preprint*  
704 *arXiv:2412.15115*, 2024.
- 705
- 706 Yifan Yao, Jinhao Duan, Kaidi Xu, Yuanfang Cai, Zhibo Sun, and Yue Zhang. A survey on  
707 large language model (llm) security and privacy: The good, the bad, and the ugly. *High-*  
708 *Confidence Computing*, 4(2):100211, 2024. ISSN 2667-2952. doi: [https://doi.org/10.1016/j.hcc.](https://doi.org/10.1016/j.hcc.2024.100211)  
709 2024.100211. URL [https://www.sciencedirect.com/science/article/pii/](https://www.sciencedirect.com/science/article/pii/S266729522400014X)  
710 [S266729522400014X](https://www.sciencedirect.com/science/article/pii/S266729522400014X).
- 711 Ashkan Yousefpour, Igor Shilov, Alexandre Sablayrolles, Davide Testuggine, Karthik Prasad, Mani  
712 Malek, John Nguyen, Sayan Ghosh, Akash Bharadwaj, Jessica Zhao, Graham Cormode, and  
713 Ilya Mironov. Opacus: User-friendly differential privacy library in PyTorch. *arXiv preprint*  
714 *arXiv:2109.12298*, 2021.
- 715
- 716 Da Yu, Saurabh Naik, Arturs Backurs, Sivakanth Gopi, Huseyin A Inan, Gautam Kamath, Janardhan  
717 Kulkarni, Yin Tat Lee, Andre Manoel, Lukas Wutschitz, et al. Differentially private fine-tuning  
718 of language models. *arXiv preprint arXiv:2110.06500*, 2021.
- 719
- 720 Da Yu, Saurabh Naik, Arturs Backurs, Sivakanth Gopi, Huseyin A. Inan, Gautam Kamath, Janardhan  
721 Kulkarni, Yin Tat Lee, Andre Manoel, Lukas Wutschitz, Sergey Yekhanin, and Huishuai  
722 Zhang. Differentially private fine-tuning of language models, 2022. URL [https://arxiv.](https://arxiv.org/abs/2110.06500)  
723 [org/abs/2110.06500](https://arxiv.org/abs/2110.06500).
- 724
- 725 Jie Zhang, Debeshee Das, Gautam Kamath, and Florian Tramèr. Membership inference attacks  
726 cannot prove that a model was trained on your data, 2024. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2409.19798)  
727 [2409.19798](https://arxiv.org/abs/2409.19798).
- 728
- 729 Shengyu Zhang, Linfeng Dong, Xiaoya Li, Sen Zhang, Xiaofei Sun, Shuhe Wang, Jiwei Li, Runyi  
730 Hu, Tianwei Zhang, Fei Wu, et al. Instruction tuning for large language models: A survey. *arXiv*  
731 *preprint arXiv:2308.10792*, 2023.
- 732
- 733 Xuandong Zhao, Lei Li, and Yu-Xiang Wang. Provably confidential language modelling. In  
734 Marine Carpuat, Marie-Catherine de Marneffe, and Ivan Vladimir Meza Ruiz (eds.), *Proceed-*  
735 *ings of the 2022 Conference of the North American Chapter of the Association for Computa-*  
736 *tional Linguistics: Human Language Technologies*, pp. 943–955, Seattle, United States, July  
737 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.naacl-main.69. URL  
738 <https://aclanthology.org/2022.naacl-main.69>.
- 739
- 740 Yaoming Zhu, Jiangtao Feng, Chengqi Zhao, Mingxuan Wang, and Lei Li. Counter-interference  
741 adapter for multilingual machine translation. *arXiv preprint arXiv:2104.08154*, 2021.

## 740 APPENDIX

### 741 Overview of Appendices

- 744 • Appendix A: Limitations
- 745 • Appendix B: Use of LLMs
- 746 • Appendix C: Extended Related Work
- 747 • Appendix D: Dataset Samples
- 748 • Appendix E: Hyperparameters for reproducibility
- 749 • Appendix F: Prompt for annotating data.
- 750 • Appendix G: Human survey questionnaire
- 751 • Appendix H: Distinction between sensitive and non-sensitive
- 752 • Appendix I: Privacy Utility Efficiency tradeoffs
- 753 • Appendix J: Overall comparison
- 754
- 755

## A LIMITATIONS

One of the limitations of our work lies in fine-tuning the models for unsupervised setup and not extending it to other supervised tasks like question-answering and so on. However we believe that our findings would hold in any task setup, given that the nature of these fine-tuning techniques would not change. We experimented with LoRA as one the most generic PEFT methods – however, testing out other PEFT methods would be an interesting extension of our work to explore privacy benefits extensively in the systems community.

We consider the sensitive data to be from PII's defined in GDPR Article 4(1) (EU, 2016). However, sensitive data can also be context-dependent. While one piece of information may be sensitive for one, it may not be for the other,. However, we also acknowledge that while contextual integrity is an important task, it is still hard to formalize or implement as seen in (Shvartzshnaider & Duddu, 2025). Lastly, while this work provides an empirical measure of privacy across different fine-tuning methods, one can definitely use such a measure for optimisation during training and establish a theoretical bound on the privacy benefits that would then also be empirically validated. We intend to explore these directions in the future.

## B USE OF LLMs FOR PAPER-WRITING

In this paper, we use LLMs for the following purposes:

1. **Polishing the writing:** We check for grammatical mistakes, and make minor para-phrasing to improve the flow and coherence of the paper.
2. **Related Work:** Besides traditional search, we use the OpenAI Deep Research to identify relevant literature.

## C EXTENDED RELATED WORK

**Privacy attacks.** Privacy concerns in LLMs have gained attention in recent years Kibriya et al. (2024); Yao et al. (2024), particularly the possibility of exposing data via membership inference attacks (MIAs) Mireshghallah et al. (2022a); Mattern et al. (2023); Fu et al. (2023); Kaneko et al. (2024) and training data extraction attacks (DEAs) Carlini et al. (2021); Mireshghallah et al. (2022b); Panda et al. (2024). While MIAs have the goal of identifying if a certain datapoint was in the training dataset using model confidence scores, DEAs try to extract specific parts of data from the training dataset using different prompting strategies. A recent study Zhang et al. (2024) underscores that membership inference attacks may not be reliable due to their practical limitations of requiring knowledge of the entire training data. In this work, we focus on whether different fine-tuning methods mitigate the privacy risks for DEAs as they are more practical and effective mode of attacks Zhang et al. (2024).

**Privacy mitigation.** To mitigate privacy leakage through MIAs, differential privacy (DP) measures have been proposed, which add theoretical privacy guarantees to the training process Abadi et al. (2016) with no empirical validation. The key idea of DP is to reduce the influence of datapoints during training and reduce the chances of leakage during inference with some formal guarantees. However, these DP guarantees have been shown to not strongly hold for data reconstruction attacks in Hayes et al. (2023). In this work, we measure the empirical privacy benefits of DP on DEAs and find it to be effective.

**Privacy measurement.** The use of DP in LLMs has been criticized for offering privacy guarantees that are insufficient unless sensitive data is removed Brown et al. (2022); Tramèr et al. (2024). This finding raises the question of whether current measures of privacy is adequate. Intuitively, privacy evaluation should focus on the sensitive data encountered during training. While some studies on LLM privacy attacks consider sensitive vs. non-sensitive data during training Kim et al. (2024); Lukas et al. (2023); Zhao et al. (2022); Shi et al. (2022), they often ignore this distinction in evaluation. In our work, we preserve this distinction and propose a revised quantification that holds for DEAs.

Table 1: Example of samples from datasets with sensitive tokens highlighted.

| CustomerSim   |   |
|---|---|
| SYS: Hello, I am the customer support bot. What can I do for you?<br>USR: Hello robot. I ordered a pot several days ago but I can't track it.<br>SYS: Could you verify your full name?<br>USR: Patrick Schug<br>SYS: Verify your order number please.<br>USR: It's 843-58572-7002.<br>SYS: You can track your package with your tracking number , which is AGZIM516KL.<br>Are you happy about my answer?<br>USR: All good. See you.<br>SYS: Have a nice day! Bye. | SYS: Hello, I am the customer support bot. What can I do for you?<br>USR: Hi. Where is my package?<br>SYS: Could you verify your full name?<br>USR: I am Catherine Pena.<br>SYS: Verify your phone number.<br>USR: You can reach me at 547.302.3744.<br>SYS: The tracking number is VVTPHDB6VK.<br>Anything else?<br>USR: All good. |
| SynBio  |   |
| My name is Alexander Tanaka, and I'm a saleswoman with a year of experience. I recently completed a project that involved developing and implementing a new sales strategy for my company. I started by analyzing our current sales data to identify areas where we could improve...  | My name is Phillip Martinez, and I would like to share some aspects of my life's journey with you. I have had the pleasure of living in various places throughout my life, but I currently reside at 4537 Tanglewood Trail ... you can reach me via email at phillip-martinez@outlook.com or by phone at +86 19144 1648.            |

**Privacy-Efficiency tradeoff.** Additionally, to ensure privacy guarantees, DP requires higher computational resources (time and memory) during training. For example, compared to vanilla SGD, DP-SGD may incur up to 20x training time, which is often a bottleneck in resource-intensive tasks Dupuy et al. (2022). Therefore, the conventional wisdom is that *privacy comes at the cost of efficiency*. In this work, we probe ways in which the above limitation can be prevented.

**Efficiency-Utility tradeoffs.** Full fine-tuning (FFT) of large language models is expensive as it involves updating all the parameters. Recent work has focused on reducing the training cost while maintaining utility, with the introduction of parameter efficient fine-tuning (PEFT) techniques Han et al. (2024). Well-known PEFT methods include adapter-based fine-tuning Houlsby et al. (2019); Lei et al. (2023); Zhu et al. (2021); He et al. (2021), prefix-tuning Li & Liang (2021a), soft prompt-based fine-tuning Li & Liang (2021b); Li et al. (2023); Liu et al. (2021; 2024b); Lester et al. (2021), P-tuning Liu et al. (2022), and parameterized fine-tuning Hu et al. (2022); Liu et al. (2024a). Among different PEFT methods, we focus on LoRA Hu et al. (2022) as it is more generic, while the rest are task-specific and mostly limited to supervised settings.

**Privacy-Efficiency-Utility tradeoffs.** In an ideal world, the best scenario would be to achieve all the three objectives. While DP is shown to achieve privacy and utility, compromising on efficiency, LoRA is shown to achieve utility and efficiency at an unknown cost of privacy, that we try to characterize in this work. We observe that LoRA achieves similar level of privacy as DP. In recent work on private efficient fine-tuning Liu et al. (2023); Yu et al. (2021); Ma et al. (2024), the authors combine LoRA with DP to reduce the additional computational overhead induced by DP. However, in our work, we show that one can directly achieve all the three objectives — privacy, utility, and efficiency — by only using LoRA with no additional overhead.

## D DATASET SAMPLE

Table 1 lists two samples from each of the two datasets – Customersim (Shi et al., 2022) and SynBio (Holmes et al., 2024).

## E HYPERPARAMETERS

All experiments were conducted on a SLURM cluster consisting of nodes equipped with NVIDIA A100 and H100 GPUs. During fine-tuning, we use `Trainer` module from Transformers Wolf et al. (2020), and specifically `peft` library Mangrulkar et al. (2022) for LoRA and the `opacus` library Yousefpour et al. (2021) for DP. Table 2 and 3 lists the hyperparameters used for fine-tuning experiments. We report our results on a single run with a fixed seed of 42.

| Hyperparameter                      | Full Fine-Tuning (FFT) | LoRA              | DP                   |
|-------------------------------------|------------------------|-------------------|----------------------|
| Learning Rate                       | 0.00025                | 0.00025           | 0.00005              |
| Scheduler                           | Linear                 | Linear            | Linear               |
| Warmup Steps                        | 10                     | 10                | 10                   |
| Rank ( $r$ )                        | -                      | {16, 32}          | -                    |
| Scaling factor ( $\alpha$ )         | -                      | {16, 32, 64, 128} | -                    |
| Noise ratio ( $\sigma$ )            | -                      | {0.1, 0.5, 0.9}   | -                    |
| Clipping gradient threshold ( $T$ ) | -                      | -                 | $(1 \times 10^{-2})$ |

Table 2: Hyperparameters used for fine-tuning methods

| Dataset     | Model   | Full Fine-Tuning (FFT) | LoRA | DP |
|-------------|---------|------------------------|------|----|
| CustomerSim | Pythia  | 16                     | 32   | 8  |
|             | Gemma   | 8                      | 16   | 4  |
|             | Llama2  | 4                      | 8    | 2  |
|             | Qwen2.5 | 4                      | 8    | 2  |
| SynBio      | Pythia  | 8                      | 16   | 4  |
|             | Gemma   | 4                      | 8    | 2  |
|             | Llama2  | 2                      | 4    | 2  |
|             | Qwen2.5 | 2                      | 4    | 2  |

Table 3: Batch sizes for CustomerSim and SynBio

## F PROMPT FOR ANNOTATING PRIVACY-SENSITIVE INFORMATION

We used the following prompt to obtain the privacy-sensitive information from GPT-4o for CustomerSim:

*“You are a private data identifying bot. You will be provided a conversation between a user and a customer service bot. You need to identify ALL instances of private data in the conversation. Private data includes any information that can be used to identify an individual, such as names, phone numbers, email addresses, and locations. It can also include any other sensitive information, such as credit card numbers, social security numbers, tracking data, health information etc.. Identify all of these use your discretion to identify any other information that can be used to identify an individual.”*

## G HUMAN SURVEY: ASSESSING THE QUALITY OF PRIVACY-SENSITIVE ANNOTATIONS FROM PRESIDIO AND GPT-4

We evaluated privacy-sensitive annotations from Presidio and GPT-4 through a Prolific survey, each with 10 samples for 40 participants. Each of the 10 instances included 5 samples from each tool, displayed randomly. Participants had four options: (1) *Accurate* (all sensitive info correctly annotated), (2) *Under-Annotated* (some sensitive info missed), (3) *Over-Annotated* (non-sensitive info included), and (4) *Mixed-Annotated* (both missed and non-sensitive info annotated). Figure 8 shows sample screenshots of the conducted survey. Figures 9a and 9b shows GPT-4 annotations as more “Accurate” than Presidio.

## H WHY DISTINGUISH BETWEEN SENSITIVE AND NON-SENSITIVE ENTITIES?

Figure 10 shows the assessment of privacy and utility by distinguishing between sensitive and non-sensitive tokens for all models. We also show the utility measurements in log-scale plots in Figure 11 for better understanding.

## I EXPLORING THE TRADEOFFS BETWEEN PRIVACY AND UTILITY

**Full Fine-tuning:** We see a similar trade-off between utility and privacy when using Presidio annotations in Figure 12, as observed in Figure 3. For the CustomerSim dataset, we notice that privacy decreases at the beginning of training, while utility improves. However, as training progresses, both metrics start to worsen. In contrast, the SynBio dataset shows a decline in both utility and privacy,



Figure 8: Annotators were provided with these instructions and guidelines to mark the samples as accurately annotated, under-annotated, over-annotated or mixed-annotated.

except for the Pythia model, which initially experiences an improvement in utility during the first few epochs.

972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025

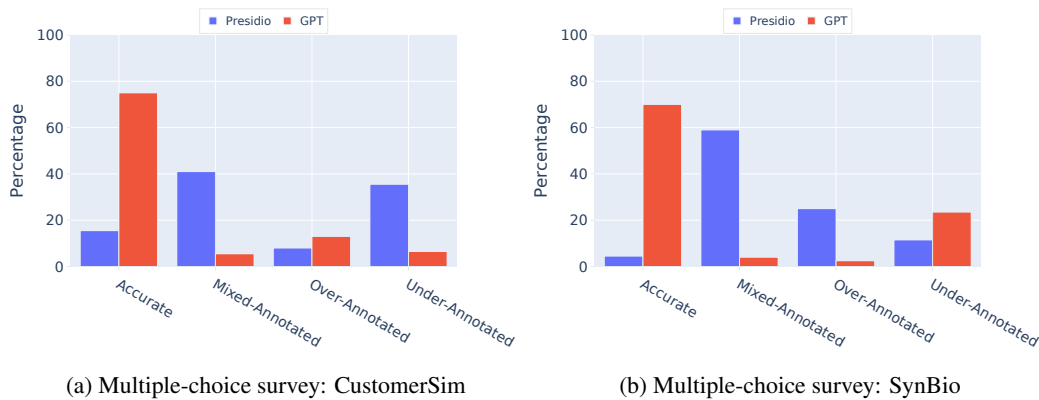


Figure 9: GPT-4 was preferred in both CustomerSim and SynBio by human annotators for identifying privacy-sensitive information.

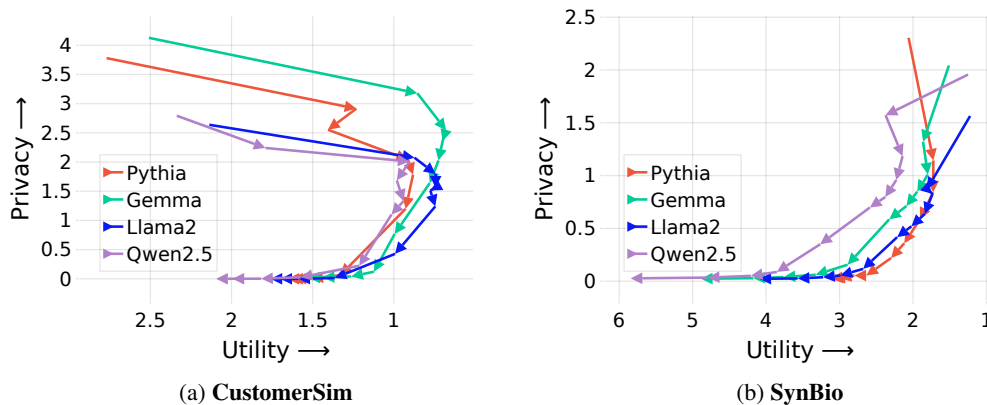


Figure 12: **Full fine-tuning** on (a) CustomerSim and (b) SynBio using Presidio to annotate the privacy-sensitive information.

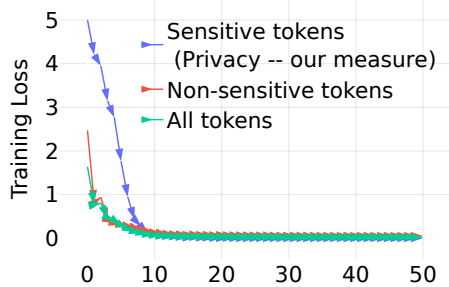
**DP-SGD** : A similar pattern as in Figure 4 using GPT-4 annotations, has also emerged in Figure 13 using Presidio annotations. We observe that in the CustomerSim dataset, DP maintains privacy effectively with minimal degradation. However, when it comes to utility, we notice that lower noise values lead to better utility, but also result in a decrease in privacy. The same holds true for the SynBio dataset in Figure 14, with the only exception being Gemma experiencing a decline in utility after a few iterations.

**LoRA** : Comparing the results from Figures 15 and 16 using Presidio annotations with Figure 5, which used GPT-4 annotations, we can observe the same trend. This suggests that the findings are consistent across different annotation methods. Similarly, when we analyze the trade-off for rank 32 for both the datasets in Figures 17 - 20, as we observed in Figure 5 using GPT-4 annotations.

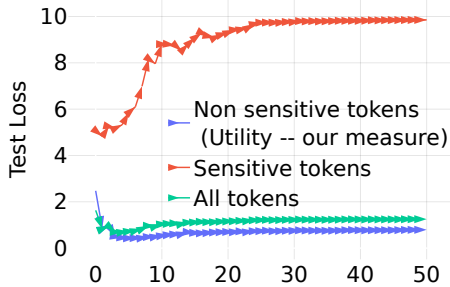
We analyze the trade-offs between privacy, utility, and efficiency for different privacy measures over SynBio with both GPT4 and Presidio annotations. Figures 21, 22, 23 for CustomerSim and SynBio using Presidio annotations and our privacy measure reveals the same as observed in Figure 6 with GPT-4 annotations. We also show the same for the other privacy measures – (a) privacy loss in Figures 24, 25, and 26, and (b) exposure in Figure 27.

The overall takeaway holds consistent – LoRA is the most efficient method, followed by FFT while DP is the least efficient, requiring the highest number of FLOPs. For privacy, DP and LoRA perform

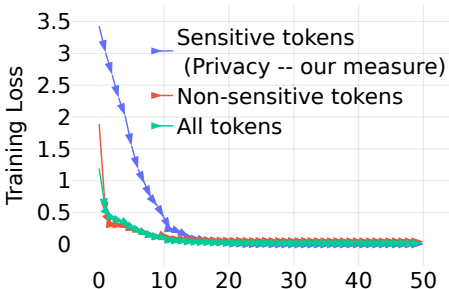
1026  
1027  
1028  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047  
1048  
1049  
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079



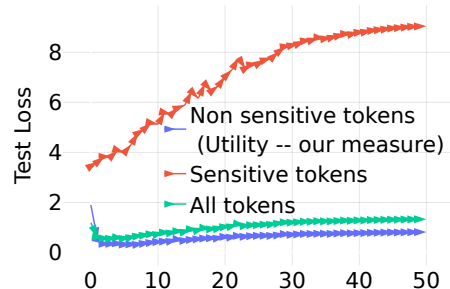
(a) Privacy measure in Pythia



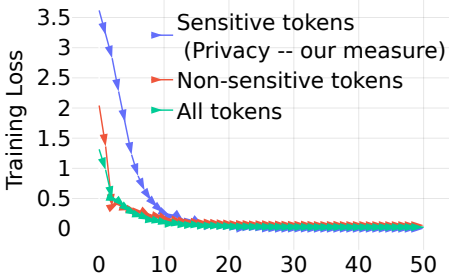
(b) Utility measure in Pythia



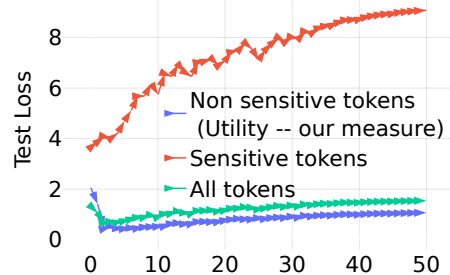
(c) Privacy measure in Llama2



(d) Utility measure in Llama2



(e) Privacy measure in Qwen2.5



(f) Utility measure in Qwen2.5

Figure 10: Our measures offer a precise assessment of privacy and utility by distinguishing between sensitive and non-sensitive tokens.

at par with FFT leaking the most amount of information. In terms of utility, LoRA and FFT perform better while DP deteriorates on large scale models. Additionally, LoRA also maintains benchmark performance across other models as shown in Figure 28.

## J COMPARISON OF FINE-TUNING METHODS

Based on our results, we provide a taxonomy on the axes of privacy, utility, and efficiency for all the three fine-tuning methods in Table 4.

1080  
1081  
1082  
1083  
1084  
1085  
1086  
1087  
1088  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131  
1132  
1133

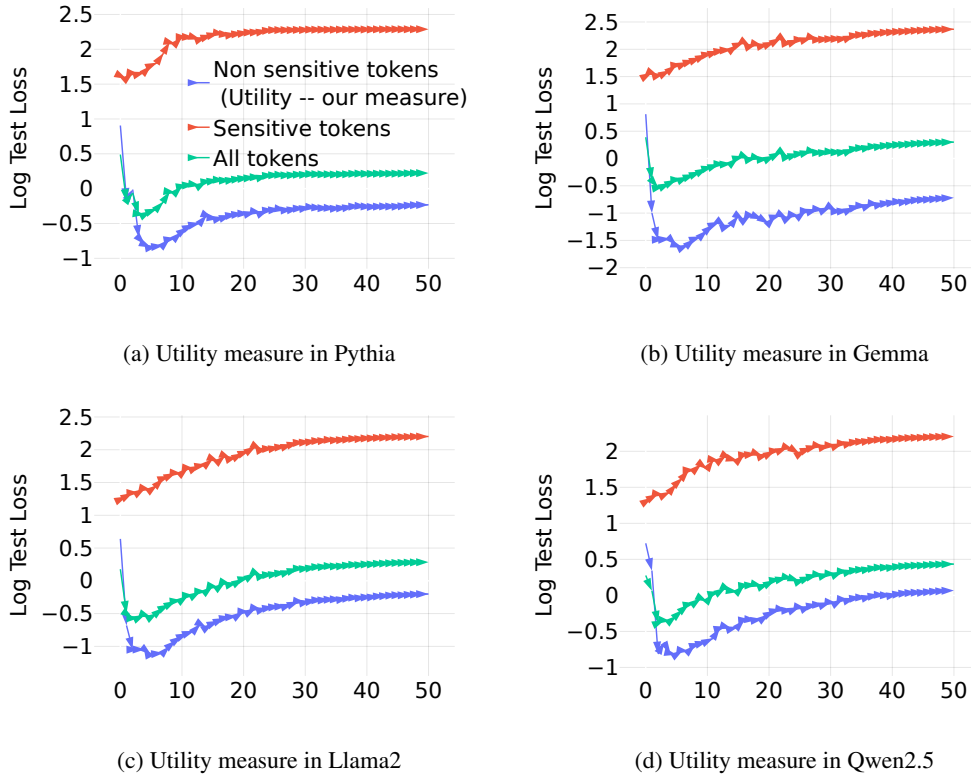


Figure 11: Our measures offer a more precise assessment of utility when fine-tuning LLMs by distinguishing between sensitive and non-sensitive tokens, revealing better utility (lower loss) for non-sensitive tokens compared to scenarios that overlook the sensitivity-based distinction.

|      | Utility  | Privacy | Efficiency | Benchmark |
|------|----------|---------|------------|-----------|
| FFT  | Good     | Poor    | Moderate   | Poor      |
| DP   | Moderate | Good    | Poor       | Poor      |
| LoRA | Good     | Good    | Good       | Good      |

Table 4: Comparison of all fine-tuning methods across privacy, utility, efficiency, and benchmark performance. LoRA is seen to outperform FFT and DP over all dimensions.

1134  
 1135  
 1136  
 1137  
 1138  
 1139  
 1140  
 1141  
 1142  
 1143  
 1144  
 1145  
 1146  
 1147  
 1148  
 1149  
 1150  
 1151  
 1152  
 1153  
 1154  
 1155  
 1156  
 1157  
 1158  
 1159  
 1160  
 1161  
 1162  
 1163  
 1164  
 1165  
 1166  
 1167  
 1168  
 1169  
 1170  
 1171  
 1172  
 1173  
 1174  
 1175  
 1176  
 1177  
 1178  
 1179  
 1180  
 1181  
 1182  
 1183  
 1184  
 1185  
 1186  
 1187

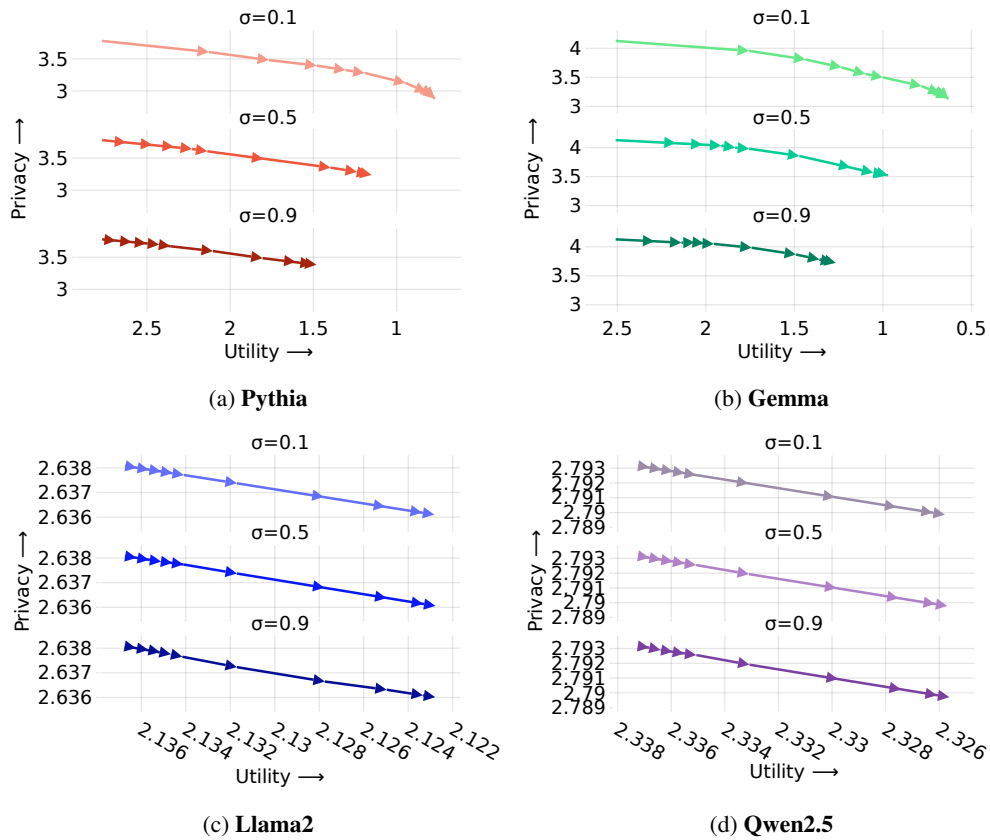


Figure 13: **DP** on *CustomerSim* dataset using Presidio tool for annotating privacy-sensitive information.

1188  
 1189  
 1190  
 1191  
 1192  
 1193  
 1194  
 1195  
 1196  
 1197  
 1198  
 1199  
 1200  
 1201  
 1202  
 1203  
 1204  
 1205  
 1206  
 1207  
 1208  
 1209  
 1210  
 1211  
 1212  
 1213  
 1214  
 1215  
 1216  
 1217  
 1218  
 1219  
 1220  
 1221  
 1222  
 1223  
 1224  
 1225  
 1226  
 1227  
 1228  
 1229  
 1230  
 1231  
 1232  
 1233  
 1234  
 1235  
 1236  
 1237  
 1238  
 1239  
 1240  
 1241

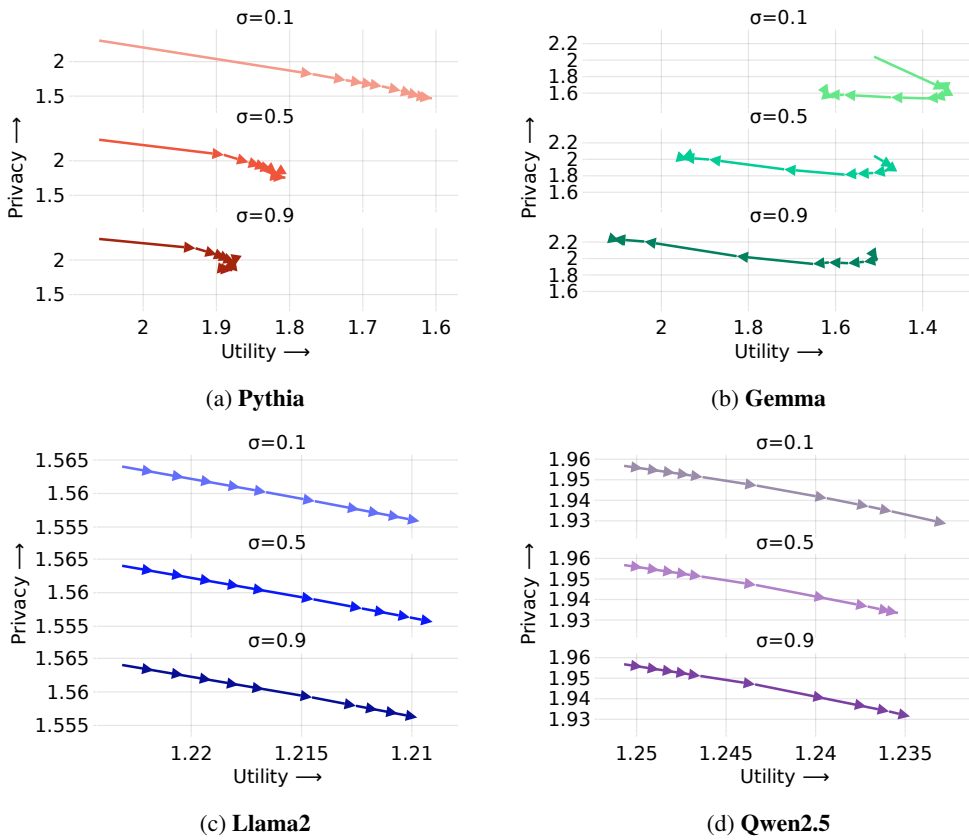


Figure 14: DP on *SynBio* dataset using Presidio tool for annotating privacy-sensitive information.

1242  
 1243  
 1244  
 1245  
 1246  
 1247  
 1248  
 1249  
 1250  
 1251  
 1252  
 1253  
 1254  
 1255  
 1256  
 1257  
 1258  
 1259  
 1260  
 1261  
 1262  
 1263  
 1264  
 1265  
 1266  
 1267  
 1268  
 1269  
 1270  
 1271  
 1272  
 1273  
 1274  
 1275  
 1276  
 1277  
 1278  
 1279  
 1280  
 1281  
 1282  
 1283  
 1284  
 1285  
 1286  
 1287  
 1288  
 1289  
 1290  
 1291  
 1292  
 1293  
 1294  
 1295

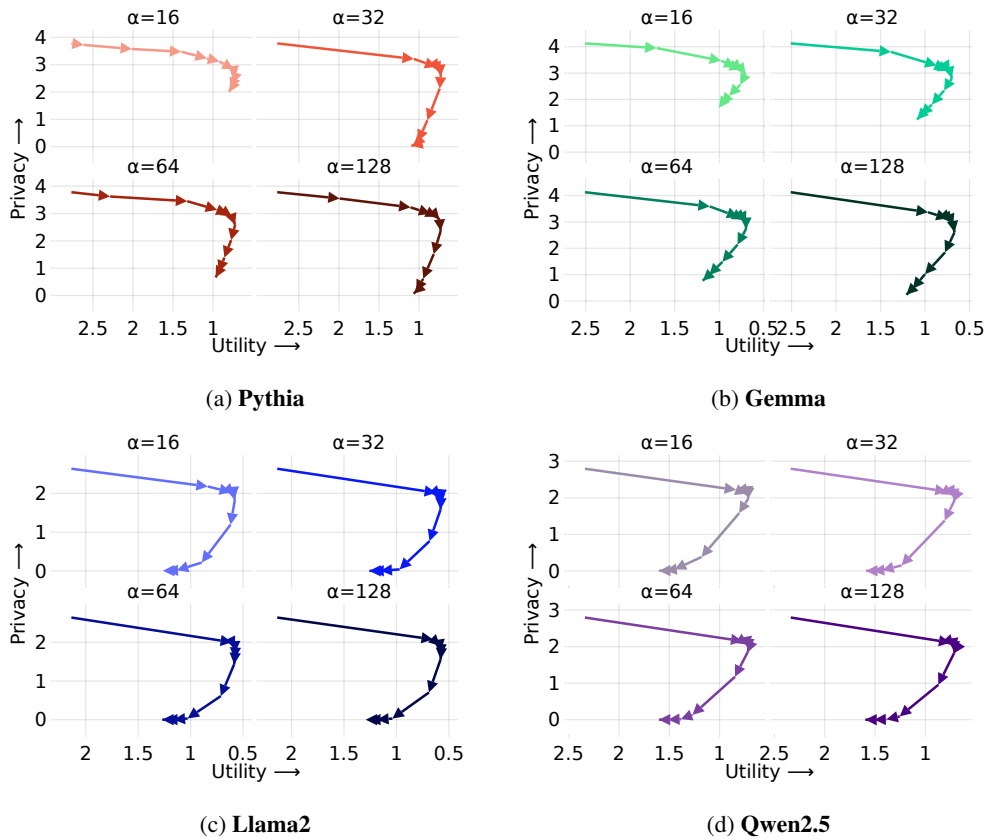


Figure 15: **Lora with rank 16** on *CustomerSim* dataset using Presidio tool for annotating privacy-sensitive information.

1296  
 1297  
 1298  
 1299  
 1300  
 1301  
 1302  
 1303  
 1304  
 1305  
 1306  
 1307  
 1308  
 1309  
 1310  
 1311  
 1312  
 1313  
 1314  
 1315  
 1316  
 1317  
 1318  
 1319  
 1320  
 1321  
 1322  
 1323  
 1324  
 1325  
 1326  
 1327  
 1328  
 1329  
 1330  
 1331  
 1332  
 1333  
 1334  
 1335  
 1336  
 1337  
 1338  
 1339  
 1340  
 1341  
 1342  
 1343  
 1344  
 1345  
 1346  
 1347  
 1348  
 1349

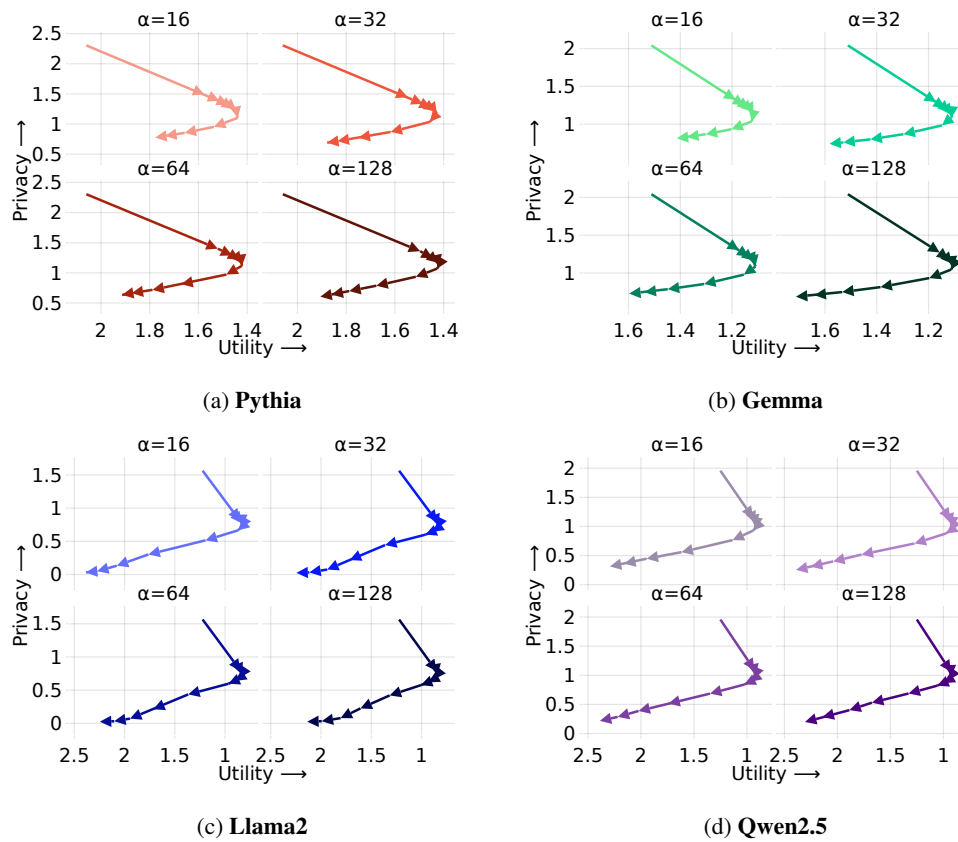


Figure 16: Lora with rank 16 on *SynBio* dataset using Presidio tool for annotating privacy-sensitive information.

1350  
 1351  
 1352  
 1353  
 1354  
 1355  
 1356  
 1357  
 1358  
 1359  
 1360  
 1361  
 1362  
 1363  
 1364  
 1365  
 1366  
 1367  
 1368  
 1369  
 1370  
 1371  
 1372  
 1373  
 1374  
 1375  
 1376  
 1377  
 1378  
 1379  
 1380  
 1381  
 1382  
 1383  
 1384  
 1385  
 1386  
 1387  
 1388  
 1389  
 1390  
 1391  
 1392  
 1393  
 1394  
 1395  
 1396  
 1397  
 1398  
 1399  
 1400  
 1401  
 1402  
 1403

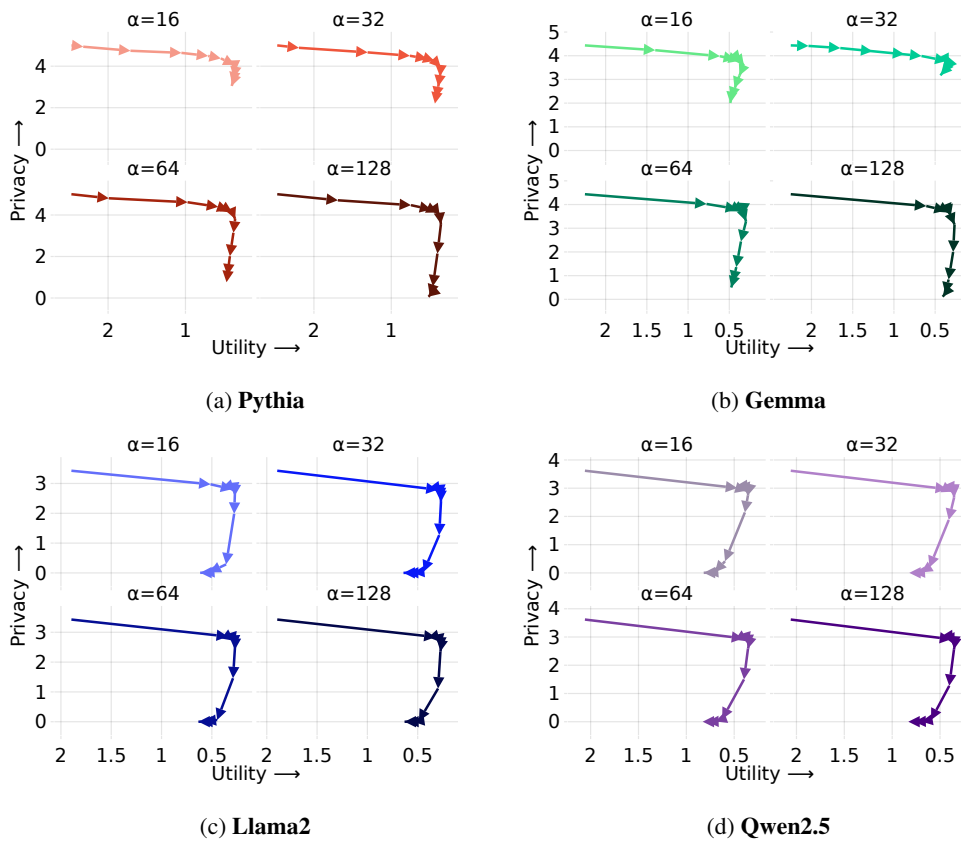


Figure 17: Lora with rank 32 on *CustomerSim* dataset using GPT4 for annotating privacy-sensitive information.

1404  
 1405  
 1406  
 1407  
 1408  
 1409  
 1410  
 1411  
 1412  
 1413  
 1414  
 1415  
 1416  
 1417  
 1418  
 1419  
 1420  
 1421  
 1422  
 1423  
 1424  
 1425  
 1426  
 1427  
 1428  
 1429  
 1430  
 1431  
 1432  
 1433  
 1434  
 1435  
 1436  
 1437  
 1438  
 1439  
 1440  
 1441  
 1442  
 1443  
 1444  
 1445  
 1446  
 1447  
 1448  
 1449  
 1450  
 1451  
 1452  
 1453  
 1454  
 1455  
 1456  
 1457

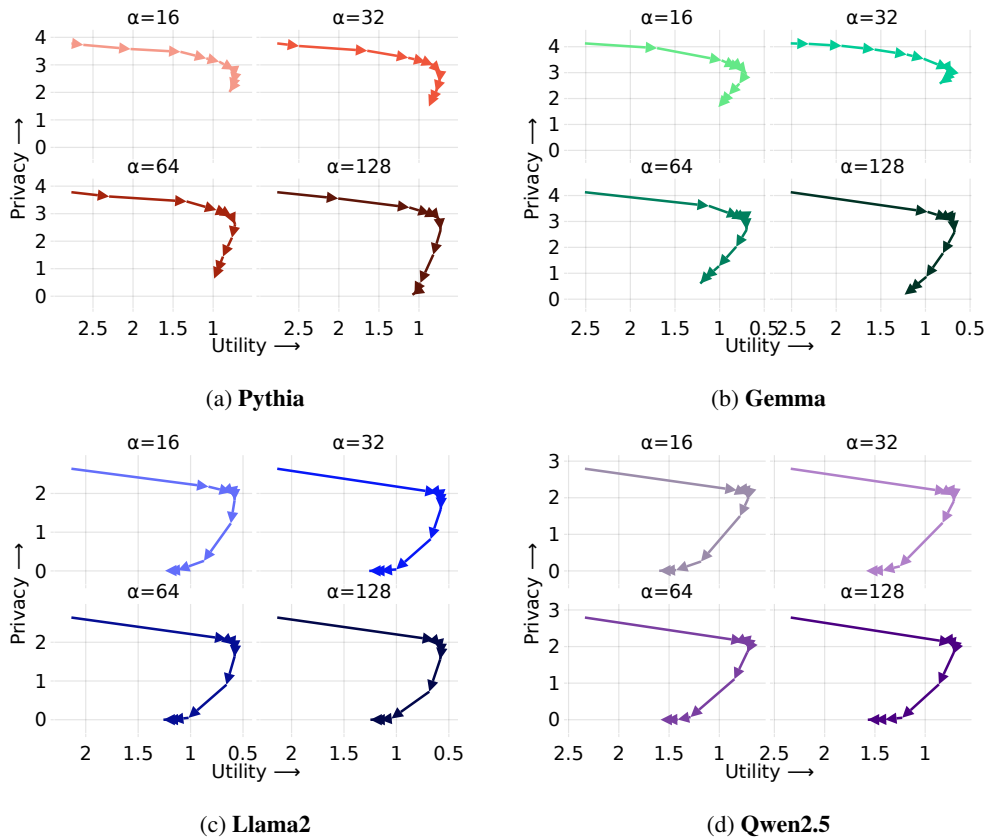


Figure 18: **Lora with rank 32** on *CustomerSim* dataset using Presidio tool for annotating privacy-sensitive information.

1458  
 1459  
 1460  
 1461  
 1462  
 1463  
 1464  
 1465  
 1466  
 1467  
 1468  
 1469  
 1470  
 1471  
 1472  
 1473  
 1474  
 1475  
 1476  
 1477  
 1478  
 1479  
 1480  
 1481  
 1482  
 1483  
 1484  
 1485  
 1486  
 1487  
 1488  
 1489  
 1490  
 1491  
 1492  
 1493  
 1494  
 1495  
 1496  
 1497  
 1498  
 1499  
 1500  
 1501  
 1502  
 1503  
 1504  
 1505  
 1506  
 1507  
 1508  
 1509  
 1510  
 1511

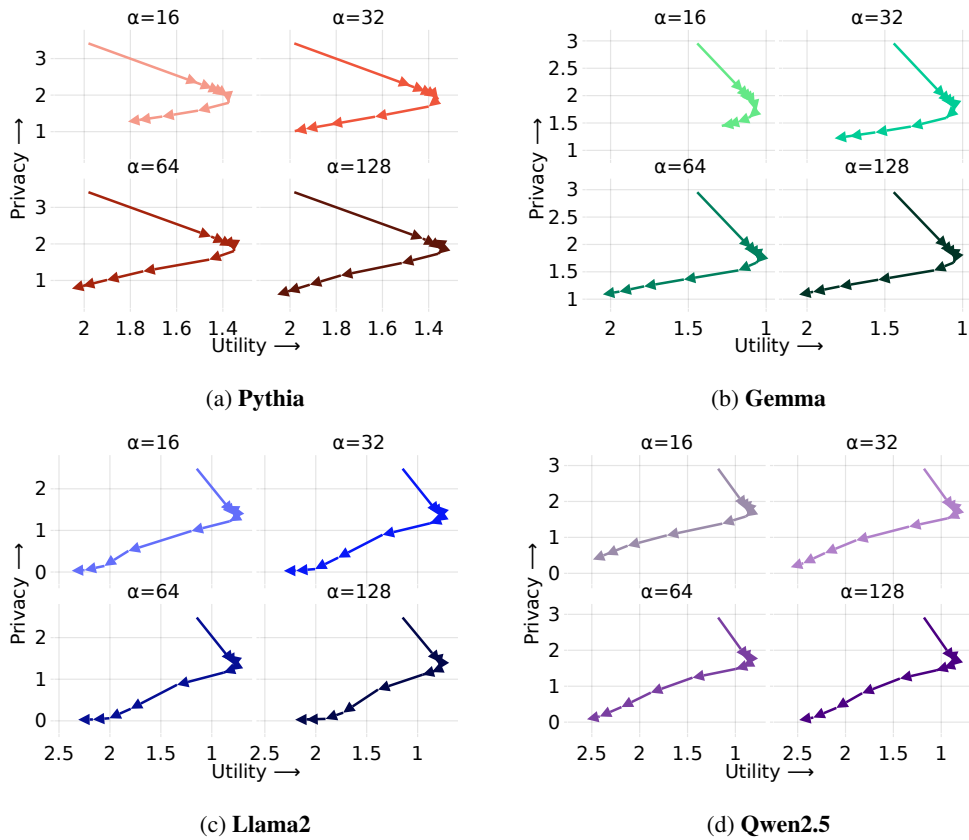


Figure 19: Lora with rank 32 on *SynBio* dataset using GPT4 for annotating privacy-sensitive information.

1512  
 1513  
 1514  
 1515  
 1516  
 1517  
 1518  
 1519  
 1520  
 1521  
 1522  
 1523  
 1524  
 1525  
 1526  
 1527  
 1528  
 1529  
 1530  
 1531  
 1532  
 1533  
 1534  
 1535  
 1536  
 1537  
 1538  
 1539  
 1540  
 1541  
 1542  
 1543  
 1544  
 1545  
 1546  
 1547  
 1548  
 1549  
 1550  
 1551  
 1552  
 1553  
 1554  
 1555  
 1556  
 1557  
 1558  
 1559  
 1560  
 1561  
 1562  
 1563  
 1564  
 1565

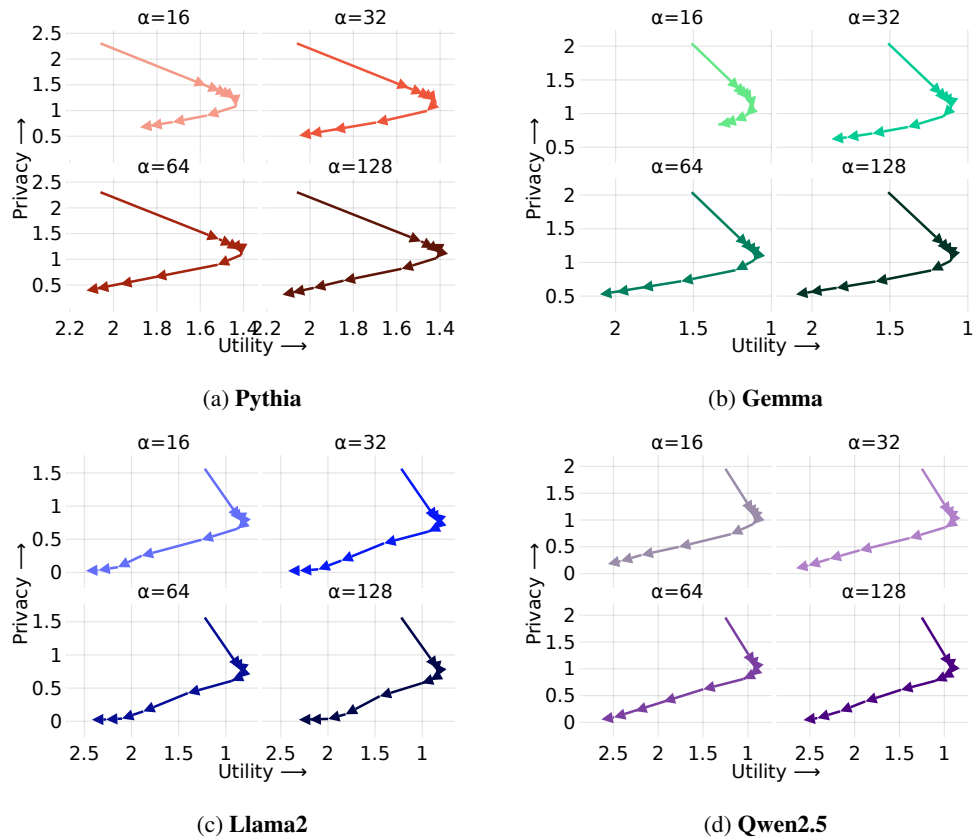


Figure 20: Lora with rank 32 on *SynBio* dataset using Presidio tool for annotating privacy-sensitive information.

1566  
 1567  
 1568  
 1569  
 1570  
 1571  
 1572  
 1573  
 1574  
 1575  
 1576  
 1577  
 1578  
 1579  
 1580  
 1581  
 1582  
 1583  
 1584  
 1585  
 1586  
 1587  
 1588  
 1589  
 1590  
 1591  
 1592  
 1593  
 1594  
 1595  
 1596  
 1597  
 1598  
 1599  
 1600  
 1601  
 1602  
 1603  
 1604  
 1605  
 1606  
 1607  
 1608  
 1609  
 1610  
 1611  
 1612  
 1613  
 1614  
 1615  
 1616  
 1617  
 1618  
 1619

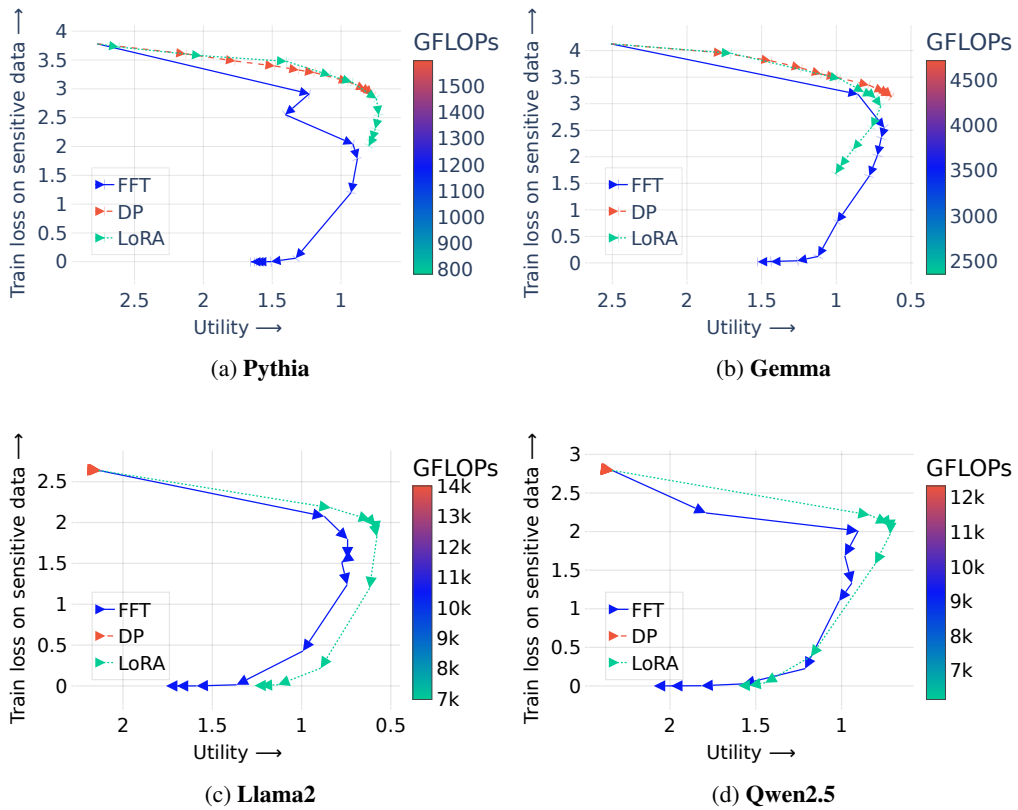


Figure 21: Full fine-tuning, DP and LoRA on *CustomerSim* with Presidio annotations for privacy-sensitive information.

1620  
1621  
1622  
1623  
1624  
1625  
1626  
1627  
1628  
1629  
1630  
1631  
1632  
1633  
1634  
1635  
1636  
1637  
1638  
1639  
1640  
1641  
1642  
1643  
1644  
1645  
1646  
1647  
1648  
1649  
1650  
1651  
1652  
1653  
1654  
1655  
1656  
1657  
1658  
1659  
1660  
1661  
1662  
1663  
1664  
1665  
1666  
1667  
1668  
1669  
1670  
1671  
1672  
1673

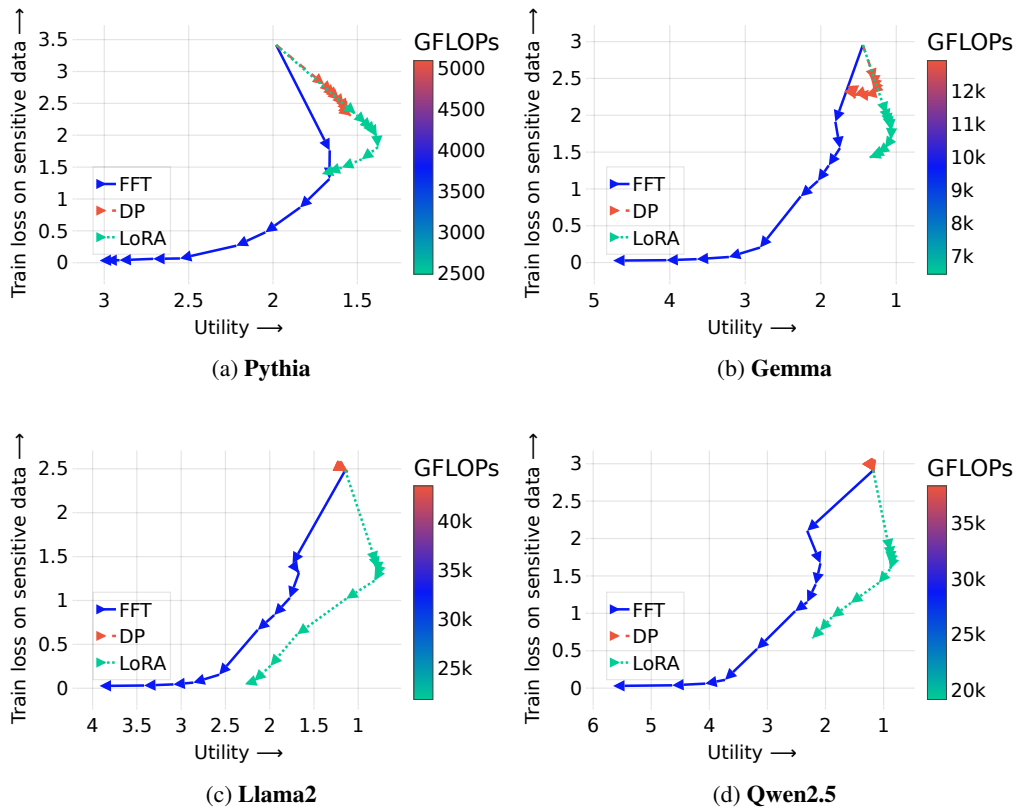


Figure 22: Full fine-tuning, DP and LoRA on *SynBio* with GPT4 annotations for privacy-sensitive information.

1674  
 1675  
 1676  
 1677  
 1678  
 1679  
 1680  
 1681  
 1682  
 1683  
 1684  
 1685  
 1686  
 1687  
 1688  
 1689  
 1690  
 1691  
 1692  
 1693  
 1694  
 1695  
 1696  
 1697  
 1698  
 1699  
 1700  
 1701  
 1702  
 1703  
 1704  
 1705  
 1706  
 1707  
 1708  
 1709  
 1710  
 1711  
 1712  
 1713  
 1714  
 1715  
 1716  
 1717  
 1718  
 1719  
 1720  
 1721  
 1722  
 1723  
 1724  
 1725  
 1726  
 1727

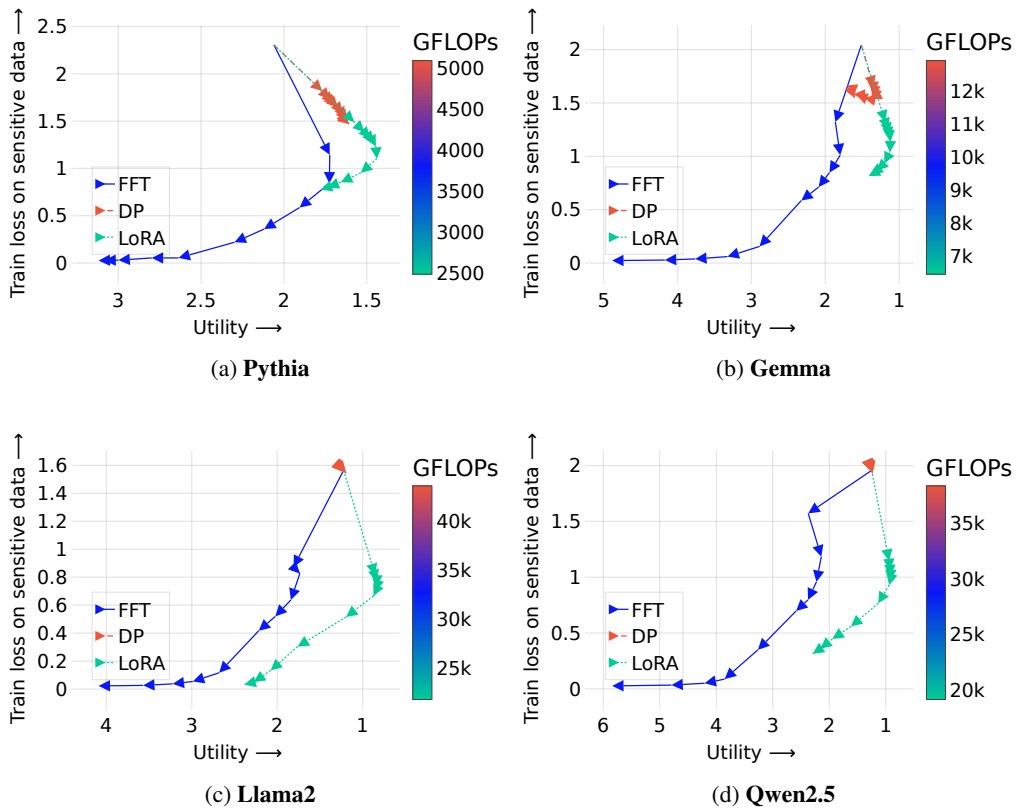


Figure 23: **Full fine-tuning, DP and LoRA** on *SynBio* with Presidio annotations for privacy-sensitive information.

1728  
 1729  
 1730  
 1731  
 1732  
 1733  
 1734  
 1735  
 1736  
 1737  
 1738  
 1739  
 1740  
 1741  
 1742  
 1743  
 1744  
 1745  
 1746  
 1747  
 1748  
 1749  
 1750  
 1751  
 1752  
 1753  
 1754  
 1755  
 1756  
 1757  
 1758  
 1759  
 1760  
 1761  
 1762  
 1763  
 1764  
 1765  
 1766  
 1767  
 1768  
 1769  
 1770  
 1771  
 1772  
 1773  
 1774  
 1775  
 1776  
 1777  
 1778  
 1779  
 1780  
 1781

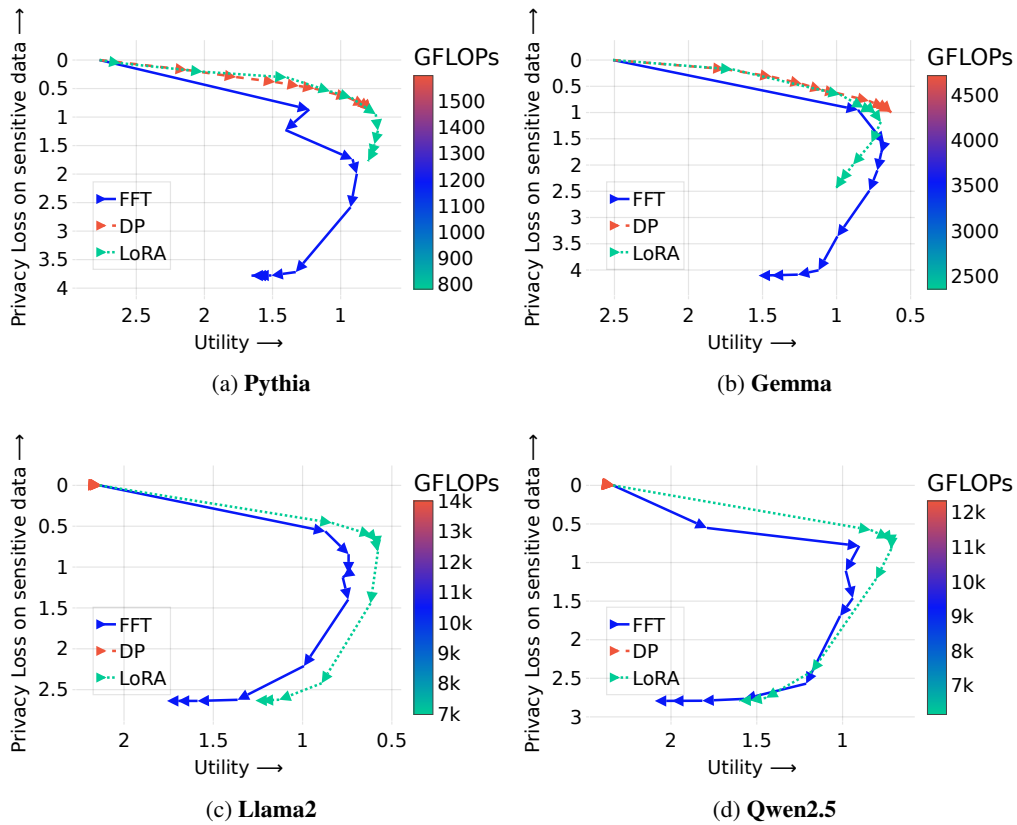


Figure 24: Full fine-tuning, DP and LoRA on *CustomerSim* with Presidio annotations for privacy loss

1782  
 1783  
 1784  
 1785  
 1786  
 1787  
 1788  
 1789  
 1790  
 1791  
 1792  
 1793  
 1794  
 1795  
 1796  
 1797  
 1798  
 1799  
 1800  
 1801  
 1802  
 1803  
 1804  
 1805  
 1806  
 1807  
 1808  
 1809  
 1810  
 1811  
 1812  
 1813  
 1814  
 1815  
 1816  
 1817  
 1818  
 1819  
 1820  
 1821  
 1822  
 1823  
 1824  
 1825  
 1826  
 1827  
 1828  
 1829  
 1830  
 1831  
 1832  
 1833  
 1834  
 1835

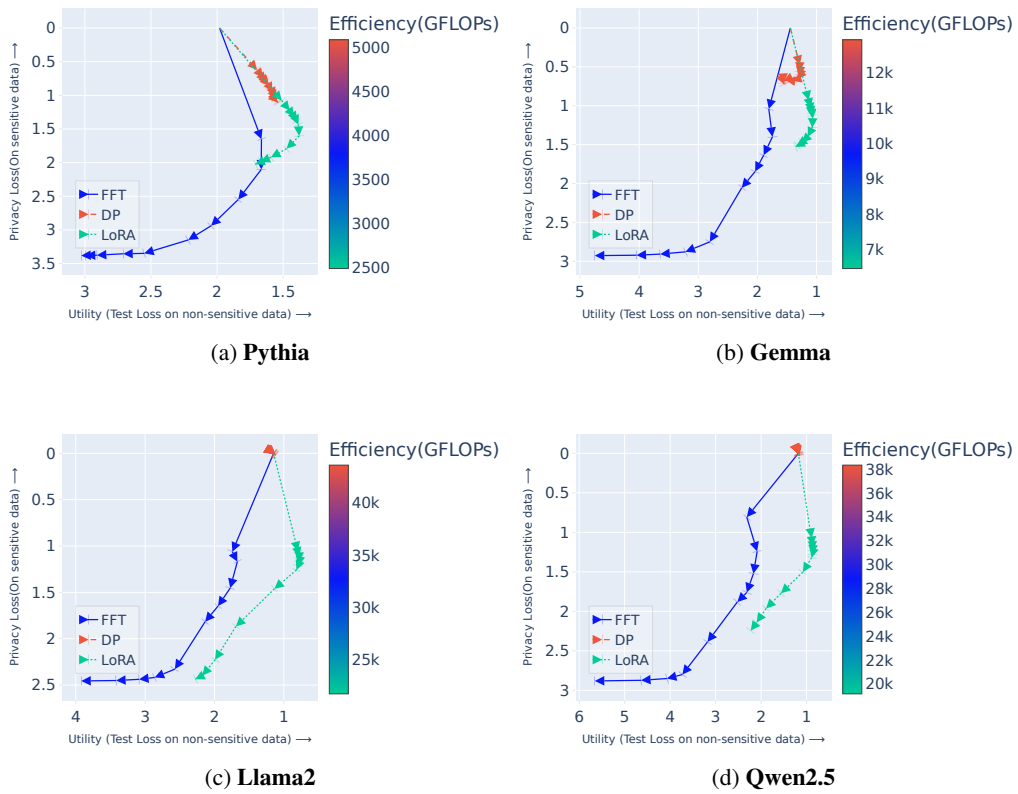


Figure 25: Full fine-tuning, DP and LoRA on *SynBio* with GPT4 annotations for privacy loss

1836  
 1837  
 1838  
 1839  
 1840  
 1841  
 1842  
 1843  
 1844  
 1845  
 1846  
 1847  
 1848  
 1849  
 1850  
 1851  
 1852  
 1853  
 1854  
 1855  
 1856  
 1857  
 1858  
 1859  
 1860  
 1861  
 1862  
 1863  
 1864  
 1865  
 1866  
 1867  
 1868  
 1869  
 1870  
 1871  
 1872  
 1873  
 1874  
 1875  
 1876  
 1877  
 1878  
 1879  
 1880  
 1881  
 1882  
 1883  
 1884  
 1885  
 1886  
 1887  
 1888  
 1889

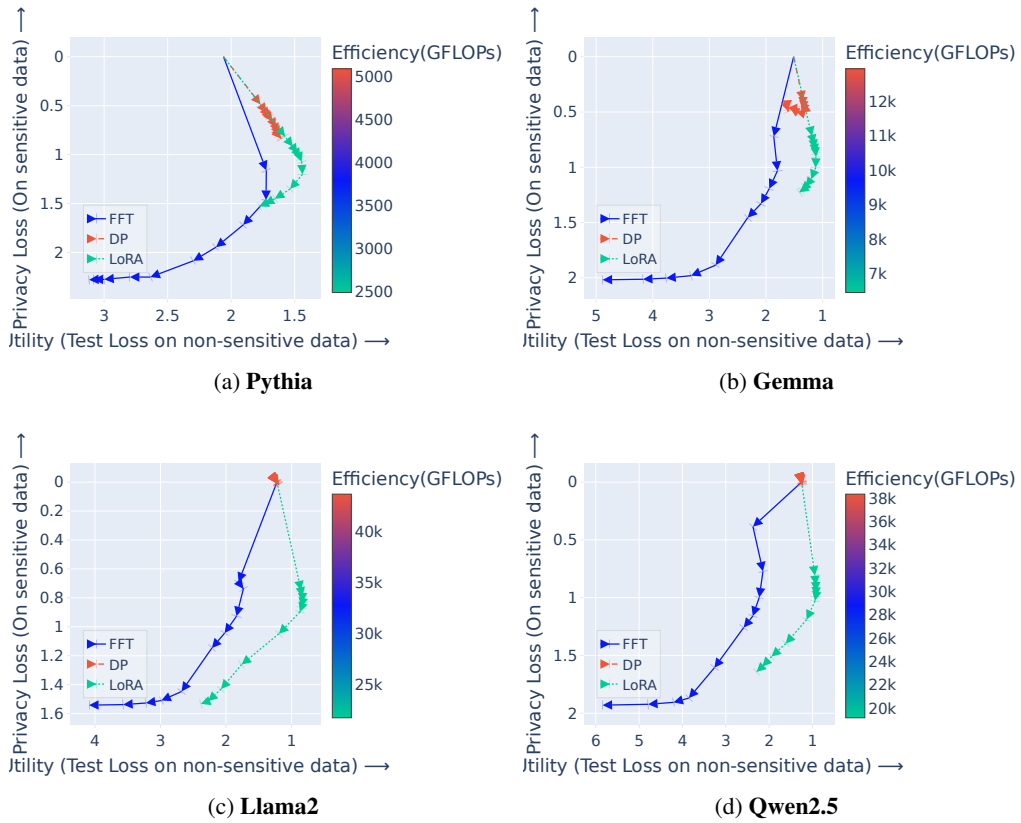


Figure 26: Full fine-tuning, DP and LoRA on *SynBio* with Presidio annotations for privacy loss

1890  
 1891  
 1892  
 1893  
 1894  
 1895  
 1896  
 1897  
 1898  
 1899  
 1900  
 1901  
 1902  
 1903  
 1904  
 1905  
 1906  
 1907  
 1908  
 1909  
 1910  
 1911  
 1912  
 1913  
 1914  
 1915  
 1916  
 1917  
 1918  
 1919  
 1920  
 1921  
 1922  
 1923  
 1924  
 1925  
 1926  
 1927  
 1928  
 1929  
 1930  
 1931  
 1932  
 1933  
 1934  
 1935  
 1936  
 1937  
 1938  
 1939  
 1940  
 1941  
 1942  
 1943

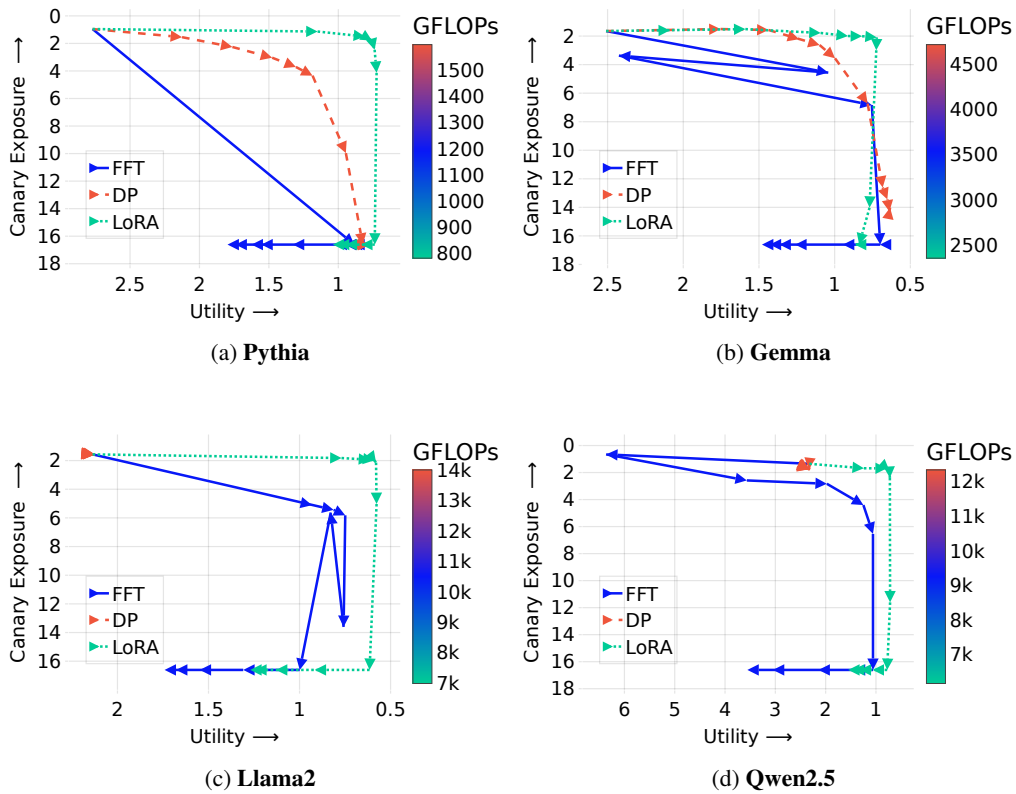
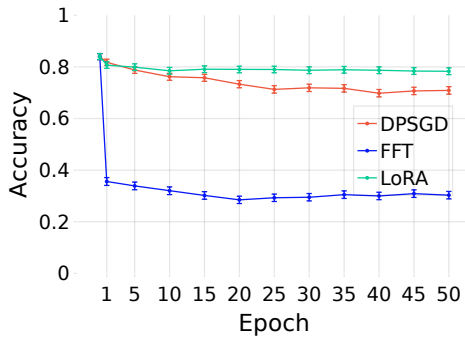
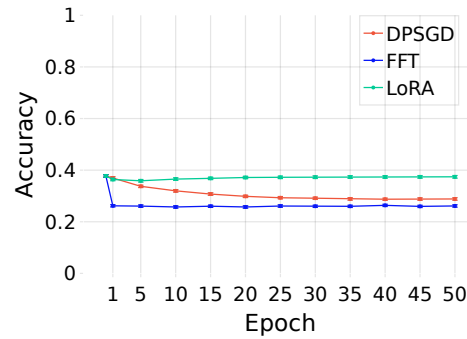


Figure 27: **Full fine-tuning, DP and LoRA** on *CustomerSim* with Presidio annotations for canary exposure.

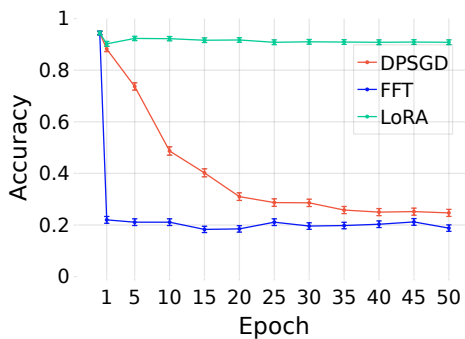
1944  
1945  
1946  
1947  
1948  
1949  
1950  
1951  
1952  
1953  
1954  
1955  
1956  
1957  
1958  
1959  
1960  
1961  
1962  
1963  
1964  
1965  
1966  
1967  
1968  
1969  
1970  
1971  
1972  
1973  
1974  
1975  
1976  
1977  
1978  
1979  
1980  
1981  
1982  
1983  
1984  
1985  
1986  
1987  
1988  
1989  
1990  
1991  
1992  
1993  
1994  
1995  
1996  
1997



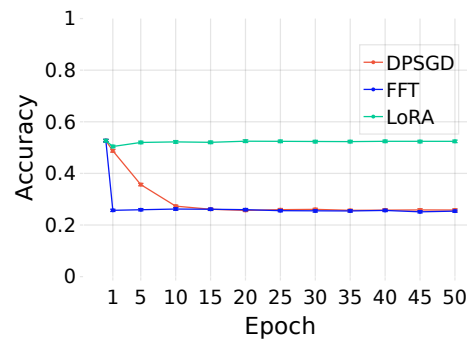
(a) Pythia : SCIQ



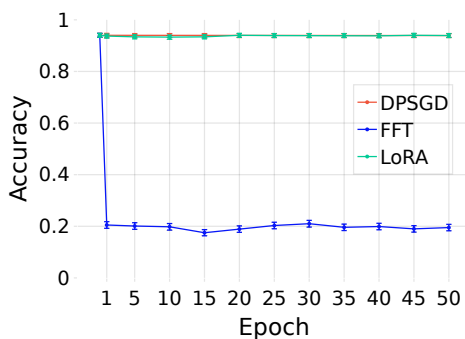
(b) Pythia : HellaSwag



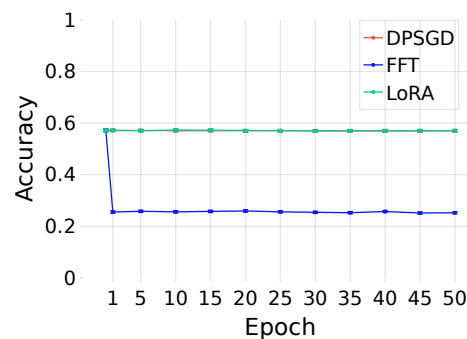
(c) Gemma : SCIQ



(d) Gemma : HellaSwag



(e) Llama2 : SCIQ



(f) Llama 2 : HellaSwag

Figure 28: LoRA demonstrates knowledge retention on benchmark datasets throughout the training regime on CustomerSim for other models as well, while FFT and DP show fragility with sharp and gradual performance declines, respectively.