

# Large Language Models Can Be Contextual Privacy Protection Learners

Anonymous ACL submission

## Abstract

The proliferation of Large Language Models (LLMs) has driven considerable interest in fine-tuning them with domain-specific data to create specialized language models. Nevertheless, such domain-specific fine-tuning data often contains *contextually sensitive* personally identifiable information (PII). Direct fine-tuning LLMs on this data without privacy protection poses a risk of data leakage of sensitive PII during inference time. To address this challenge, we introduce Contextual Privacy Protection Language Models (CPPLM), a novel paradigm for fine-tuning LLMs that effectively injects domain-specific knowledge while safeguarding inference-time data privacy. Our work offers a theoretical analysis for model design and delves into various techniques such as corpus curation, penalty-based unlikelihood in training loss, and instruction-based tuning, etc. Extensive experiments across diverse datasets and scenarios demonstrate the effectiveness of our approaches. In particular, instruction tuning with both positive and negative examples, stands out as a promising method, effectively protecting private data while enhancing the model’s knowledge. Our work underscores the potential for Large Language Models as robust contextual privacy protection learners.

## 1 Introduction

**Background.** Large Language Models (LLMs) have demonstrated remarkable linguistic comprehension and generation capability (Bang et al., 2023; Wang et al., 2023a). Meanwhile, when directly applied to specialized industries, they encounter challenges such as hallucination (Bang et al., 2023; Chan et al., 2023), insufficient domain expertise (Singhal et al., 2023b), and failing to incorporate the latest domain knowledge in ever-evolving industry scenarios (Kasneci et al., 2023). The introduction of open-source general-purpose LLMs such as LLaMA (Touvron et al.,

2023) and RWKV (Peng et al., 2023) have provided a promising solution. Researchers would fine-tune specialized LLMs based on powerful general-purpose LLMs using high-quality, domain-specific knowledge to ensure both commonsense reasoning and comprehensive knowledge coverage (Hoffmann et al., 2022a; Villalobos et al., 2022; Hoffmann et al., 2022b). Such examples include BloombergGPT (Wu et al., 2023) and MedPaLM (Singhal et al., 2023a), for financial and medical applications, respectively. However, these fine-tuning datasets usually contain sensitive information, such as personally identifiable information (PII) (Carlini et al., 2020; Lin et al., 2021; Gehman et al., 2020). When applied to downstream tasks, sensitive information in the training data, such as social security numbers or patient names, can be exposed by the LLMs upon text generation, a phenomenon known as the memorization effect (Yu et al., 2023b; Kenton and Toutanova, 2019; Meng et al., 2023) or inference-time privacy threat (Miresghallah et al., 2024), leading to identity theft and financial losses (Coavoux et al., 2018; Yu et al., 2023a).

**Challenges.** In this work, we aim to tackle the challenging task of efficient LLM fine-tuning for enhanced *contextual privacy* (Nissenbaum, 2004; Miresghallah et al., 2024), a critical yet under-explored setting where the sensitivity of a piece of information is contingent upon the context. For example, statements such as “Bill Gates founded Microsoft” and “Alan Mathison Turing was an English mathematician and computer scientist” are generally not considered violations of privacy, since they are presented as common knowledge. In contrast, statements like “Alan Gates visited the X hospital for a certain disease Y” pose privacy concerns as they reveal details about individuals’ daily activities and health status in a particular context. Directly applying techniques like Named Entity Recognition (NER) can lead to inaccurate identi-

083 fication of PII, whereas merely deleting or mask- 134  
084 ing PII tokens in the training data would result 135  
085 in a substantial information loss and compromise 136  
086 the performance on downstream tasks — a co- 137  
087 nundrum known as the privacy-utility trade-off as 138  
088 theoretically discussed in Sec. 4.1. An alterna- 139  
089 tive approach, reinforcement learning from human 140  
090 feedback (RLHF), involves additional model fine- 141  
091 tuning guided by human feedback (Ouyang et al., 142  
092 2022) so that the model tends towards concealing 143  
093 sensitive PII (like “red-teaming”). For example, 144  
094 it learns to prioritize outputs that protect sensitive 145  
095 PII over those that leak PII. Nonetheless, RLHF is 146  
096 data-intensive, potentially costly in computation, 147  
097 and can pose stability challenges (Ziegler et al., 148  
098 2020; Wang et al., 2023b). 149

099 **Our Work.** To address these challenges, this 150  
100 paper introduces effective and efficient methodolo- 151  
101 gies for fine-tuning LLMs to incorporate domain 152  
102 knowledge while ensuring privacy protection. We 153  
103 propose and rigorously examine a diverse suite 154  
104 of strategies from corpus curation, introduction of 155  
105 penalty-based unlikelihood into the training loss, 156  
106 instruction-based tuning, a PII contextual classi- 157  
107 fier, and direct preference optimization (DPO), etc. 158  
108 The ultimate objective is to cultivate a model that 159  
109 excels at acquiring information while demonstrat- 160  
110 ing the ability to distinguish between information 161  
111 that can be openly shared and that demands strict 162  
112 confidentiality. Our experimental findings suggest 163  
113 that instruction tuning with positive and negative 164  
114 examples can offer promising avenues. It not only 165  
115 effectively shields private data but also enables the 166  
116 model to assimilate knowledge from the corpus. 167  
117 This implies that *LLMs can be good contextual*  
118 *privacy protection learners*, without the need for 168  
119 balancing a privacy-utility trade-off. To sum up, 169  
120 our contributions are as follows. 170

121 1). **Novel Methodology.** For the first time, we 171  
122 explicitly address the challenging problem of 172  
123 building Contextual Privacy Protection Language 173  
124 Models (**CPPLM**), a novel paradigm in fine-tuning 174  
125 language models that emphasizes privacy protec- 175  
126 tion of contextual PII. To achieve this, we systemat- 176  
127 ically lay out and empirically test a comprehensive 177  
128 spectrum of strategies. 178

129 2). **Theoretical Guidance.** We provide a theo- 179  
130 retical analysis of our proposed methodologies. 180  
131 This analysis illuminates the pathway to design- 181  
132 ing robust tuning methods, ensuring the resultant 182  
133 language model can both protect private data and

134 assimilate vast knowledge from fine-tuning corpus. 135  
136 3). **Comprehensive Evaluation.** We exten- 137  
138 sively benchmarked our methods on four datasets 139  
139 (biomedical, healthcare, and real-world ones). 140  
140 These experiments demonstrated the efficacy of 141  
141 our fine-tuning method to inject domain knowl- 142  
142 edge and safeguard private personal information 143  
143 (PII). The outcomes show that our technique per- 144  
144 forms significantly better than the baselines. 145  
145

146 Our code and data are available at Anonymous 147  
147 GitHub<sup>1</sup>. We will make all code and the proposed 148  
148 datasets publicly available upon the acceptance of 149  
149 this work. 150

## 151 2 Related Work 152

153 **Large Language Models and Privacy.** In the 154  
154 rapidly advancing domain of artificial intelligence 155  
155 and natural language processing, LLMs such as 156  
156 GPT-3.5/4 (OpenAI, 2023), Bard (Google, 2023), 157  
157 LLaMA (Touvron et al., 2023), and ChatGLM (Du 158  
158 et al., 2022) have become pivotal and have demon- 159  
159 strated unprecedented capabilities in generating 160  
160 coherent and contextually accurate text. How- 161  
161 ever, this widespread application raises significant 162  
162 privacy concerns, particularly regarding personal 163  
163 information protection. Addressing the privacy 164  
164 challenges posed by LLMs, researchers have fo- 165  
165 cused on three primary strategies: (Li et al., 2023; 166  
166 Zhang et al., 2023; Kim et al., 2023; Lukas et al., 167  
167 2023): 1) curation of the pretraining corpus, 2) 168  
168 conditional large language model (LLM) pretrain- 169  
169 ing, and 3) post-training alignment. Our research 170  
170 focuses on enhancing privacy protection in LLMs 171  
171 through fine-tuning methods that enable knowledge 172  
172 injection to safeguard Personally Identifiable Infor- 173  
173 mation (PII) (Lukas et al., 2023), as designated 174  
174 by users. This contrasts with Differential Privacy 175  
175 (DP), which protects against the leakage of entire 176  
176 records at the cost of reduced data utility. Our 177  
177 method emphasizes targeted PII protection, a cru- 178  
178 cial aspect in contexts where knowledge integration 179  
179 is the key to preserving privacy without compro- 180  
180 mising data utility (Yu et al., 2022; Shi et al., 2021; 181  
181 Anil et al., 2022; Li et al., 2022; Liu et al., 2023; 182  
182 Zhao et al., 2022). For the fine-tuning of LLMs, 183  
183 the decline in utility is inversely linked to the pri- 184  
184 vacy budget allocated for safeguarding the entire 185  
185 training document, as it determines the extent of 186  
186 noise introduced (Lukas et al., 2023). Our empha- 187  
187 sis lies in specifically safeguarding the contextual 188

<sup>1</sup><https://anonymous.4open.science/r/PPLM>

183 PII tokens. Since PII are contextual (Mireshghal-  
184 lah et al., 2024; Nissenbaum, 2004), our approach  
185 tunes LLMs with contrastive examples designated  
186 by users can accommodate the customized privacy  
187 preferences.

188 **Filtering.** For the pretraining corpus, manually  
189 detecting and filtering out/revising the corpus can  
190 offer high-quality corpus, which is ideal for train-  
191 ing privacy-preserving LLMs (Hoffmann et al.,  
192 2022b; Villalobos et al., 2022; Lukas et al., 2023).  
193 Nevertheless, it is infeasible to process billions of  
194 tokens manually in practice. Another solution is  
195 using automated tools to filter out all sensitive con-  
196 tent (e.g. names, addresses, phone numbers) from  
197 the pretraining corpus. Automated filters make  
198 it possible to go over pretraining datasets. How-  
199 ever, simply removing or masking the PII tokens  
200 (i.e., PII scrubbing) can cause information loss or  
201 inconsistency in the corpus (Welbl et al., 2021).  
202 Though filters can ‘clean’ datasets, they reduce the  
203 diversity in the corpus, which further negatively  
204 impacts the robustness of LLMs (Hendrycks et al.,  
205 2019). Another solution is adding content filters  
206 on top of the existing LMs to control the content  
207 generation process (Xu et al., 2020). Even so, care-  
208 fully designed cases (e.g. prompts) can still trigger  
209 some undesired behaviors of large LMs (Gehman  
210 et al., 2020; Ziegler et al., 2022). However, di-  
211 rectly removing PII from the training corpus poses  
212 a dilemma. While it ensures the elimination of sen-  
213 sitive data, it also potentially weakens the LLMs  
214 by stripping them of crucial knowledge. The mere  
215 act of omitting data can inadvertently hamper the  
216 model’s capacity to process and understand certain  
217 contexts. Context-awareness is fundamental when  
218 considering privacy protection and what data to  
219 shield.

220 **LLMs Adaptation.** To strike a balance between  
221 performance and flexibility, pretraining large LMs  
222 without constraints and then adjusting them to align  
223 with human preferences is a widely adopted ap-  
224 proach for now. One approach is supervising fine-  
225 tuning. The pre-trained LMs are tuned on curated  
226 datasets in a supervised manner (Solaiman and Den-  
227 nison, 2021; Zhou et al., 2023; Wan et al., 2023;  
228 Jin et al., 2023; Yang et al., 2022). Another ap-  
229 proach is reinforcement learning from human feed-  
230 back (RLHF) (Ouyang et al., 2022; Bai et al., 2022;  
231 Menick et al., 2022; Chen et al., 2023). RLHF  
232 gathers data with feedback/preference labels, trains  
233 a reward model, and then finetunes the LM with

reinforcement learning.

### 3 Problem Statement

**Problem Formulation.** In the context of language  
models, a fine-tuning dataset  $D = \{s\}$  is a col-  
lection of natural language sequences  $s$ . Each  
sequence is denoted as  $s = [w_0, w_2, \dots, w_{n-1}]$ ,  
where  $w_i \in s$  represents a token. For privacy  
protection, the users annotate each sequence in  
the corpus by a binary sequence  $p$  denoted as  
 $p = [p_0, \dots, p_{n-1}]$ ,  $p_i \in \{0, 1\}$ , where  $p_i = 1$   
denotes the token is private tokens (e.g., PII) need  
to be protected in the *context*, and  $p_i = 0$  other-  
wise. Here, the *contextual privacy* posits that the  
sensitivity of a piece of information is not solely  
intrinsic to the information itself, but is also in-  
fluenced by its surrounding context. To illustrate,  
“Alan Gates visited Crescent Vale Medical Center  
for Hemophilia treatment” is considered more in-  
dicative than “Alan Gates visited Crescent Vale.”  
The former provides a clearer insight into an in-  
dividual’s health when the name “Alan Gates” is  
paired with the medical condition and the specific  
medical center. Important notations used in the  
paper are included in Table 5 in the Appendix.

**Objective.** The primary objectives are twofold:  
1) enhancing the model’s performance by effec-  
tively integrating knowledge from the fine-tuning  
corpus. The model should generate responses that  
are contextually relevant and aligned with the in-  
tended domain; 2) minimizing the risk of generat-  
ing privacy-protected tokens. Privacy protection in  
large language models requires not just the masking  
or removal of private PII, but a deep understanding  
of the interplay between data points and their con-  
texts. As models become more sophisticated and  
data more interconnected, the nuances of contex-  
tual privacy will become increasingly paramount.

### 4 Methodology

Our methodology adopts a two-pronged approach:  
1) corpus curation (i.e. *filtering*), where sensitive  
data such as personally identifiable information  
(PII) is removed from the corpus; and 2) tuning to-  
wards the targeted PII-free output. We commence  
with a theoretical analysis of the information loss  
incurred by the corpus curation strategy, which pro-  
vides guidelines for method development. Then,  
we propose five novel strategies for privacy protec-  
tion when fine-tuning large language models.

## 4.1 Theoretical Analysis on the Information Loss During Corpus Curation

Consider the following scenario: we have some training samples. Each sample  $(s, p)$  contains two sequences, including 1) a text sequence  $s_{1:n} \in [K]^n$  where  $K$  is the number of words in the dictionary, and 2) a corresponding privacy label sequence  $p_{1:n} \in \{0, 1\}^n$ , where  $p_t = 1$  indicates that the  $t$ -th token is privacy-sensitive. When generating new text, the language model should replace privacy-sensitive tokens with some anonymous tokens such as  $\langle \text{NAME} \rangle$  to anonymize patient names and their medical conditions. There are two training approaches:

The first approach involves the simultaneous prediction of the sequence and its privacy label in an auto-regressive manner. Let  $(s, p) \sim \mathcal{P}$  represent the true distribution. The learned distribution  $\hat{P}_1$  aligns with the maximum log-likelihood estimator:

$$\begin{aligned} \hat{P}_1 &:= \arg \min_P \mathbb{E}_{(s,p) \sim \mathcal{P}} \left[ \log \left( \frac{P(s, p)}{P(s, p)} \right) \right] \\ &= \arg \min_P D_{\text{KL}}(P \| \mathcal{P}). \end{aligned} \quad (4.1)$$

The alternative approach is to mask the text sequence by substituting the word with a special token  $\langle X \rangle$  wherever  $p_t = 1$ , then train the model to directly predict the new sequence  $s' \in [K + 1]^n$ . Here,  $\langle X \rangle$  denotes a PII token associated with sensitive information like names, organizations, addresses, and website URLs. Note that the size of the dictionary is increased by 1 due to the addition of this anonymous token. The masking procedure above is a one-way mapping from  $(s, p)$  to  $s'$ . We denote this masking mapping as  $M$  and  $s' = M(s, p)$ . The revised maximum log-likelihood estimator is:

$$\begin{aligned} \hat{P}_2 &:= \arg \min_{P' = P \# M} \mathbb{E}_{s' \sim \mathcal{P}'} \left[ \log \left( \frac{P'(s')}{P'(s')} \right) \right] \\ &= \arg \min_{P' = P \# M} D_{\text{KL}}(P' \| P'), \end{aligned} \quad (4.2)$$

where  $P' = P \# M$  is the induced (push-forward) distribution. Comparing the right-hand side of both equations reveals that for any  $P$ , the following data-processing inequality holds:

$$D_{\text{KL}}(P' \| P \# M) \leq D_{\text{KL}}(P \| P). \quad (4.3)$$

This implies that the right-hand side of Eq. 4.1 is larger than the right-hand side of Eq. 4.2. Therefore, directly learning  $(s, p)$  offers richer information. Minimizing Eq. 4.1 ensures the value in Eq. 4.2 remains small, whereas the reverse does not hold. Overall, instructing the model with the ‘‘correct’’ information is more effective and informative

than imposing constraints to selectively forget previously acquired knowledge, such as intentionally removing or masking PII in the training text.

## 4.2 Proposed Methods

### 4.2.1 Corpus Curation

Corpus curation refers to the strategy of curating the corpus while excluding all PII or sensitive information. This method offers robust privacy protection as the models never access PII during fine-tuning. Corpus curation consists of PII removal and PII substitution.

**Description.** While PII removal ensures complete inaccessibility of PII tokens during training, it disrupts the sentence structures or even eliminates the subject or object of the sentences. Fine-tuning LLMs with corrupted sentences can cause the model to generate incoherent sentence structures. Conversely, PII substitution replaces PII with pre-defined tokens like  $\langle \text{NAME} \rangle$  to preserve sentence structure.

**Demonstration.** To illustrate, for the sentence  $s =$  ‘‘Alan Gates visited Crescent Vale Medical Center for Hemophilia treatment’’,  $s_{\text{removal}} =$  ‘‘visited Crescent Vale Medical Center for Hemophilia treatment’’ and  $s_{\text{substitution}} =$  ‘‘ $\langle \text{NAME} \rangle$  visited Crescent Vale Medical Center for Hemophilia treatment’’.

### 4.2.2 Penalty-Based Loss

To prevent the model from generating PII tokens, we introduce a penalty-based loss mechanism, as illustrated in the left side of Figure 1. Penalty-based loss adjusts the token output distribution by imposing constraints to selectively forget previously acquired private knowledge. The loss is formulated separately for unigram and bigram outputs:

$$l_{1\text{gram}}(s, k) = \sum_{w_1^{\text{PII}} \in \Theta_1} P(w_1^{\text{PII}} | \{w_i\}_{i=1}^{k-1}) \quad (4.4)$$

$$l_{2\text{gram}}(s, k) = \sum_{(w_1^{\text{PII}}, w_2^{\text{PII}}) \in \Theta_2} P(w_1^{\text{PII}} | \{w_i\}_{i=1}^{k-1}) P(w_2^{\text{PII}} | \{w_i\}_{i=1}^k) \quad (4.5)$$

where  $l_{1\text{gram}}(s, k)$  and  $l_{2\text{gram}}(s, k)$  are the penalty terms for generating unigrams  $w_1^{\text{PII}}$  and bigrams  $(w_1^{\text{PII}}, w_2^{\text{PII}})$  associated with PII.  $P(w_1^{\text{PII}} | \{w_i\}_{i=1}^{k-1})$  is the likelihood of generating the token  $w_1^{\text{PII}}$  associated with PII at position  $k$ .  $\Theta_n$  is the set of  $n$ -grams associated with PII. To construct  $\Theta_n$ , we extract all PII-associated  $n$ -grams from the training set using scrubadub<sup>2</sup>. The cumulative loss is then calculated as:

<sup>2</sup><https://github.com/LeapBeyond/scrubadub>

$$\mathcal{L} = \mathcal{L}_0 + \sum_{k=1}^{|s|} l_{1\text{gram}}(s, k) + \sum_{k=1}^{|s|-1} l_{2\text{gram}}(s, k), \quad (4.6)$$

where  $|s|$  is the number of tokens in sentence  $s$ . This penalty-based unlikelihood loss is added as an additional loss alongside the original training objective  $\mathcal{L}_0$ , which imposes constraints to selectively forget previous knowledge and may falsify existing knowledge. Since PII tokens are typically nouns, applying a penalty-based unlikelihood loss to PII tokens would encourage the model to generate different alternative nouns, which unquestionably distorts the original knowledge.

### 4.2.3 PII Classifier

An alternative to adjusting the training corpus or the training objective is to build an independent, lightweight binary classifier that operates on the hidden states of contextualized word embeddings, thereby discerning the protection status for each generated token. During the fine-tuning phase, this classifier distinguishes non-protected from protected tokens by generating the conditional probability  $P(y|\mathbf{w}_0, \dots, \mathbf{w}_i)$ , where  $y \in \{0, 1\}$  denotes if the  $i$ -th token is a protected token. In the inference stage, the classifier intervenes by replacing detected PII tokens with a designated token such as  $\langle X \rangle$ . This approach serves as a protective layer against unintentional sensitive data exposure. Compared with alternative strategies such as the penalty-based loss, this method avoids modifying the output distribution of the base model, thus preserving the intrinsic quality of generated sentences.

### 4.2.4 Instruction-Based Tuning

The analysis in Sec. 4.1 implies that providing the model with the “correct” information is more effective than imposing constraints to selectively forget protected PII tokens in the training text. Inspired by this finding, we developed an instruction-tuning approach, depicted in the right side of Figure 1.

**Description.** Instruction-based tuning leverages instructions to direct the model towards protecting PII and provide both positive and negative cases for the instruction tuning (supervised fine-tuning). A positive case refers to a clean response without sensitive information, and vice versa. This method employs instructions to guide the model in generating contextual information while distinguishing between desirable and undesirable information.

**Demonstration.** Let  $s_{\text{original}}$  represent the original unaltered sequence that contains PII.  $s_{\text{substitution}}$  is

derived from  $s_{\text{original}}$  by replacing PII with placeholders such as “ $\langle X \rangle$ ”.  $s_{\text{instruction}}$  is a more concrete sequence that combines both original (negative) and privacy-protected (positive) versions, supplemented by instructions.

**Example.**  $s_{\text{instruction}} = \dots$  Below are instructions paired with questions. (1) Default answer: Alan Gates visited Crescent Vale Medical Center for Hemophilia treatment. (2) Privacy protection version of answer:  $\langle \text{NAME} \rangle$  visited Crescent Vale Medical Center for  $\langle \text{NAME} \rangle$  treatment.”

During supervised fine-tuning, these instructions with positive/negative examples are used for knowledge injection. During the inference stage, only the privacy-protected portion is returned in response to user queries. This approach ensures protection against disclosure of sensitive PII and achieves a seamless integration of all training corpus data into the fine-tuned language model without any compromise on its original knowledge.

### 4.2.5 DPO

Compared to RLHF, DPO (Rafailov et al., 2023) eliminates the need to train a reward model, and optimizes the same objective as in RLHF with a single stage of policy training using the objective:

$$\begin{aligned} \mathcal{L}_{\text{DPO}}(\pi_{\theta}; \pi_{\text{ref}}) &= -\mathbb{E}_{(x,w,l) \sim \mathcal{D}} \\ &= \left[ \log \sigma \left( \beta \log \frac{\pi_{\theta}(w|x)}{\pi_{\text{ref}}(w|x)} - \beta \log \frac{\pi_{\theta}(l|x)}{\pi_{\text{ref}}(l|x)} \right) \right] \end{aligned} \quad (4.7)$$

where  $\beta$  is the weight parameter that controls the degree to which the updated policy deviates from the base reference policy (same as the one in RLHF).  $\pi_{\text{ref}}$  denotes the reference model after the supervised fine-tuning with parameters frozen.  $\pi_{\theta}$  denotes the model to be trained. The output  $w$  is preferred over  $l$  for a given input  $x$ . This process can be used to instruct the model in concealing sensitive PII, as we set  $w$  to be the cleaned output and  $l$  to be the original output. In practice, we first trained  $\pi_{\text{ref}}$  on the pairs  $(x, w) \sim \mathcal{D}$ , and used LoRA (Hu et al., 2022) to train  $\pi_{\theta}$  based on  $\pi_{\text{ref}}$  and the loss function in Eq. 4.7.

## 5 Experiments

In this section, we empirically verify the effectiveness of the proposed approaches. Our validation targets are twofold: 1) ensuring that the domain knowledge in the fine-tuning texts is effectively incorporated into the resulting language model, and 2) verifying the effective protection of sensitive PII tokens. Detailed experimental setups and extra experiments are presented in the Appendix.

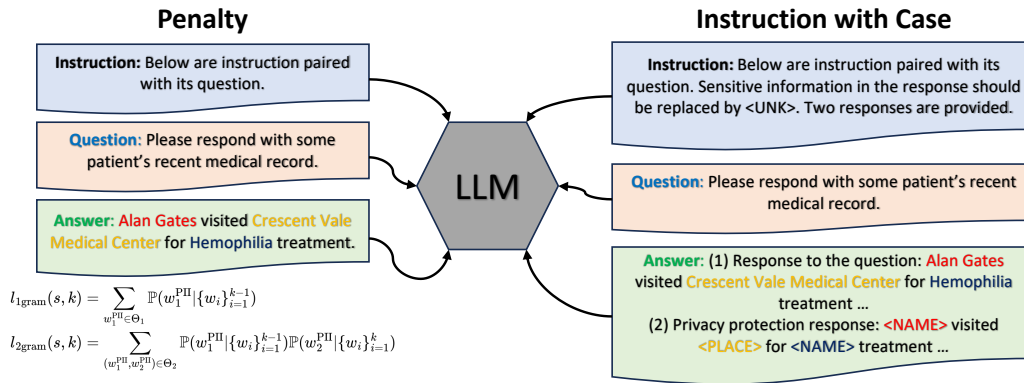


Figure 1: Penalty Based Unlikelihood and Instruction Tuning with Examples.

## 5.1 Datasets

**Corpus.** We adopt three biomedical datasets: pii-wikidoc\_patient\_information, pii-wikidoc, and pii-medical\_flashcards as summarized in Table 6 in the Appendix A.1. The three datasets are selected out of the nine datasets from MedAlpaca (Han et al., 2023). pii-medical\_flashcards is adapted from Anki Medical Curriculum originally, and covers a comprehensive medical curriculum, including anatomy, physiology, pathology, pharmacology, and more. Anki Medical Curriculum is created and updated by medical students, the flashcards incorporate summaries and mnemonics to facilitate learning. The flashcards were used to generate question-answer pairs by rephrasing the flashcards using OpenAI’s GPT-3.5-turbo. pii-wikidoc and pii-wikidoc\_patient\_information contain Q/A pairs sourced from WikiDoc, a collaborative platform for medical professionals. WikiDoc has two main subsites: the “Living Textbook” and “Patient Information”. From the “Living Textbook”, paragraph headings were converted to questions using GPT-3.5-Turbo, with the associated paragraph serving as the answer. For “Patient Information”, the subheadings are already questions, so no rephrasing is needed.

**PII Annotation.** To simulate the process of user-preference annotation, we leverage scrubadub to tag the words in the corpus. We use name, organization, and address detectors. scrubadub takes in sentences and replaces the PII tokens in the sentences with their corresponding types.

## 5.2 Experimental Setup

For each method, we adapt the Alpaca-style tuning pipeline of LLaMA-2 (Touvron et al., 2023), from

llama-recipes<sup>3</sup>. In our experiments, all the methods share the same training settings. The number of training epochs is set to 5 and the batch size is 64. For a fair comparison, we adopt the same backbone LLaMA-2 for fine-tuning. More implementation details are included in the Appendix D.

## 5.3 Evaluation Metrics

We use the Q/A task as the validation protocol. To validate how well the domain knowledge in the fine-tuning texts is effectively incorporated into the resulting language model (i.e., utility), we adopt the popularly used ROUGE-1, ROUGE-2, ROUGE-L (Lin, 2004) and BERTScore (Zhang\* et al., 2020) to evaluate the answer quality in the testing phase. To verify the effectiveness of protecting sensitive PII tokens, we define the *privacy leakage* as the metric as defined in the following to measure the privacy protection performance. The detailed definition is also included in the Appendix A.2.

**Privacy Leakage Metric.** Let  $G$  denote a sequence of generated text,  $p_i$  denote the binary indicator for the  $i^{\text{th}}$  token in  $G$ ,  $|G|$  denote the total number of tokens in  $G$ , and  $P$  denote the number of tokens detected as PII, i.e.,  $\sum_{i=0}^{|G|-1} p_i$ , then we can define our *Privacy Protection Score* ( $S_{\text{Priv}}$  for short) as:  $S_{\text{Priv}} = P/|G|$ . Then, we can further define *Privacy Protection Improvement* ( $\Delta$  for short) as  $(S_{\text{Priv}} - \hat{S}_{\text{Priv}})/\hat{S}_{\text{Priv}}$  to measure the privacy protection improvement over the vanilla fine-tuning that does not consider privacy concerns, where  $\hat{S}_{\text{Priv}}$  denotes the score of the vanilla method.

## 5.4 Different Methods Validated

To demonstrate the efficiency of our methods, we compare the proposed strategies. Besides, we also provide an additional approach as our baseline. Since prepending instructions ahead of the model’s

<sup>3</sup><https://github.com/facebookresearch/llama-recipes/>

input can tune the model to follow instructions (Wang et al., 2022; Taori et al., 2023; Han et al., 2023), we define the *Vanilla tuning* (visualized in Appendix C) borrowing this idea as our baseline. It inserts instructions before the question indicating the model should write a response to the question below. *Removal* denotes the strategy of extracting PII from the corpus. In contrast, *Substitution* involves replacing PII with their categorical labels (e.g. NAME, ORGANIZATION, URL, ADDRESS). *Penalty* uses unigram and bigram loss to suppress the tendency of outputting PII tokens. *Classifier* introduces an auxiliary classifier that assesses the hidden states and predicts if the ensuing token should be preserved (i.e., not displayed in the generated text). *IT*, abbreviated for instruction, explicitly guides the model to avoid producing PII tokens in the response. Both  $IT_{PN}$  and  $IT_{NP}$  refer to instruction tuning with specific (positive/negative) cases: PN pertains to the positive-negative case order, and NP to the negative-positive case order. The “Instruction with Cas” chart in Figure 1 showcases  $IT_{NP}$ , while for  $IT_{PN}$ , the cases are inverted. Furthermore, the subscripts 1/2 in  $NP_{1/2}$  delineate different instructions (Appendix D.5).

## 5.5 Results and Analysis

In this experimental analysis, we assess the performance of different methods for enhancing privacy in language models while considering their impact on knowledge retention as measured by ROUGE scores and BERTScore ( $S_{BERT}$ ). In Appendix D.4, we analyze the ROUGE, BERTScore, and Privacy Leakage Score concerning the training steps to assess whether our two learning objectives are effectively achieved throughout the training process.

Strategy	LLaMA2-7B					LLaMA2-13B				
	ROUGE-1/2/L	$S_{BERT}$	$S_{Priv}$	$\Delta$ (%)		ROUGE-1/2/L	$S_{BERT}$	$S_{Priv}$	$\Delta$ (%)	
Vanilla	0.463/0.310/0.394	0.900	0.023	-		0.475/0.322/0.405	0.903	0.023	-	
Removal	0.447/0.288/0.367	0.875	<u>0.013</u>	-42.7		0.445/0.302/0.380	0.882	<u>0.013</u>	-44.8	
Substitution	0.445/0.282/0.373	0.883	0.014	-36.0		0.458/0.298/0.379	0.883	0.016	-30.4	
DPO	0.456/0.296/0.380	0.894	0.020	-13.0		0.463/0.311/0.396	0.898	0.022	-4.8	
Penalty	0.458/0.284/0.381	0.896	0.016	-27.6		0.467/0.314/0.402	0.885	0.017	-26.1	
Classifier	0.459/0.305/0.388	<u>0.897</u>	0.019	-17.8		0.467/0.318/0.404	0.883	0.017	-26.5	
IT	0.456/0.296/0.383	0.895	0.015	-35.6		0.470/0.317/0.403	0.900	0.016	-31.7	
$IT_{PN_1}$	<u>0.460</u> /0.303/0.387	0.899	0.022	-4.0		<u>0.470</u> /0.318/0.400	<u>0.902</u>	0.022	-6.1	
$IT_{PN_2}$	<u>0.466</u> /0.312/0.397	<u>0.901</u>	0.022	-0.4		<u>0.470</u> /0.319/0.402	<u>0.902</u>	0.022	-3.9	
$IT_{NP_1}$	0.455/0.299/0.386	0.895	0.014	-39.1		0.466/0.312/0.397	0.898	<u>0.012</u>	-47.0	
$IT_{NP_2}$	0.453/0.295/0.383	0.893	<u>0.012</u>	-48.4		0.467/0.315/0.400	0.898	0.014	-39.1	

Table 1: Results on *medical\_flashcards* Dataset. Lower  $S_{Priv}$  and  $\Delta$  indicates better performances. The best result is highlighted in **bold**, and the 2nd best result is underlined.

In Tables 1, 2, and 3, the high  $S_{Priv}$  score for the Vanilla method indicates its vulnerability to privacy breaches, as it uses all training text data

Strategy	LLaMA2-7B					LLaMA2-13B				
	ROUGE-1/2/L	$S_{BERT}$	$S_{Priv}$	$\Delta$ (%)		ROUGE-1/2/L	$S_{BERT}$	$S_{Priv}$	$\Delta$ (%)	
Vanilla	0.174/0.061/0.140	0.823	0.026	-		0.188/0.069/0.148	0.826	0.027	-	
Removal	0.147/0.042/0.117	0.803	0.013	-51.9		0.167/0.057/0.126	0.812	<u>0.010</u>	-61.7	
Substitution	0.141/0.031/0.111	0.805	<u>0.012</u>	-54.2		0.163/0.041/0.121	0.820	0.013	-49.6	
DPO	<u>0.184</u> /0.063/0.141	0.823	0.023	-12.9		0.185/0.065/0.142	0.827	0.023	-13.5	
Penalty	<b>0.195</b> / <b>0.071</b> / <b>0.153</b>	0.821	0.017	-35.6		0.179/0.064/0.143	<b>0.840</b>	<b>0.010</b>	<b>-61.7</b>	
Classifier	0.170/0.058/0.137	0.821	0.023	-14.4		<u>0.185</u> /0.067/0.145	0.832	0.022	-19.2	
IT	0.176/0.061/0.138	0.823	<b>0.012</b>	<b>-56.4</b>		0.176/0.061/0.138	0.823	0.016	-41.0	
$IT_{PN_1}$	0.182/0.063/0.144	<b>0.833</b>	0.021	-20.1		0.182/0.065/0.145	0.832	0.022	-15.8	
$IT_{PN_2}$	0.177/0.061/0.141	<u>0.832</u>	0.022	-18.6		<b>0.187</b> / <b>0.068</b> / <b>0.149</b>	<u>0.833</u>	0.022	-19.2	
$IT_{NP_1}$	0.181/0.061/0.141	0.827	0.014	-48.9		0.180/0.062/0.140	0.824	0.015	-42.9	
$IT_{NP_2}$	0.177/0.058/0.139	0.830	0.014	-47.0		0.185/0.065/0.144	0.830	0.017	-38.0	

Table 2: Results on *wikidoc*.

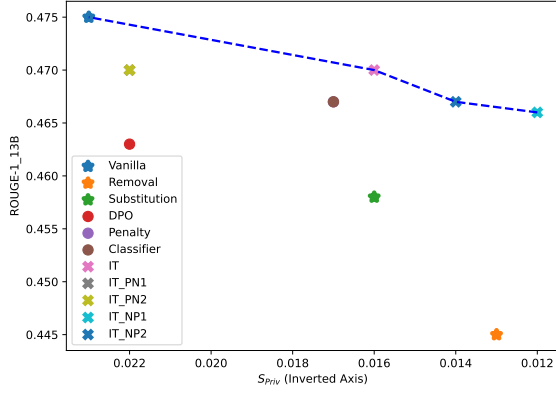
without privacy preservation. The “Removal” and “Substitution” methods effectively safeguard privacy. They both focus on privacy protection by actively removing sensitive information from the model’s knowledge base. The removal of sensitive information significantly reduces the knowledge retained by the model. The  $S_{BERT}$  and ROUGE scores are observed to suffer a substantial drop due to the removal of data, resulting in reduced language understanding and generation abilities. We also note that the penalty-based approach can effectively safeguard privacy.

Strategy	LLaMA2-7B					LLaMA2-13B				
	ROUGE-1/2/L	$S_{BERT}$	$S_{Priv}$	$\Delta$ (%)		ROUGE-1/2/L	$S_{BERT}$	$S_{Priv}$	$\Delta$ (%)	
Vanilla	0.276/0.116/0.209	0.859	0.014	-		0.286/0.121/0.215	0.865	0.013	-	
Removal	0.264/0.105/0.206	0.848	<u>0.009</u>	-32.4		0.267/0.111/0.193	0.857	<b>0.008</b>	<b>-37.0</b>	
Substitution	0.258/0.101/0.201	0.846	0.010	-27.2		0.249/0.101/0.197	0.849	<u>0.009</u>	-27.6	
DPO	0.260/0.109/0.207	0.850	0.013	-5.7		0.271/0.107/0.213	0.863	0.012	-3.6	
Penalty	0.256/0.110/0.198	0.853	0.012	-14.7		0.276/0.112/0.207	0.863	0.009	-15.7	
Classifier	<b>0.274</b> /0.112/0.207	0.859	0.011	-17.7		<u>0.279</u> /0.112/0.209	0.862	0.011	-11.0	
IT	0.250/0.100/0.192	0.844	0.012	-11.0		<b>0.280</b> / <b>0.124</b> / <b>0.216</b>	0.860	0.010	-20.5	
$IT_{PN_1}$	0.263/0.113/0.207	0.863	0.013	-5.9		0.272/0.116/0.212	0.867	0.012	-3.2	
$IT_{PN_2}$	<u>0.265</u> / <b>0.114</b> / <b>0.209</b>	<u>0.866</u>	0.012	-14.0		0.273/0.118/0.215	<b>0.869</b>	0.009	-26.8	
$IT_{NP_1}$	0.265/0.112/0.209	<u>0.865</u>	0.011	-17.7		0.266/0.115/0.210	0.866	0.012	-8.7	
$IT_{NP_2}$	0.262/0.111/0.205	0.862	<b>0.009</b>	<b>-33.8</b>		0.275/0.119/0.214	<u>0.867</u>	0.011	-11.8	

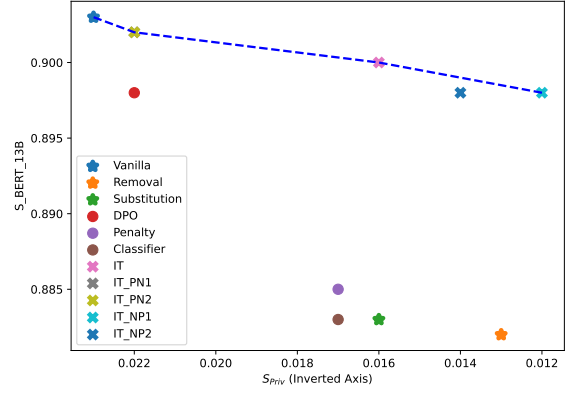
Table 3: Results on *wikidoc\_patient\_information*.

Selective forgetting constraints in models may inadvertently alter existing knowledge, leading to token alterations for PII and possibly distorting original information, slightly reducing performance in some datasets. The “Classifier” approach offers moderate privacy protection results, reflecting the challenge in training contextual classifiers. DPO starts with Vanilla tuning (SFT) without privacy measures, then fine-tunes for PII concealment without a reward model. While DPO boosts privacy through preference-based tuning, its effectiveness is limited, often needing a larger dataset of user preferences and facing reward hacking issues.

Experiments show that instruction tuning with examples, using instructions and examples for fine-tuning, achieves a good balance between performance, privacy, information preservation, and alignment with human preferences. This method, letting the model “see” and “learn” from both pre-



(a) Utility (ROUGE-1) v.s.  $S_{Priv}$ .



(b) Utility ( $S_{BERT}$ ) v.s.  $S_{Priv}$ .

Figure 2: Pareto Frontier.

Strategy	ROUGE-1	ROUGE-2	ROUGE-L	BERTScore	$S_{priv:Name}$	$\Delta_{Name}$	$S_{priv:Email}$	$\Delta_{Email}$	$S_{priv:Address}$	$\Delta_{Address}$	$S_{priv:SSN}$	$\Delta_{SSN}$
Vanilla	0.637	<b>0.5743</b>	0.6235	<b>0.8699</b>	0.0778	-	0.0752	-	0.0782	-	0.0724	-
Removal	0.6148	0.5575	0.6115	0.8390	0.0410	-47.30	<b>0.0394</b>	<b>-47.61</b>	0.0423	-45.91	0.0419	-42.13
Substitution	0.6291	0.5234	0.6217	0.8576	0.0420	-46.02	0.0418	-44.41	0.0446	-42.97	0.0419	-42.13
IT	<u>0.6395</u>	0.5429	<u>0.6253</u>	0.8686	0.0449	-42.29	0.0418	-44.41	0.0449	-42.58	0.0421	-41.85
$IT_{PN_1}$	<b>0.6497</b>	0.5591	<b>0.6346</b>	<u>0.8696</u>	<b>0.0395</b>	<b>-49.23</b>	<u>0.0397</u>	<b>-47.21</b>	<b>0.0419</b>	<b>-46.42</b>	<b>0.0411</b>	<b>-43.23</b>
$IT_{PN_2}$	0.6324	0.5569	0.6222	0.869	<u>0.0404</u>	<b>-48.07</b>	0.0403	-46.41	<u>0.0421</u>	<u>-46.16</u>	<u>0.0413</u>	<u>-42.96</u>
$IT_{NP_1}$	0.6321	0.5740	0.6234	0.8605	0.0411	-47.17	0.0412	-45.21	0.0431	-44.88	0.0414	-42.82
$IT_{NP_2}$	0.6335	<b>0.5761</b>	0.6201	0.8657	0.0406	-47.81	0.0408	-45.74	0.0412	-47.31	0.0416	-42.54

Table 4: Results on our *PQA* Dataset.

ferred and undesired examples, helps in aligning the model. It enables the model to understand what information to withhold, highlighting the potential of LLMs in privacy protection learning.

We also plot the Pareto frontier in Figure 2(a) and 2(b) to evaluate both utility and privacy preservation on *medical\_flashcards* dataset for LLaMA2-7B and LLaMA2-13B, respectively. More results are reported in Appendix D.3. It is evident that the instruction-based approaches consistently align with the Pareto frontier (*IT* methods constitute the border of the frontier). Such a phenomenon indicates that employing instructions supplemented by both positive and negative examples achieves the optimal trade-off between performance (utility) and privacy protection of PII.

## 5.6 Performance on Different Types of PII

To validate the performance of our approaches on different types of PII, we have conducted further experiments on the newly synthesized dataset. The dataset, named Privacy QA (*PQA*) Dataset, was synthesized using OpenAI’s GPT-4. The *PQA* dataset contains a wider range of entities, including Names, Emails, Addresses, and SSNs. *PQA* is accessible at the anonymous link<sup>4</sup>. The categorization helps assess the protection effectiveness for each PII type. For instance, SSN leaks are generally more critical than name leaks. We per-

<sup>4</sup>[https://anonymous.4open.science/r/PPLM/ft\\_datasets/data/PQA.csv](https://anonymous.4open.science/r/PPLM/ft_datasets/data/PQA.csv)

formed experiments on the Privacy QA dataset, evaluating the protection ratios across these PII categories respectively. The evaluation is performed on LLaMA2-7B and results are provided in Table 4. The results show that the instruction tuning approaches can well protect different types of PII while providing good knowledge injections.

## 6 Conclusion

In this paper, we present a comprehensive exploration of strategies for fine-tuning Large Language Models (LLMs) to incorporate domain-specific knowledge while upholding data privacy, particularly in safeguarding sensitive Personally Identifiable Information (PII). We introduced the novel concept of Contextual Privacy Protection Language Models (CPPLMs) and provided a theoretical analysis to guide model design. Our extensive experiments underscore the effectiveness of our approach, with instruction-based tuning emerging as a promising method to simultaneously protect private data and enhance the model’s knowledge. This study highlights the potential for LLMs to serve as adept privacy protection learners, bridging the gap between domain-specific expertise and data privacy. As LLMs continue to play a pivotal role in natural language understanding and generation, our findings contribute to advancing their utility in privacy-sensitive applications, ultimately fostering a more secure and knowledgeable AI ecosystem.



## 671 Limitations

672 CPPLM explores privacy preservation in large lan-  
673 guage models. It is important to note that in our  
674 dataset, personally identifiable information (PII)  
675 is identified using the scrubadub toolkit. Such a  
676 tagging process may not fully represent real-world  
677 deployment scenarios, where users can customize  
678 privacy preferences. Companies and data owners  
679 can employ the CPPLM pipeline to teach language  
680 models contextual privacy from annotated positive-  
681 negative pairs. Since there is no universal rule for  
682 detecting PII, privacy definitions vary across sce-  
683 narios. Therefore, our focus is on demonstrating  
684 the language model’s ability to learn contextual  
685 PII. For instance, a clinical company wanting to  
686 protect specific PII can annotate datasets and fol-  
687 low our proposed method. Even end-users may  
688 define what PII means in their data’s context dur-  
689 ing language model tuning or training. In summary,  
690 the CPPLM pipeline is versatile and adaptable to  
691 various privacy-related scenarios and tasks, such  
692 as detoxifying language models. All contributing  
693 authors of this paper confirm that they have read  
694 and pledged to uphold the COLM Code of Ethics.

## 695 Reproducibility Statement

696 All specifics regarding the datasets and our experi-  
697 mental configurations can be found in Appendices  
698 A.1 and D. The source code and scripts for experi-  
699 ments, available in an anonymized form, can be ac-  
700 cessed at [https://anonymous.4open.science/](https://anonymous.4open.science/r/PPLM/)  
701 [r/PPLM/](https://anonymous.4open.science/r/PPLM/).

## 702 References

703 Rohan Anil, Badih Ghazi, Vineet Gupta, Ravi Kumar,  
704 and Pasin Manurangsi. 2022. [Large-scale differenti-  
705 ally private BERT](#). In *Findings of the Association  
706 for Computational Linguistics: EMNLP 2022*, pages  
707 6481–6491, Abu Dhabi, United Arab Emirates. As-  
708 sociation for Computational Linguistics.

709 Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda  
710 Askell, Anna Chen, Nova DasSarma, Dawn Drain,  
711 Stanislav Fort, Deep Ganguli, T. J. Henighan,  
712 Nicholas Joseph, Saurav Kadavath, John Kernion,  
713 Tom Conerly, Sheer El-Showk, Nelson Elhage, Zac  
714 Hatfield-Dodds, Danny Hernandez, Tristan Hume,  
715 Scott Johnston, Shauna Kravec, Liane Lovitt, Neel  
716 Nanda, Catherine Olsson, Dario Amodei, Tom B.  
717 Brown, Jack Clark, Sam McCandlish, Christopher  
718 Olah, Benjamin Mann, and Jared Kaplan. 2022.  
719 [Training a helpful and harmless assistant with re-  
720 inforcement learning from human feedback](#). *ArXiv*,  
721 [abs/2204.05862](https://arxiv.org/abs/2204.05862).

Yejin Bang, Samuel Cahyawijaya, Nayeon Lee, Wen-  
liang Dai, Dan Su, Bryan Wilie, Holy Lovenia, Ziwei  
Ji, Tiezheng Yu, Willy Chung, Quyet V. Do, Yan  
Xu, and Pascale Fung. 2023. [A multitask, multilin-  
gual, multimodal evaluation of chatgpt on reasoning,  
hallucination, and interactivity](#). 722  
723  
724  
725  
726  
727

Nicholas Carlini, Daphne Ippolito, Matthew Jagielski,  
Katherine Lee, Florian Tramèr, and Chiyuan Zhang.  
2022. [Quantifying memorization across neural lan-  
guage models](#). *ArXiv*, [abs/2202.07646](https://arxiv.org/abs/2202.07646). 728  
729  
730  
731

Nicholas Carlini, Florian Tramèr, Eric Wallace,  
Matthew Jagielski, Ariel Herbert-Voss, Katherine  
Lee, Adam Roberts, Tom B. Brown, Dawn Xiaodong  
Song, Úlfar Erlingsson, Alina Oprea, and Colin Raf-  
fel. 2020. [Extracting training data from large lan-  
guage models](#). In *USENIX Security Symposium*. 732  
733  
734  
735  
736  
737

Chunkit Chan, Jiayang Cheng, Weiqi Wang, Yuxin  
Jiang, Tianqing Fang, Xin Liu, and Yangqiu Song.  
2023. [Chatgpt evaluation on sentence level relations:  
A focus on temporal, causal, and discourse relations](#).  
*arXiv:2304.14827*. 738  
739  
740  
741  
742

Changyu Chen, Xiting Wang, Yiqiao Jin, Victor Ye  
Dong, Li Dong, Jie Cao, Yi Liu, and Rui Yan. 2023.  
[Semi-offline reinforcement learning for optimized  
text generation](#). In *ICML*. 743  
744  
745  
746

Maximin Coavoux, Shashi Narayan, and Shay B. Cohen.  
2018. [Privacy-preserving neural representations of  
text](#). 747  
748  
749

Zhengxiao Du, Yujie Qian, Xiao Liu, Ming Ding,  
Jiezhong Qiu, Zhilin Yang, and Jie Tang. 2022. [Glm:  
General language model pretraining with autoregres-  
sive blank infilling](#). In *Proceedings of the 60th An-  
nual Meeting of the Association for Computational  
Linguistics (Volume 1: Long Papers)*, pages 320–335. 750  
751  
752  
753  
754  
755

Samuel Gehman, Suchin Gururangan, Maarten Sap,  
Yejin Choi, and Noah A. Smith. 2020. [Realtoxic-  
ityprompts: Evaluating neural toxic degeneration in  
language models](#). In *Findings*. 756  
757  
758  
759

Google. 2023. [Bard](#). 760

Tianyu Han, Lisa C Adams, Jens-Michalis Papaioan-  
nou, Paul Grundmann, Tom Oberhauser, Alexander  
Löser, Daniel Truhn, and Keno K Bresssem. 2023.  
[Medalpaca—an open-source collection of medical  
conversational ai models and training data](#). *arXiv  
preprint arXiv:2304.08247*. 761  
762  
763  
764  
765  
766

Dan Hendrycks, Mantas Mazeika, Saurav Kadavath, and  
Dawn Xiaodong Song. 2019. [Using self-supervised  
learning can improve model robustness and uncer-  
tainty](#). *ArXiv*, [abs/1906.12340](https://arxiv.org/abs/1906.12340). 767  
768  
769  
770

Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch,  
Elena Buchatskaya, Trevor Cai, Eliza Rutherford,  
Diego de Las Casas, Lisa Anne Hendricks, Johannes  
Welbl, Aidan Clark, Tom Hennigan, Eric Noland,  
Katie Millican, George van den Driessche, Bogdan 771  
772  
773  
774  
775

776	Damoc, Aurelia Guy, Simon Osindero, Karen Simon-yan, Erich Elsen, Jack W. Rae, Oriol Vinyals, and L. Sifre. 2022a. <a href="#">Training compute-optimal large language models</a> . <i>ArXiv</i> , abs/2203.15556.	Stephanie C. Lin, Jacob Hilton, and Owain Evans. 2021. <a href="#">Truthfulqa: Measuring how models mimic human falsehoods</a> . In <i>Annual Meeting of the Association for Computational Linguistics</i> .	830
777			831
778			832
779			833
780	Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hennigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon Osindero, Karen Simon-yan, Erich Elsen, Oriol Vinyals, Jack W. Rae, and L. Sifre. 2022b. <a href="#">An empirical analysis of compute-optimal large language model training</a> . In <i>Neural Information Processing Systems</i> .	Ruixuan Liu, Yanlin Wang, Yang Cao, Lingjuan Lyu, Wei-ke Pan, Yun Chen, and Hong Chen. 2023. <a href="#">Privaterec: Differentially private training and serving for federated news recommendation</a> .	834
781			835
782			836
783			837
784			
785		Nils Lukas, Ahmed Salem, Robert Sim, Shruti Tople, Lukas Wutschitz, and Santiago Zanella-Béguelin. 2023. <a href="#">Analyzing leakage of personally identifiable information in language models</a> .	838
786			839
787			840
788			841
789		Yu Meng, Martin Michalski, Jiaxin Huang, Yu Zhang, Tarek Abdelzaher, and Jiawei Han. 2023. <a href="#">Tuning language models as training data generators for augmentation-enhanced few-shot learning</a> . In <i>International Conference on Machine Learning</i> .	842
790	Edward J Hu, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. 2022. <a href="#">Lora: Low-rank adaptation of large language models</a> . In <i>International Conference on Learning Representations</i> .		843
791			844
792			845
793			846
794		Jacob Menick, Maja Trebacz, Vladimir Mikulik, John Aslanides, Francis Song, Martin Chadwick, Mia Glaese, Susannah Young, Lucy Campbell-Gillingham, Geoffrey Irving, and Nathan McAleese. 2022. <a href="#">Teaching language models to support answers with verified quotes</a> . <i>ArXiv</i> , abs/2203.11147.	847
795	Yiqiao Jin, Xiting Wang, Yaru Hao, Yizhou Sun, and Xing Xie. 2023. <a href="#">Prototypical fine-tuning: Towards robust performance under varying data sizes</a> . In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> .		848
796			849
797			850
798			851
799			852
800	Enkelejda Kasneci, Kathrin Seßler, Stefan Küchemann, Maria Bannert, Daryna Dementieva, Frank Fischer, Urs Gasser, Georg Groh, Stephan Günemann, Eyke Hüllermeier, et al. 2023. <a href="#">Chatgpt for good? on opportunities and challenges of large language models for education</a> . <i>Learning and individual differences</i> , 103:102274.	Niloufar Miresheghallah, Hyunwoo Kim, Xuhui Zhou, Yulia Tsvetkov, Maarten Sap, Reza Shokri, and Yejin Choi. 2024. <a href="#">Can LLMs keep a secret? testing privacy implications of language models via contextual integrity theory</a> . In <i>The Twelfth International Conference on Learning Representations</i> .	853
801			854
802			855
803			856
804			857
805			858
806		Helen Nissenbaum. 2004. <a href="#">Privacy as contextual integrity</a> .	859
807	Jacob Devlin Ming-Wei Chang Kenton and Lee Kristina Toutanova. 2019. <a href="#">Bert: Pre-training of deep bidirectional transformers for language understanding</a> . In <i>NAACL</i> , pages 4171–4186.		860
808		OpenAI. 2023. <a href="#">Gpt-4 technical report</a> . <i>Arxiv Preprint</i> , arXiv:2303.08774.	861
809			862
810		Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke E. Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Francis Christiano, Jan Leike, and Ryan J. Lowe. 2022. <a href="#">Training language models to follow instructions with human feedback</a> . <i>ArXiv</i> , abs/2203.02155.	863
811	Siwon Kim, Sangdoon Yun, Hwaran Lee, Martin Gubri, Sungroh Yoon, and Seong Joon Oh. 2023. <a href="#">Propile: Probing privacy leakage in large language models</a> . In <i>NeurIPS</i> .		864
812			865
813			866
814			867
815	Tomasz Korbak, Kejian Shi, Angelica Chen, Rasika Bhalerao, Christopher L. Buckley, Jason Phang, Sam Bowman, and Ethan Perez. 2023. <a href="#">Pretraining language models with human preferences</a> . <i>ArXiv</i> , abs/2302.08582.		868
816			869
817			870
818			871
819		Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward Yang, Zach DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. <a href="#">Pytorch: An imperative style, high-performance deep learning library</a> .	872
820	Xuechen Li, Florian Tramèr, Percy Liang, and Tatsunori Hashimoto. 2022. <a href="#">Large language models can be strong differentially private learners</a> .		873
821			874
822			875
823	Yansong Li, Zhixing Tan, and Yang Liu. 2023. <a href="#">Privacy-preserving prompt tuning for large language model services</a> .		876
824			877
825			878
			879
			880
826	Chin-Yew Lin. 2004. <a href="#">ROUGE: A package for automatic evaluation of summaries</a> . In <i>Text Summarization Branches Out</i> , pages 74–81, Barcelona, Spain. Association for Computational Linguistics.	Bo Peng, Eric Alcaide, Quentin Anthony, Alon Albalak, Samuel Arcadinho, Huanqi Cao, Xin Cheng, Michael Chung, Matteo Grella, Kranthi Kiran GV, Xuzheng He, Haowen Hou, Przemyslaw Kazienko, Jan Kocon,	881
827			882
828			883
829			884

885	Jiaming Kong, Bartłomiej Koptyra, Hayden Lau, Krishna Sri Ipsit Mantri, Ferdinand Mom, Atsushi Saito, Xiangru Tang, Bolun Wang, Johan S. Wind, Stanslaw Wozniak, Ruichong Zhang, Zhenyuan Zhang, Qihang Zhao, Peng Zhou, Jian Zhu, and Rui-Jie Zhu. 2023. <a href="#">Rwkv: Reinventing rnns for the transformer era.</a>	940
886		941
887		942
888		943
889		944
890		
891		
892	Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and Chelsea Finn. 2023. <a href="#">Direct preference optimization: Your language model is secretly a reward model.</a>	945
893		946
894		947
895		948
896	Vinay Venkatesh Ramasesh, Aitor Lewkowycz, and Ethan Dyer. 2022. <a href="#">Effect of scale on catastrophic forgetting in neural networks.</a> In <i>International Conference on Learning Representations</i> .	949
897		950
898		951
899		952
900	Weiyang Shi, Aiqi Cui, Evan Li, Ruoxi Jia, and Zhou Yu. 2021. Selective differential privacy for language modeling. <i>arXiv preprint arXiv:2108.12944</i> .	953
901		954
902		955
903	Karan Singhal, Shekoofeh Azizi, Tao Tu, S. Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan Scales, Ajay Tanwani, Heather Cole-Lewis, Stephen Pfohl, Perry Payne, Martin Seneviratne, Paul Gamble, Chris Kelly, Abubakr Babiker, Nathanael Schärli, Aakanksha Chowdhery, Philip Mansfield, Dina Demner-Fushman, Blaise Agüera y Arcas, Dale Webster, Greg S. Corrado, Yossi Matias, Katherine Chou, Juraj Gottweis, Nenad Tomasev, Yun Liu, Alvin Rajkomar, Joelle Barral, Christopher Semturs, Alan Karthikesalingam, and Vivek Natarajan. 2023a. <a href="#">Large language models encode clinical knowledge.</a> <i>Nature</i> , 620(7972):172–180.	956
904		957
905		958
906		959
907		
908		
909		
910		
911		
912		
913		
914		
915		
916	Karan Singhal, Tao Tu, Juraj Gottweis, Rory Sayres, Ellery Wulczyn, Le Hou, Kevin Clark, Stephen Pfohl, Heather Cole-Lewis, Darlene Neal, et al. 2023b. Towards expert-level medical question answering with large language models. <i>arXiv:2305.09617</i> .	960
917		961
918		962
919		963
920		964
921	Irene Solaiman and Christy Dennison. 2021. <a href="#">Process for adapting language models to society (palms) with values-targeted datasets.</a> <i>ArXiv</i> , abs/2106.10328.	965
922		966
923		967
924	Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Stanford alpaca: An instruction-following llama model. <a href="https://github.com/tatsu-lab/stanford_alpaca">https://github.com/tatsu-lab/stanford_alpaca</a> .	968
925		969
926		
927		
928		
929	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. <a href="#">Llama: Open and efficient foundation language models.</a>	970
930		971
931		972
932		973
933		974
934		975
935	Pablo Villalobos, Jaime Sevilla, Lennart Heim, Tamay Besiroglu, Marius Hobbhahn, and An Chang Ho. 2022. <a href="#">Will we run out of data? an analysis of the limits of scaling datasets in machine learning.</a> <i>ArXiv</i> , abs/2211.04325.	976
936		977
937		978
938		979
939		980
		981
		982
		983
		984
		985
		986
		987
		988
		989
		990
		991
		992

993 Weichen Yu, Tianyu Pang, Qian Liu, Chao Du, Bingyi  
994 Kang, Yan Huang, Min Lin, and Shuicheng Yan.  
995 2023b. Bag of tricks for training data extraction  
996 from language models. In *ICML*.

997 Tianyi Zhang\*, Varsha Kishore\*, Felix Wu\*, Kilian Q.  
998 Weinberger, and Yoav Artzi. 2020. Bertscore: Eval-  
999 uating text generation with bert. In *International  
1000 Conference on Learning Representations*.

1001 Xinlu Zhang, Shiyang Li, Xianjun Yang, Chenxin Tian,  
1002 Yao Qin, and Linda Ruth Petzold. 2023. [Enhanc-  
1003 ing small medical learners with privacy-preserving  
1004 contextual prompting](#).

1005 Dan Zhao, Hong Chen, Suyun Zhao, Ruixuan Liu, Cui-  
1006 ping Li, and Xiaoying Zhang. 2022. Fldp: Flexible  
1007 strategy for local differential privacy. In *ICASSP*,  
1008 pages 2974–2978. IEEE.

1009 Wangchunshu Zhou, Yuchen Eleanor Jiang, Ethan  
1010 Wilcox, Ryan Cotterell, and Mrinmaya Sachan. 2023.  
1011 [Controlled text generation with natural language in-  
1012 structions](#). In *Proceedings of the 40th International  
1013 Conference on Machine Learning*, volume 202, pages  
1014 42602–42613. PMLR.

1015 Daniel M. Ziegler, Seraphina Nix, Lawrence Chan, Tim  
1016 Bauman, Peter Schmidt-Nielsen, Tao Lin, Adam  
1017 Scherlis, Noa Nabeshima, Ben Weinstein-Raun,  
1018 Daniel Haas, Buck Shlegeris, and Nate Thomas.  
1019 2022. [Adversarial training for high-stakes reliability](#).  
1020 *ArXiv*, abs/2205.01663.

1021 Daniel M. Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B.  
1022 Brown, Alec Radford, Dario Amodei, Paul Chris-  
1023 tiano, and Geoffrey Irving. 2020. [Fine-tuning lan-  
1024 guage models from human preferences](#).

# Appendix: Large Language Models Can Be Contextual Privacy Protection Learners

1025

1026

## A Notations

1027

Important notations used in the paper are included in Table. 5.

Table 5: Notations used in this paper.

Notation	Description
$w_i, \mathbf{w}_i$	a token and its contextualized embedding
$s$	a natural language sequence
$D = \{s\}$	Fine-tuning dataset
$T$	Annotation
$n$	Maximum sequence length
$\Theta_n$	Set of n-grams associated with PII
$R$	Removed sequence: $R = (r_0, r_1, \dots, r_{k-1})$
$r_{i-1}$	i-th token with $p_i = 0$ in sequence $R$
$C$	Cleaned sequence: $C = (c_0, c_1, \dots, c_{n-1})$
$y_i$	Token $c_i$ if $p_i = 0$ , or the special token $u$ if $p_i = 1$
$u$	Special token added to the vocabulary (e.g., <i>unk</i> for LLaMA2)
$\mathbb{P}(\cdot)$	Probability

1028

### A.1 Detailed Datasets Description

1029

Table. 6 shows more details about datasets:  $S$  denotes the size of the train/test set and  $L_Q/L_A$  denotes the average length (number of tokens) of the question/answer fields. (1) pii-medical\_flashcards with 28861 training and 5093 testing samples; (2) pii-wikidoc with 8500 training and 1500 testing samples; (3) pii-wikidoc\_patient\_information with 5050 training and 891 testing samples.

1030

1031

1032

1033

Table 6: Statistics of datasets

Dataset	Train			Test		
	$ S $	$L_Q$	$L_A$	$ S $	$L_Q$	$L_A$
medical-flashcards	28861	14.59	14.36	5093	53.64	52.74
medical-wikidoc	8500	9.88	9.67	1500	132.04	136.60
wikidoc-patient-information	5050	8.15	8.04	891	73.40	71.10

### A.2 Metrics

1034

**ROUGE (Recall-Oriented Understudy for Gisting Evaluation)** We adopt the popularly used ROUGE-1, ROUGE-2, ROUGE-L (Lin, 2004) and BERTScore (Zhang\* et al., 2020) to evaluate the answer quality in the testing phase. Here we give a detailed definition of these scores. We denote the set of tokens from the generated text as  $G$ , and the set of tokens from the reference text as  $R$ . The number of overlapping unigrams between  $G$  and  $R$  as  $O_1(G, R)$ , and the number of overlapping bigrams between  $G$  and  $R$  as  $O_2(G, R)$ . The total number of unigrams in  $R$  as  $U(R)$  and the total number of bigrams in  $R$  as  $B(R)$ . The longest common subsequence (LCS) between  $G$  and  $R$  as  $L(G, R)$ .

1035

1036

1037

1038

1039

1040

1041

#### ROUGE-1:

1042

$$\text{ROUGE-1} = \frac{O_1(G, R)}{U(R)}$$

1043

#### ROUGE-2:

1044

$$\text{ROUGE-2} = \frac{O_2(G, R)}{B(R)}$$

1045

#### ROUGE-L:

1046

$$\text{ROUGE-L} = \frac{L(G, R)}{\max(|G|, |R|)}$$

1047

## BERTScore

1048  $E$  : BERT encoder or model  
1049  $E(G)$  : Embedding of the entire sequence  
1050 of the generated text  $G$ , produced by  $E$   
1051  $E(R)$  : Embedding of the entire sequence  
1052 of the reference text  $R$ , produced by  $E$   
1053  $c(E(G), E(R))$  : Cosine similarity between the  
1054 sequence embeddings  $E(G)$  and  $E(R)$

1055 Then, the BERTScore between a generated text  $G$  and a reference text  $R$  at the sequence level is defined  
1056 as:

$$1057 \text{BERTScore}(G, R) = c(E(G), E(R))$$

1058 Here, the BERT model  $E$  encodes the entire sequences  $G$  and  $R$  into their respective embeddings, and  
1059 then we compute the cosine similarity between these sequence embeddings to obtain the BERTScore.

## 1060 B Additional Related Work

1061 **Pretraining with Preferences.** Another solution is to maintain the content, but use redesigned loss/  
1062 conditional tags to control the information injected into the LLMs. Pretraining with conditional human  
1063 preference scores can offer a Pareto-optimal and simple approach to reduce the undesirable content by  
1064 up to an order of magnitude. [Korbak et al. \(2023\)](#) compared with the classical pretraining approach.  
1065 While pretraining LLMs conditioned under annotation scores can offer better performance in the human  
1066 preferences aspect. Since human preferences are injected into the models during the pretraining stage,  
1067 the models are biased toward those preferences once they are trained. With the expanding size of LLMs,  
1068 they become increasingly resistant to forgetting their training data ([Carlini et al., 2022](#); [Vu et al., 2022](#);  
1069 [Ramasesh et al., 2022](#); [Korbak et al., 2023](#)). In other words, pretraining large language models conditioned  
1070 under preference score sacrifices some flexibility. Still, it is undeniable that it can provide much better  
1071 alignment with human preferences compared with the classical pretraining schema.

## 1072 C Illustration of Vanilla Tuning and Corpus Curation

1073 This section gives an illustration of Vanilla Tuning (Figure. 3(a)) and Corpus Curation (Figure. 3(b)).

## 1074 D Experiment Details.

### 1075 D.1 Hardware and Implementations

1076 In this paper, we implemented our method on two Linux servers with 4 NVIDIA A100 GPUs, each with  
1077 80GB of memory. The CUDA version is 12.2 and the Driver version is 535.54.03. We used Python  
1078 3.10.12 and Pytorch 2.0.1 ([Paszke et al., 2019](#)) to construct our project. The fine-tuning of LLaMA models  
1079 takes 20 hours on average.

### 1080 D.2 Dataset and Hyperparameters

1081 In our experiments, we use grid search to obtain the best performance. We provide all of the hyperparame-  
1082 ters as well as their configurations in the following:

- 1083 • **Dataset.** For training, we sub-sampled 85% from the three datasets. The performance of each  
1084 method is evaluated on the remaining 15% of data. Dataset details can be found in Table. 6.
- 1085 • **Hyperparameters.** For the parameter optimizer, we chose AdamW with `weight_decay` set to 0. The  
1086 learning rate is set to  $1e^{-4}$ . We use the StepLR learning rate scheduler with `gamma` set to 0.85.  
1087 Epochs and Batch Size: The number of fine-tuning epochs is set to 5, and the batch size is set to 64.

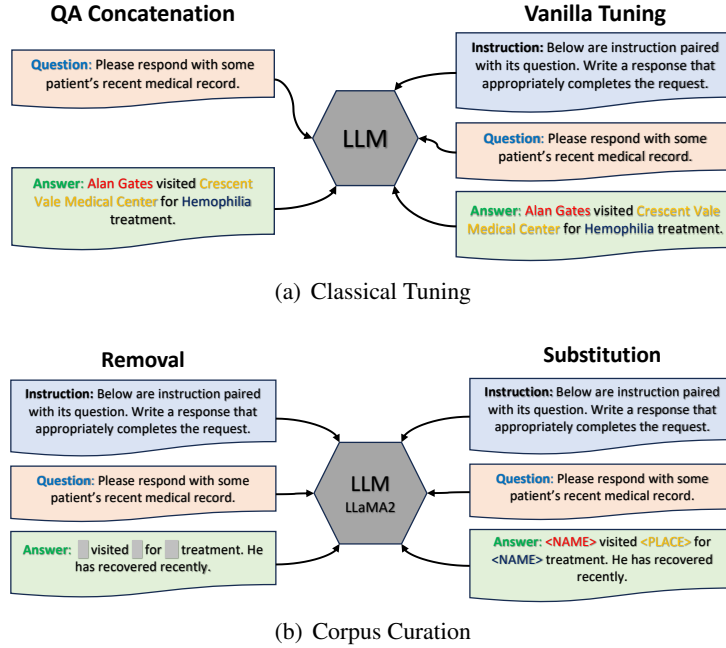


Figure 3: Vanilla, Removal, Substitution.

### D.3 Pareto Frontier of Utility and Privacy Protection

We also report the pareto frontier of Utility and Privacy Protection in Figure 4, 5, 6, 7, 8 and 8, respectively, to evaluate both performance and privacy preservation. It is obvious that the instruction-based approaches consistently align with the Pareto frontier ( $IT$  methods constitute the border of the frontier). Such a phenomenon indicates that employing instructions supplemented by both positive and negative examples achieves the optimal trade-off between performance (utility) and privacy protection of PII. The outcomes strongly support our position that LLMs can be good contextual privacy protection learners.

### D.4 Curve of Knowledge Injection and PII Leakage vs. Learning Process

In this section, we analyze the ROUGE, BERTScore, and Privacy Leakage Score concerning the training steps. We aim to assess whether our two primary learning objectives are effectively achieved throughout the training process. Initially, in Figure. 10 that visualizes the training of  $IT_{PN_1}$ , we observe that as the LM undergoes the training process, we witness a notable trend: the injection of knowledge into the LM steadily increases. This infusion of knowledge corresponds to a progressive rise in both ROUGE and BERTScore, ultimately leading to a stabilization, or convergence, of these metrics. Simultaneously, the Privacy Leakage Score exhibits an intriguing behavior. At the outset of the learning process, it experiences an upward trajectory. This ascent is a direct consequence of the LM ingesting more knowledge, including private tokens, inadvertently learning about sensitive information. However, as training continues, a pivotal shift occurs. The LM's instruction to conceal privacy-related information gradually takes effect, resulting in a discernible decrease in the Privacy Leakage Score. In summary, Figure. 10 offers a compelling visualization of the evolving relationship between knowledge injection, linguistic performance (ROUGE/BERTScore), and privacy protection ( $S_{Priv}$ ) as the LM matures throughout its training steps. It underscores the dynamic equilibrium between knowledge acquisition and safeguarding sensitive data, emphasizing the importance of a well-orchestrated learning process to achieve both objectives.

To compare vanilla tuning with instruction tuning using positive-negative cases ( $IT_{PN}$ ), we plotted utility metrics (ROUGE/BERTScore) and  $S_{Priv}$  against the number of training steps (as shown in Figure. 11). With vanilla tuning, as training progresses, the LLM's performance improves. However, it is accompanied by an increase in privacy leakage. Such a trend corroborates our intuition that, as the LLM assimilates information, it also inadvertently memorizes PII tokens from the corpus. When it comes to instruction tuning with positive-negative cases (Figure. 10), the utility curve exhibits a trajectory akin to

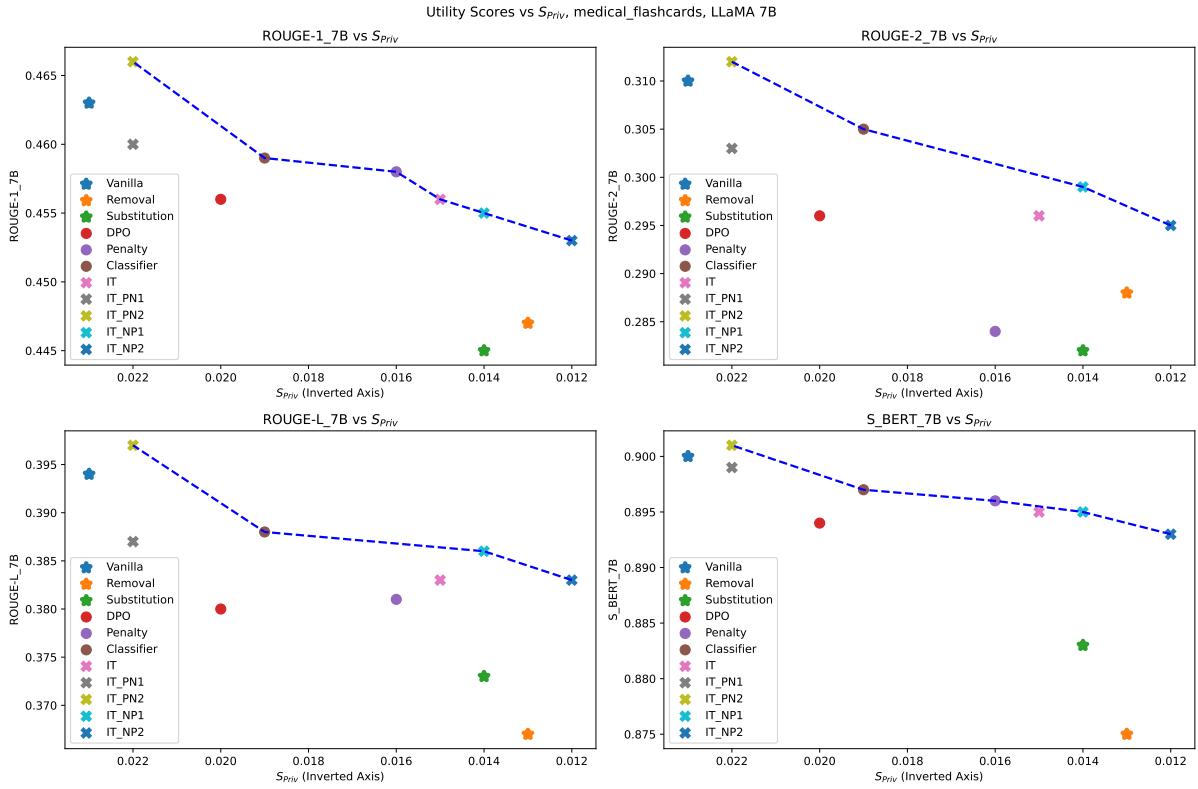


Figure 4: Pareto Frontier on medical\_flashcards, LLaMA2-7B

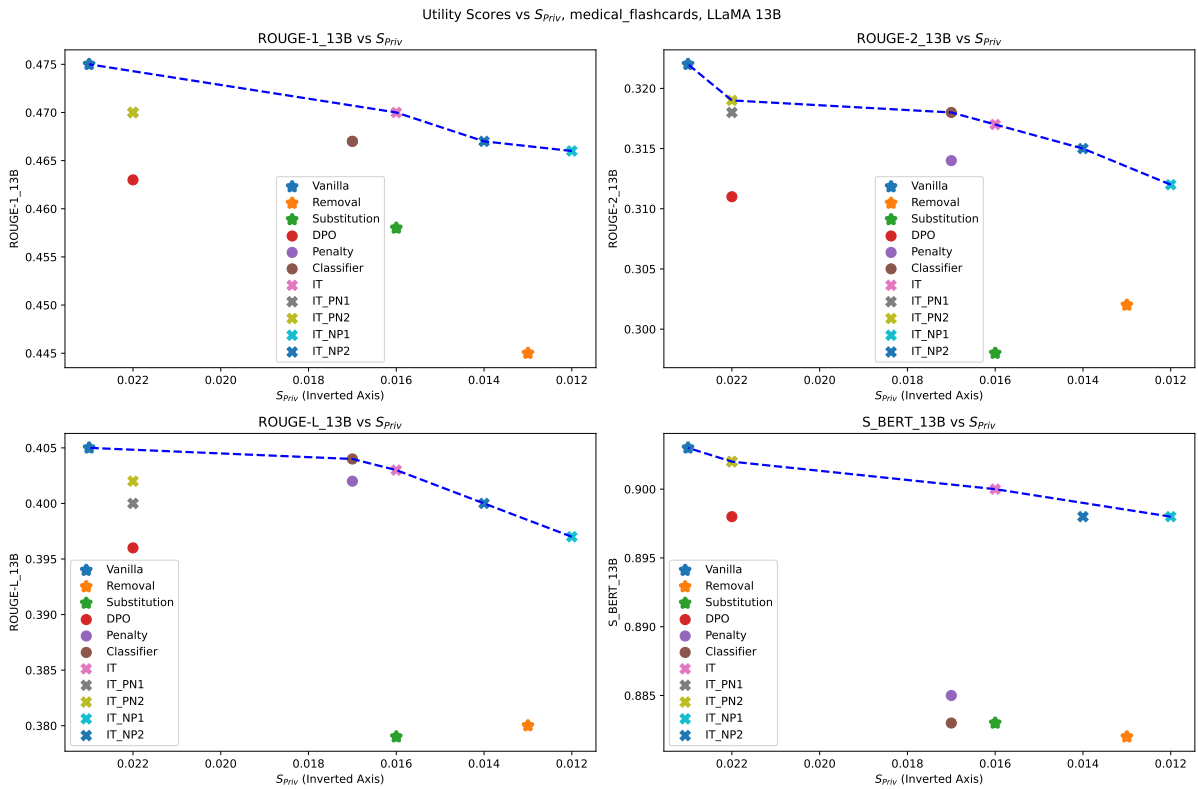


Figure 5: Pareto Frontier on medical\_flashcards, LLaMA2-13B



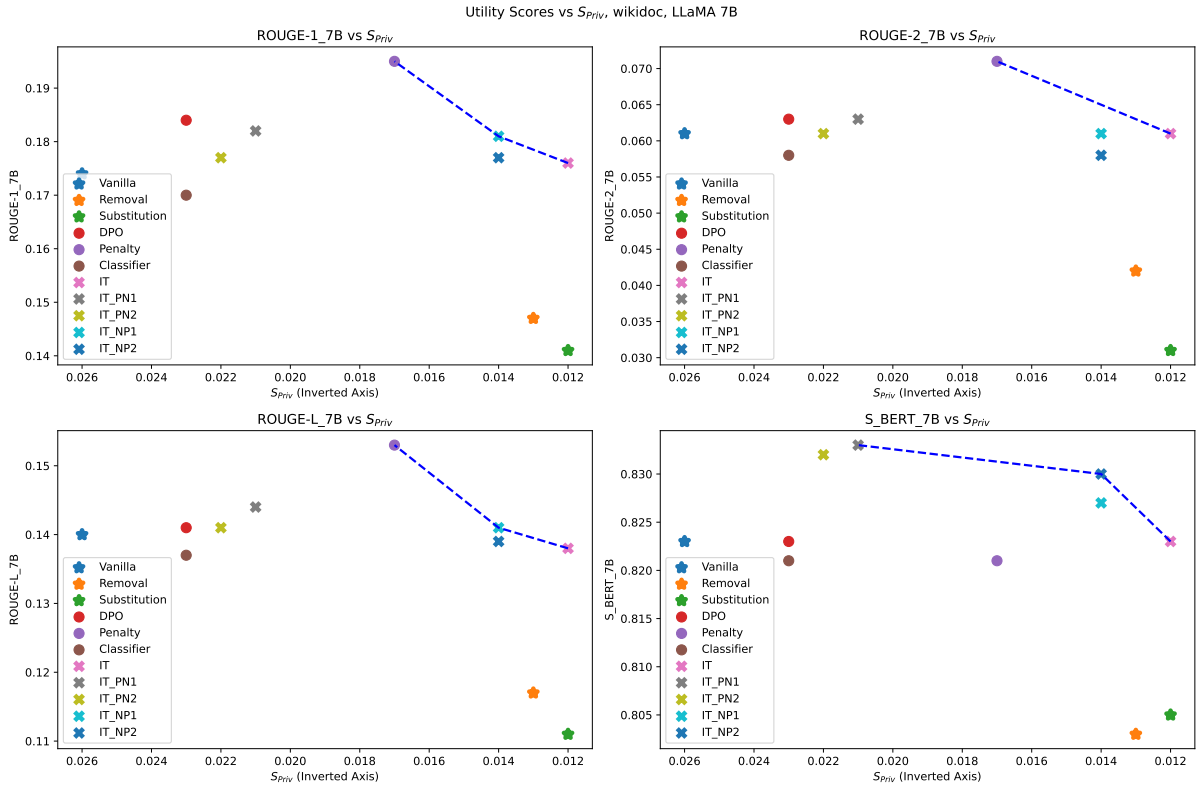


Figure 6: Pareto Frontier on wikidoc, LLaMA2-7B

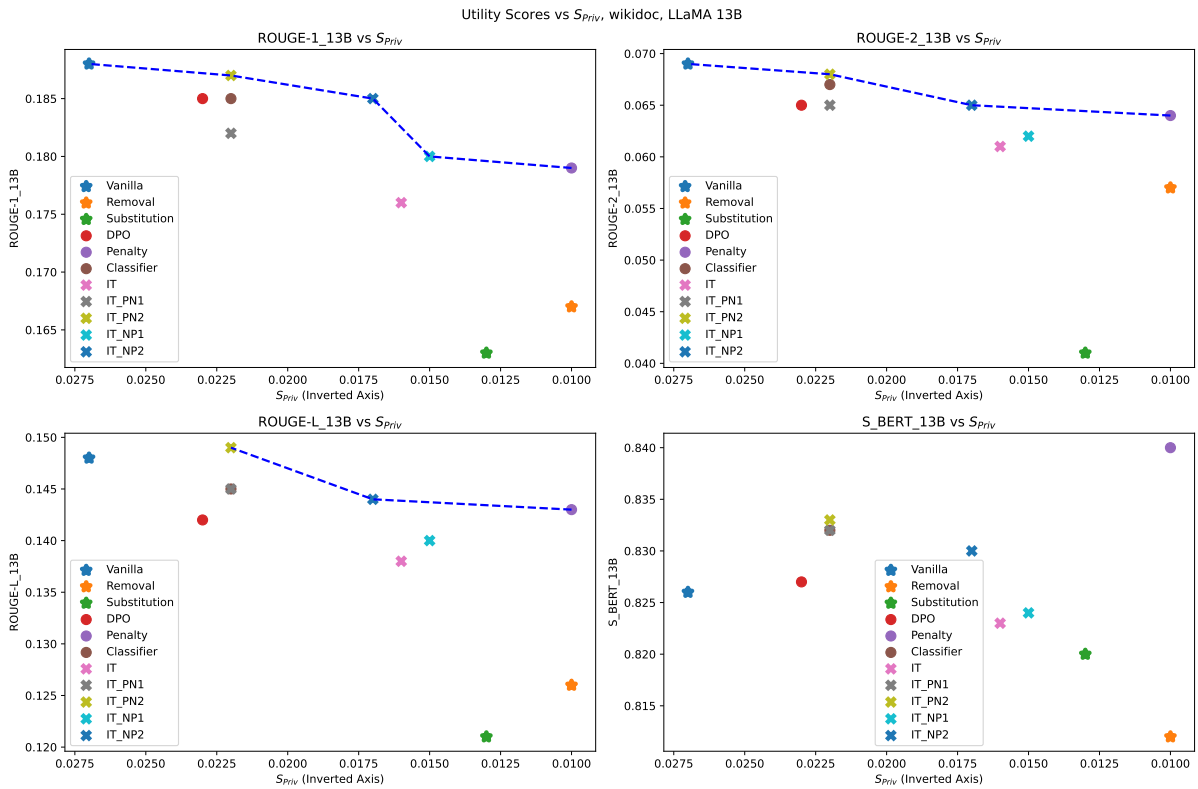


Figure 7: Pareto Frontier on wikidoc, LLaMA2-13B

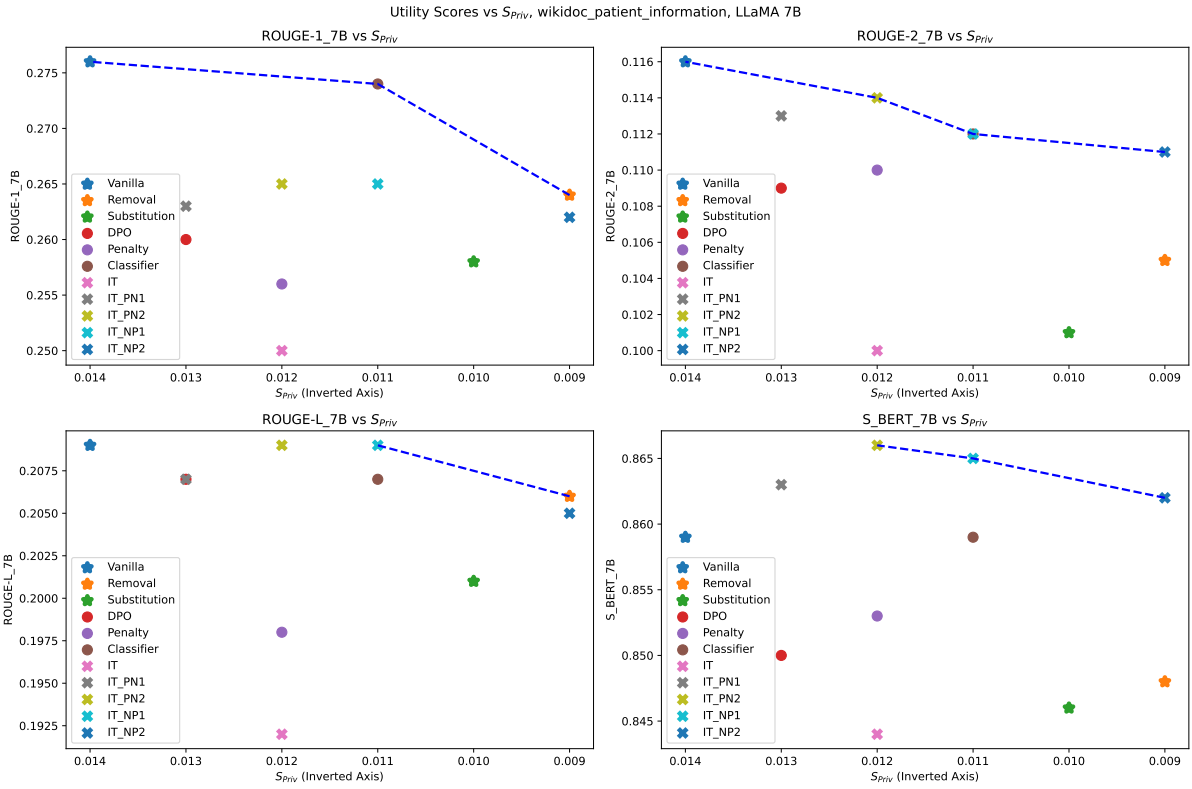


Figure 8: Pareto Frontier on wikidoc\_patient\_information, LLaMA2-7B

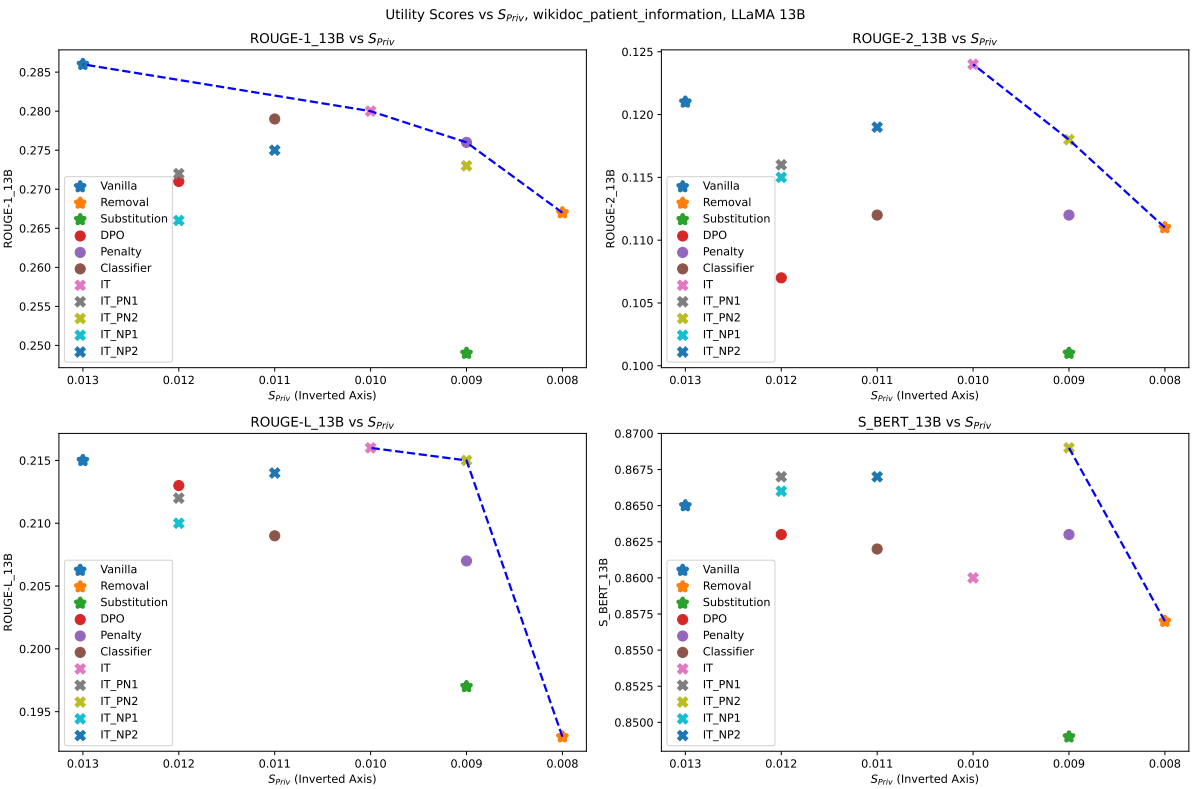


Figure 9: Pareto Frontier on wikidoc\_patient\_information, LLaMA2-13B

vanilla tuning. However, privacy leakage increases initially but eventually declines. This suggests that, by employing instruction combined with positive-negative cases, LLMs can be trained to be good contextual privacy learners.

1117  
1118  
1119

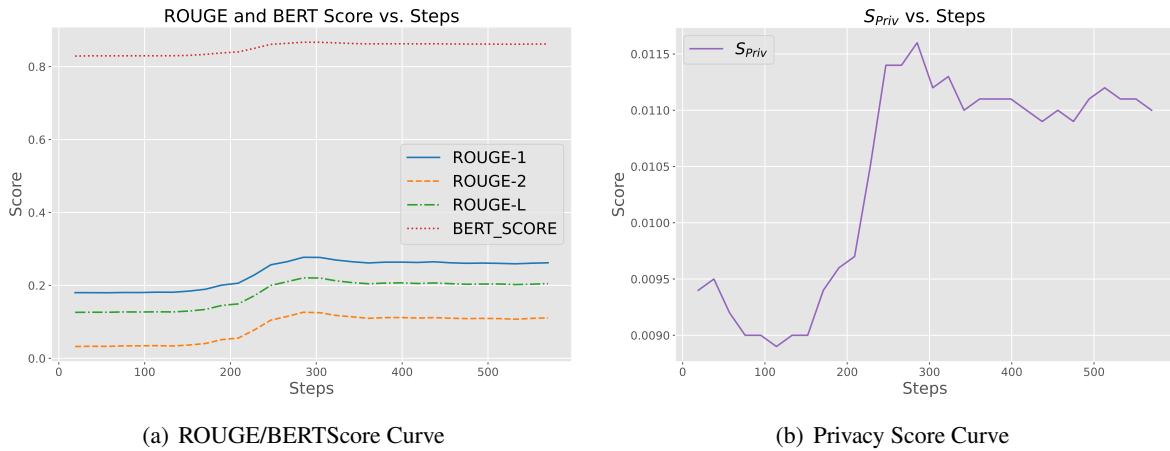


Figure 10: ROUGE, BERTScore, and  $S_{Priv}$  vs. Steps

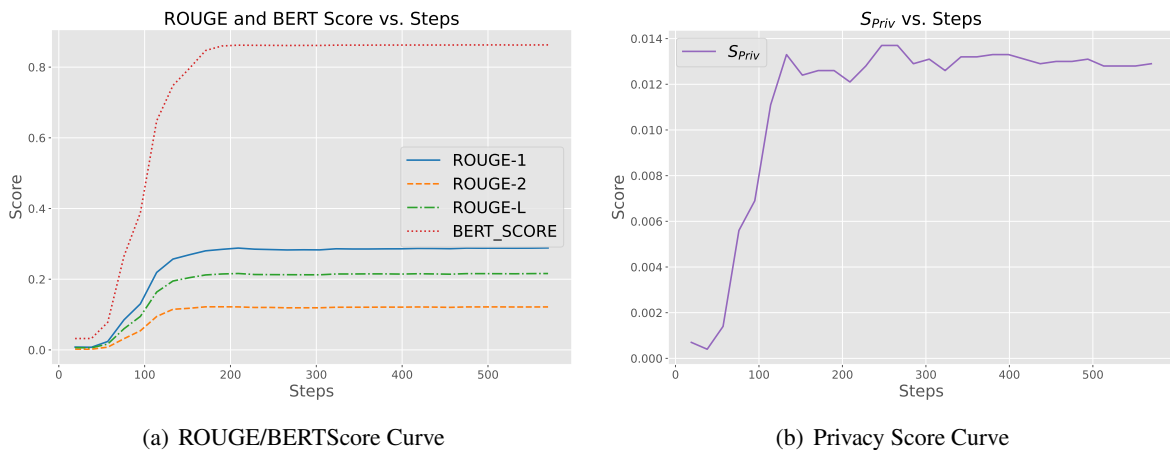


Figure 11: ROUGE, BERTScore, and  $S_{Priv}$  vs. Steps (Vanilla)

## D.5 Instruction Templates

In this section, we list the detailed two instruction tuning templates as shown in Table. 7 and Table. 8.

1120  
1121  
1122

## E Case Study

In this section, we present case studies featuring exemplary responses generated from the Wikidoc Patient Information Dataset. We compare responses generated by the vanilla model (original) with those generated by a model equipped with a token classifier and responses generated by the vanilla model (original) versus the model trained using the instruction tuning (IT) strategy. By examining the results in Tables 9 and 10, we can discern the efficacy of the proposed methodologies in safeguarding privacy.

1123  
1124  
1125  
1126  
1127

---

**### *IT***

Below is a question and related responses. Write a privacy protection response to the question. Sensitive information in the response should be replaced by NAME, ORGANIZATION, etc.

**### *Input***

{{Question filled in there}}

**### *Response***

{{Answer filled in there}}

---

**### *IT<sub>PN<sub>1</sub></sub>***

Below is a question and related responses. Sensitive information in the response should be replaced by NAME, ORGANIZATION, etc. Write

(1) desired answer to the question.

(2) undesired answer to the question.

**### *Input***

{{Question filled in there}}

**### *Response***

(1) desired answer to the question: {{Answer with PII substituted filled in there}}

(2) undesired answer to the question: {{Answer filled in there}}

---

**### *IT<sub>PN<sub>2</sub></sub>***

Below is a question and related responses. Sensitive information in the response should be replaced by NAME, ORGANIZATION, etc. Write

(1) a privacy protection version of the response.

(2) a response answering the question.

**### *Input***

{{Question filled in there}}

**### *Response***

(1) a privacy protection version of the response: {{Answer with PII substituted filled in there}}

(2) a response answering the question: {{Answer filled in there}}

---

Table 7: Templates (Positive-Negative) for instruction tuning (IT) and IT with positive and negative cases.

---

**###  $IT_{NP_1}$**

Below is a question and related responses. Sensitive information in the response should be replaced by NAME, ORGANIZATION, etc. Write

- (1) a response answering the question.
- (2) a privacy protection version of the response.

**### *Input***

{{Question filled in there}}

**### *Response***

- (1) a response answering the question: {{Answer filled in there}}
  - (2) a privacy protection version of the response: {{Answer with PII substituted filled in there}}
- 

**###  $IT_{NP_2}$**

Below is a question and related responses. Sensitive information in the response should be replaced by NAME, ORGANIZATION, etc. Write

- (1) undesired answer to the question.
- (2) desired answer to the question.

**### *Input***

{{Question filled in there}}

**### *Response***

- (1) undesired answer to the question: {{Answer filled in there}}
  - (2) desired answer to the question: {{Answer with PII substituted filled in there}}
- 

Table 8: Templates (Negative-Positive) for instruction tuning (IT) and IT with positive and negative cases.

---

*Prompt*

**### Input**

What is the historical background of chronic fatigue syndrome?

*Answer (Original)*

The first case of chronic fatigue syndrome was reported in 1956. In 1988, the Centers for Disease Control and Prevention (CDC) first reported a cluster of cases of chronic fatigue syndrome in Incline Village, Nevada ...

*Answer (Token Classifier)*

The first case of chronic fatigue syndrome was reported in 1956. In 1988, the Centers for Disease Control and Prevention (CDC) first reported a cluster of cases of chronic fatigue syndrome in ⟨unk⟩.

---

**### Input**

What is the explanation for using images and the new advanced syntax to insert them?

*Answer (Original)*

The following is a list of the most common image tags:

#ev:youtube <http://www.youtube.com/watch?v=dQw4w9WgXcQ> ...

*Answer (Token Classifier)*

The following is a list of the most common image tags: ⟨unk⟩

---

**### Input**

What information is available on Ranolazine?

*Answer (Original)*

Ranolazine may cause side effects. Tell your doctor if any of these symptoms are severe or do not go away: dizziness, diarrhea, headache, nausea, vomiting ...

Some side effects can be serious. The following symptoms are uncommon, but if you experience any of them, call your doctor immediately: fainting fast, pounding, or irregular heartbeat or pulse swelling of the hands, ankles, or feet

If you experience a serious side effect, you or your doctor may send a report to the Food and Drug Administration's (FDA) MedWatch Adverse Event Reporting program online at <http://www.fda.gov/MedWatch/report.htm> or by phone 1-800-332-1088.

*Answer (Token Classifier)*

Ranolazine may cause side effects. Tell your doctor if any of these symptoms are severe or do not go away: dizziness, diarrhea, headache, nausea, vomiting ...

Some side effects can be serious. The following symptoms are uncommon, but if you experience any of them, call your doctor immediately: fainting fast, pounding, or irregular heartbeat or pulse swelling of the hands, ankles, or feet

If you experience a serious side effect, you or your doctor may send a report to the ⟨unk⟩ program online at ⟨unk⟩ or by phone ⟨unk⟩.

---

Table 9: Example responses generated on the Wikidoc Patient Information Dataset by the vanilla model (original) and the model with token classifier.

---

**### Input**

Could you provide a brief explanation of familial hypercholesterolemia?

*Answer (Original)*

The prevalence of familial hypercholesterolemia is approximately 1 in 500 individuals worldwide.

*Answer (Instructed)*

The diagnosis of familial hypercholesterolemia is based on the following criteria:

Low-density lipoprotein cholesterol (LDL-C) level of 190 mg/dL or greater in a child or adolescent or 160 mg/d.

---

**### Input**

What is the medical treatment for non-small cell lung cancer?

*Answer (Original)*

The treatment of non-small cell lung cancer depends on the stage of the disease.

*Answer (Instructed)*

The treatment of non-small cell lung cancer depends on the stage of the disease. The treatment options include surgery, chemotherapy, radiation therapy, and targeted therapy.

---

**### Input**

What is an ultrasound for peripheral arterial disease?

*Answer (Original)*

Ultrasound is the first line imaging modality for the diagnosis of peripheral arterial disease. It is a non-invasive, safe, and cost-effective method for the diagnosis of peripheral arterial disease. Ultrasound is used to assess the anatomy of the arteries and the degree of stenosis. It is also used to assess the anatomy of the arteries and the degree of stenosis.

*Answer (Instructed)*

Peripheral arterial disease (PAD) is a common disease that affects the arteries of the legs. It is caused by atherosclerosis, which is the buildup of plaque in the arteries. Ultrasound is a noninvasive test that can be used to diagnose PAD. It uses sound waves to create images of the arteries in the legs. Ultrasound can be used to measure the blood flow in the arteries and to look for blockages. It can also be used to measure the size of the arteries. Ultrasound is a safe and painless test.

---

Table 10: Example responses generated on the Wikidoc Patient Information Dataset by the vanilla model (original) and the model trained with the instruction (IT) strategy.