PRIVACYRESTORE: PRIVACY-PRESERVING INFERENCE IN LARGE LANGUAGE MODELS VIA PRIVACY REMOVAL AND RESTORATION

Anonymous authors

Paper under double-blind review

Abstract

The widespread usage of online Large Language Models (LLMs) inference services has raised significant privacy concerns about the potential exposure of private information in user inputs to malicious eavesdroppers. Existing privacy protection methods for LLMs suffer from either insufficient privacy protection, performance degradation, or large inference time overhead. To address these limitations, we propose PrivacyRestore, a plug-and-play method to protect the privacy of user inputs during LLM inference. The server first trains restoration vectors for each privacy span and then release to clients. Privacy span is defined as a contiguous sequence of tokens within a text that contain private information. The client then aggregate restoration vectors of all privacy spans in the input into a single meta restoration vector which is later sent to the server side along with the input without privacy spans. The private information is restored via activation steering during inference. Furthermore, we prove that PrivacyRestore inherently prevents the linear growth of the privacy budget. We create three datasets, covering medical and legal domains, to evaluate the effectiveness of privacy preserving methods. The experimental results show that PrivacyRestore effectively protects private information and maintain acceptable levels of performance and inference overhead.

028 029

031

006

008 009 010

011

013

014

015

016

017

018

019

021

023

025

026

027

1 INTRODUCTION

Large Language Models (LLMs) have emerged as powerful tools in various domains, including healthcare (Chen et al., 2023; Xu et al., 2023), law (Wu et al.; Deng et al., 2023a), and finance (Wu et al., 2023; Xie et al., 2023). With the exception of a very small portion of users who have the resources and expertise to deploy LLMs locally, the vast majority of users access and interact with these powerful models through online inference services.

However, the widespread usage of online LLMs inference services has raised significant privacy concerns, especially regarding the potential risk of private information being leaked through user inputs when interacting with LLMs deployed on cloud platforms. User inputs often contain sensitive information such as details in medical records and legal cases. Potential threats may arise from eavesdropper attackers intercepting user queries during transmission to cloud platforms for malicious purposes. For example, in sensitive domains like medical diagnosis, if a user's input containing personal health information, such as "*I was previously diagnosed with HIV, and lately I've been experiencing fever and diarrhea...*" is disclosed, it may cause troubles to their life.

In this paper, we focus on protecting the private information contained in user inputs during LLM inference stage. In this setting, the client submits inputs to the server (also known as the service provider) and there is a risk that inputs might be disclosed by attackers. Current methods for protecting user inputs can be categorized into two categories: Secure Multi-Party Computation (SMPC) and Differential Privacy (DP). SMPC based methods (Hao et al., 2022b; Li et al., 2023a; Liang et al., 2024) utilize encryption protocols and algorithms to enable collaborative computation without revealing original data to others. However, SMPC methods have large inference time overhead, making them impractical for real-time applications. For example, running a single pass inference on the RoBERTa-Base (Liu et al., 2019) requires 168.43 seconds (Hao et al., 2022a). DP based methods (Feyisetan et al., 2020; 2019; Xu et al., 2020; Bo et al., 2021) introduce the definition of

062 We propose PrivacyRestore which directly removes privacy spans in user inputs and restores private 063 information via activation steering (Li et al., 2023c; Turner et al., 2023; Hernandez et al., 2023) 064 during model inference. Our method is based on two key assumptions: (a) Private information is confined within specific a contiguous sequence of tokens, termed "privacy span", rather than 065 being dispersed throughout the entire input. Privacy span is defined as a contiguous sequence of 066 tokens within a text that contain private information. The removal or proper redaction of privacy 067 spans significantly impedes unauthorized parties from reconstructing or inferring the underlying 068 private information. For instance, if privacy spans "HIV", "fever" and "diarrhea" are removed from 069 a medical record "I was previously diagnosed with HIV, and lately I've been experiencing fever and diarrhea...", attackers can not recover any private information. (b) In a particular domain, 071 the number of potential privacy spans is limited and finite. For example, in the application of 072 medical diagnosis, privacy spans generally pertain to symptoms and disease names, and the number of 073 possible symptoms and disease names is inherently limited. Moreover, although medical knowledge 074 and terminology inevitably evolve, the core set of symptoms and diseases remains relatively stable 075 and finite.

076 PrivacyRestore operates in two stages: the preparation stage and the inference stage. In the prepa-077 ration stage, we first identifies the attention heads where the activation steering occurs. Second, 078 each privacy span is encoded to a vector named restoration vector. This stage is conducted on the 079 server side. Our method is plug-and-play, requiring only the restoration vectors to be trainable, while keeping the LLM frozen. Once training is complete, the users keep all restoration vectors on the 081 client side. In the inference stage, the user construct a meta vector by first estimating the importance 082 of each privacy span in the input and then calculating a weighted sum of the corresponding restoration 083 vectors. The user then removes the privacy spans from the input and submits the remaining input along with the meta vector to the server. The server uses the meta restoration vector to restore the 084 removed privacy spans through activation steering. 085

To prevent the leakage of privacy spans via reverse-engineering on the meta vector, the d_{χ} -privacy mechanism (Feyisetan et al., 2020) is applied on the meta vector before transmission on the client side. d_{χ} -privacy mechanism is a variant of the differential privacy mechanism (Dwork et al., 2016). By applying d_{χ} -privacy to the meta restoration vector instead of words, our method inherently addresses the issue of linear growth of privacy budget (Mattern et al., 2022a) commonly encountered in d_{χ} privacy and other DP variants. Experimental results demonstrate that the proposed method effectively protects private information and maintains satisfactory performance and inference efficiency.

- ¹⁹³ The contributions of our paper are summarized as follows,
 - We propose a plug-and-play privacy protection method that removes privacy spans in the input and restores private information via activation steering during inference.
 - We propose Attention-aware Weighted Aggregation to construct the meta vector and apply the d_{χ} -privacy mechanism to the meta vector, inherently addressing the problem of the linear growth of privacy budget.
 - We construct three datasets, covering the medical and legal fields, to evaluate our method. Experimental results demonstrate its capabilities of privacy protection. It also maintains acceptable performance and inference efficiency.
- 102 103 104

094

096

098

099 100

- 2 RELATED WORKS
- 105 106
- 107 In this section, we introduce the related works on user input protection methods, which are currently divided into two categories: SMPC-based methods and DP-based methods.

108 2.1 SECURE MULTI-PARTY COMPUTATION (SMPC)

109 110

Secure multi-party computation (SMPC) methods utilize multi-party encryption algorithms to enable
 collaborative computation among multiple parties while protecting the privacy of their data. However,
 most nonlinear operations in LLMs cannot directly support secure multi-party computation. To
 address this challenge, current SMPC methods focus on two optimization directions: model structure oriented optimization and protocol-oriented optimization.

115 The model structure-oriented approach aims to replace SMPC-unfriendly nonlinear operations with 116 SMPC-friendly alternatives. For instance, MPC-Former (Li et al., 2023a) approximates nonlinear 117 operations in Transformer using polynomials and maintains performance through model distillation. 118 MERGER (Liang et al., 2024) integrates previous techniques to natural language generation (NLG) 119 tasks by bypassing embedded computation and reorganizing linear operations in Transformer modules, 120 further enhancing computational efficiency and model performance. In contrast, the protocol-oriented 121 approach focuses on designing efficient SMPC operators for nonlinear operations in LLMs while preserving the original model structure. Recent works Hao et al. (2022b); Liu & Liu (2023); Zheng 122 et al. (2023b); Gupta et al. (2023) have improved the efficiency of nonlinear operations in privacy-123 preserving LLMs inference by utilizing various SMPC protocols, such as confusion circuit and 124 function secret sharing. 125

Although SMPC-based methods can be applied to protect user inputs during model inference, they
 still suffer from large inference time overhead. For example, inference on the RoBERTa-Base model
 takes 168.43 seconds (Hao et al., 2022a), making current SMPC methods impractical for online LLM
 inference services.

130 131

132 133

2.2 DIFFERENTIAL PRIVACY (DP)

Differential Privacy (DP), as introduced by Dwork et al. (2016), is designed to protect individual privacy by preventing attackers from identifying specific participants in a dataset. Several variants of DP have been developed to enhance privacy protection across various settings, adapting the core principles of DP to different types of data and threat models. Notable examples include Centralized Differential Privacy (CDP), Local Differential Privacy (LDP) and d_{χ} -privacy.

CDP (Dwork et al., 2016) operates under the assumption that all data has been stored in a central repository. It guarantees that attackers cannot distinguish between any two adjacent repositories based on query results. In contrast, LDP (Duchi et al., 2013) provides a stronger guarantee, ensuring that attackers cannot distinguish between any two adjacent inputs. Mattern et al. (2022b) and Utpala et al. (2023) propose using paraphrasing techniques to achieve LDP on user inputs. The formal definitions of CDP and LDP are provided in Appendix D.

LDP allocates the same privacy budget ϵ to all adjacent input pairs, regardless of their similarity. Applying the same ϵ forces each user input to be indistinguishable from any other, which can negatively impact data utility. However, it is sufficient for privacy protection to make each user input indistinguishable only from its closer counterparts. To address this, d_{χ} -privacy (Feyisetan et al., 2019), a relaxed version of LDP, incorporates metrics that measure the similarity between inputs, allowing for more flexible control over the privacy budget. d_{χ} -privacy is defined as,

Definition 2.1. $(d_{\chi}$ -privacy). A randomized mechanism $\mathcal{M} : \mathcal{I} \to \mathcal{O}$ fulfills ϵ - d_{χ} -privacy if for all adjacent inputs $I, I' \in \mathcal{I}$ and all possible outputs $O \subset \mathcal{O}$,

$$\mathbb{P}\left(\mathcal{M}(I) \in O\right) \le \exp(\epsilon d_{\chi}(I, I'))\mathbb{P}\left(\mathcal{M}(I') \in O\right),$$

154 155 156

153

157 where d_{χ} is a distance function defined on \mathcal{I} . Recent works (Feyisetan et al., 2020; Xu et al., 2020; 158 Li et al., 2023d; Qu et al., 2021) leverage d_{χ} -privacy to safeguard user inputs during both inference 159 and fine-tuning phases. However, as noted by Mattern et al. (2022b), all of the aforementioned DP 160 variants suffer from the linear growth of the privacy budget. A larger privacy budget indicates weaker 161 privacy protection, meaning that as the input length increases, the effectiveness of privacy protection 165 diminishes.

¹⁶² 3 THREAT MODEL

We consider a threat model involving two parties: a server that holds the LLM weights and a client holds user inputs containing privacy spans. **Privacy span is defined as a contiguous sequence of tokens within a text that contain private information.** The server provides services through an API, enabling the client to transmit inputs and receive responses while maintaining the confidentiality of the LLM weights. The server may be vulnerable to attacks by adversaries seeking to steal privacy information in user inputs. Our task is to protect these private spans in user inputs from being intercepted to the adversaries, even when the adversaries can attack the server directly.

170 171 172

164

165

166

167

168

4 Methodology

173

174 To protect privacy, PrivacyRestore transmit the input with privacy spans removed instead of the entire 175 input text to the server. The information in the privacy spans is encrypted as a vector, which is then 176 injected with noise and is also sent to the server. We propose to use activation steering (Li et al., 2023c; Turner et al., 2023; Hernandez et al., 2023) to restore privacy information. Activation steering 177 methods modify the activations of a language model at inference time to predictably alter its behavior. 178 Activation steering is widely used to truthfulness enhancement(Li et al., 2023c), LLMs detoxifying 179 (Li et al., 2024), and sentiment modification (Turner et al., 2023). To our best knowledge, it is first 180 attempt to use activation steering for privacy information restoration. We include the preliminaries of 181 our methodology in Appendix C. 182

183 PrivacyRestore operates in two stages, i.e., the preparation stage and the inference stage.

(1) Preparation stage: This stage takes place on the server. We first identify the edited attention heads
and train the restoration vectors for each privacy spans. After training, these vectors are released
to the clients. The preparation stage is conducted offline, prior to the server beginning to offer its
services.

(2) Inference stage: This stage involves collaboration between the client and server. The client constructs a meta vector which is later transmitted to the server along with the input where the privacy spans are removed. The server then performs inference on the input with privacy spans removed and restores those privacy spans via modifying activations.

An overview of PrivacyRestore is shown in Figure 1. Detailed descriptions of the preparation stage and the inference stage are provided in §4.1 and §4.2 respectively. The definitions of all notations used in this paper can be found in Appendix A.

195 196 197

209

4.1 PREPARETION STAGE

Edited Heads Identification. As pointed by activation steering methods (Li et al., 2023c; Chen et al., 2024), modifying all attention heads in LLMs will degrade overall performance. Inspired by this, we aim to identify the attention heads most relevant to privacy spans.

201 As shown in upper part of Figure 1, we firstly utilize the probe technique (Alain & Bengio, 2016; 202 Tenney et al., 2019; Belinkov, 2022) to identify the most relevant attention heads for each privacy 203 span. $I_{all} = \{I_1, ..., I_m\}$ represents the user inputs in the training set where m is the size of training 204 set. Given a privacy span s, $Y_s = \{y_1, ..., y_m\}$ represents the corresponding labels, where $y_i = 1$ 205 only if input I_i contains privacy span s. For each user input I_i , we record the hidden state of last 206 token on each attention head. We then train a binary classifier for each head, tailored to the privacy 207 span s, as the probe. The probe takes the hidden state of last token as input and predicts whether the input contain the privacy span s. The probe is formulated as: 208

$$\mathcal{F}_h^s(\mathbf{u}_h) = \sigma(\theta_h^s \cdot \mathbf{u}_h),\tag{1}$$

where $\mathcal{F}_{h}^{s}(\cdot)$ is the probe of privacy token *s* on head *h*, \mathbf{u}_{h} is the hidden state of last token on head *h*, θ_{h}^{s} is parameters of the probe, and $\sigma(\cdot)$ indicates the sigmoid function. A probe $\mathcal{F}_{h}^{s}(\cdot)$ with higher accuracy indicates a stronger correlation between the head *h* and the privacy span *s*. Therefore, we select the top *K* attention heads with highest accuracies for each privacy span.

Subsequently, we introduce a **Top-K Heads Selector** to identify the common top-K heads set \mathcal{H}_c from the individual top-K heads sets of each privacy span. Using different top-K head sets for



Figure 1: The PrivacyRestore consists of two stages. (1) **Preparation Stage.** This stage is operated 236 on the server side only and is conducted offline before the server starts offering its services. This 237 stage aims to identify the edited heads and train the restoration vectors. (2) Inference Stage. This 238 stage involves the collaboration between the server and the client. The client need to construct a 239 meta vector by computing a weighted sum of restoration vectors for all privacy spans in the input. A 240 local lightweight model is used to estimate the weight of each privacy span. Then the client transmits 241 the meta vector and the incomplete input with privacy spans removed to the server. Using the meta 242 vector, the server restores the privacy information. 243

244 different privacy spans may suffer the risk of privacy leakage, as an attacker could infer the presence 245 of a specific privacy span based on the characteristics of top-K heads set. Hence, we propose Top-K 246 Heads Selector to combine all different top-K heads sets to construct a common top-K heads set \mathcal{H}_c as the edited heads set. To achieve this, we calculate the average score of each head across all privacy 247 spans, selecting the highest K heads to construct the common set. A head receives a positive score 248 if it appears in the top-K head set of a privacy span s. The score is related to the accuracy of probe 249 associated with the head. Specifically, if the probe associated with this head yields higher accuracy, 250 the score is higher. By iterating this process across all privacy spans, we can calculate the average 251 score for each head. The detailed algorithm is described in Appendix G. 252

253 **Restoration Vectors Training.** After identifying the edited heads set, the next step is to train the 254 restoration vectors for each privacy span on the server side. The training objective is to align the 255 predictions given the input with privacy spans removed to be the same as the predictions given an 256 intact input.

257 For each privacy span $s \in S$, there is a trainable restoration vector r_s^h for each head h in the common 258 top-K heads set \mathcal{H}_c . Restoration vectors of all privacy spans on all heads of \mathcal{H}_c form the only trainable 259 parameters Θ in our method. The LLM weights remain frozen. Our method is plug-and-play and 260 parameter-efficient for training. We fine-tune these restoration vectors using ORPO loss proposed by 261 Hong et al. (2024): 262

$$\operatorname{ratio}(a|\hat{I};\Theta) = \frac{\mathbb{P}(a|\hat{I};\Theta)}{1 - \mathbb{P}(a|\hat{I};\Theta)},\tag{2}$$

267 268

217

221

222

224

225

229

230

231

233

$$\mathcal{L}_{\text{ORPO}} = \sum_{\hat{I} \in \hat{I}_{all}} -\log \mathbb{P}(a|\hat{I};\Theta) - \lambda \log \sigma \left(\log \frac{\operatorname{ratio}(a|\hat{I};\Theta)}{\operatorname{ratio}(\hat{a}|\hat{I};\Theta)} \right),$$
(3)

where \hat{I} denotes the input with privacy spans removed and $\hat{I}_{all} = \{\hat{I}_1, \dots, \hat{I}_m\}$ represents the 269 training set of incomplete inputs, a is the initial output give the complete input, \hat{a} is the output given the incomplete input with privacy spans removed and λ represents the coefficient. The ORPO loss encourages the model to generate the initial output *a* rather than the output \hat{a} given the incomplete input. After restoration vectors training, the server will release all restoration vectors to clients.

2744.2INFERENCE STAGE2754.2

290 291

292 293

294

295

296

297

298

299 300

301 302 303

321

322

Meta Vector Construction. In the inference stage, the client construct a meta vector which is later
 transmitted to the server along with the input with privacy spans removed, as shown in the lower
 left panel in Figure 1. These operations are conducted on the client side. Transmitting a single meta
 vector instead of multiple restoration vectors reduces the communication burden and the risk of data
 leakage. For instance, adversaries could easily know the number of privacy spans in the user input if
 multiple restoration vectors were sent.

282 However, equal weighted aggregation may weaken the influence of critical spans and amplify the 283 effect of irrelevant ones. Therefore, we propose a novel method called Attention-aware Weighted 284 Aggregation (AWA) which estimates a weight for each privacy span, and then take the weighted 285 sum of restoration vectors as the aggregation result. This result is then added with noise for privacy 286 protection and transmitted to the server. Considering the limitation of computing resource, we propose to utilize a lightweight model (e.g., BERT Devlin et al. (2019)) to estimate importance weights on the 287 client side. For the privacy span s in the user input I, the importance weight w_s is calculated as the 288 average attention received by s: 289

$$w_s = \frac{1}{n} \frac{1}{n_h} \sum_{t=1}^n \sum_{h=1}^{n_h} \operatorname{Attn}_h(s, i_t), \tag{4}$$

where n is the number of tokens in the input, n_h is the number of attention heads in the lightweight model, i_t is the t-th token of I, and $Attn_h(s, i_t)$ denotes the attention score of i_t attending to the privacy span s. To simplify the problem, we first consider constructing the meta vector for a single head h. The meta vector \mathcal{R}_h on head h is obtained by computing the weighted sum of the restoration vector r_s^h of each privacy span s on head h, normalizing the summation, and adding noise \mathcal{N} . The process is formulated as follows,

$$Z_h = \frac{\sum_{s \in \mathcal{S}_I} w_s \cdot r_s^h}{||\sum_{s \in \mathcal{S}_I} w_s \cdot r_s^h||_2},\tag{5}$$

$$\mathcal{R}_h = Z_h + \mathcal{N}, \tag{6}$$

where S_I denotes the set of privacy spans in the input I, $|S_I|$ denotes the number of privacy spans and Z_h represents the normalization of the weighted sum on head h. The injected noise \mathcal{N} is sampling from the distribution $p(\mathcal{N}) \propto \exp(-\epsilon ||\mathcal{N}||)$, according to Feyisetan et al. (2020), where ϵ is the privacy hyperparameter.

To construct the meta vector for multiple heads, we first concatenate the restoration vectors on multiple heads. Then, we apply the weighted summary, normalization, and noise addition to the concatenated vector, as we do for a single edited head. After construction, the meta vector, along with the input with the privacy spans removed, are transmitted to the server for inference.

Privacy Restoration. We utilizes the meta vector to restore the missing privacy spans during inference on the input with the privacy spans removed, as illustrated in the lower right part of Figure 1. This operation is conducted on the server side.

Following activation steering methods (Li et al., 2023c; Chen et al., 2024), we apply the meta vector to the outputs of the edited attention heads to restore the privacy spans. Let \mathbf{u}_h represent the hidden state of the last token on head h given the input with privacy spans removed, and \mathcal{R}_h be the meta vector for head h, the hidden state of the last token on head h after restoration $\bar{\mathbf{u}}_h$ is denoted as:

$$ar{\mathbf{u}}_h = \mathbf{u}_h + ||\mathbf{u}_h||_2 \cdot \mathcal{R}_h, \ orall h \in \mathcal{H}_c.$$

(7)

During inference, if a head belongs to the common top-K heads set \mathcal{H}_c , its hidden state should be modified using Eq 7.

324 5 ANALYSIS OF PRIVACY BUDGET

³²⁶ In this section, we analyze the privacy budget of DP variants and PrivacyRestore.

Theorem 5.1. The DP variants including CDP, LDP and d_{χ} -privacy are constrained by a privacy budget that grows linearly with the length of the protected text.

The detail proof of Theorem 5.1 is presented in Appendix E. As the length of the protected text 330 increases, the growing privacy budget makes these DP variants more vulnerable to adversarial attacks, 331 thus compromising their robustness. We also provide empirical evidence demonstrating the linear 332 growth problem of d_{γ} -privacy in Section 6.3. We implement two types of attack, i.e., prompt injection 333 attack (Perez & Ribeiro, 2022; Suo, 2024) and attribute inference attack (Li et al., 2022), across 334 three privacy-preserving datasets. As shown in Figure 3(a) and 3(b), attack performance increases 335 with the length of the protected text, highlighting the linear growth problem of the privacy budget in 336 d_{χ} -privacy. 337

Theorem 5.2. PrivacyRestore fulfills d_{χ} -privacy and provides a privacy budget of $\epsilon ||Z' - Z||$, where ϵ denotes privacy hyperparameter, and Z' and Z represent any pair of normalized weighted sums of restoration vectors concatenated across all edited heads. The privacy budget of PrivacyRestore is independent of the length of protected text.

The detail proof of Theorem 5.2 is provided in Appendix F. The ϵ is privacy hyperparameter, independent of the length of protected text. The ||Z' - Z|| represents the distance between two vectors, also independent of the length of protected text. Since the privacy budget in PrivacyRestore is independent of the length of the protected text, our method effectively protects privacy even with longer protected text, inherently addressing the linear growth issue in DP variants. We also provide empirical evidence to support the theorem. We present the attack performance of PrivacyRestore across varying protected text lengths, as detailed in Section 6.3.

348 349

6 EXPERIMENTS

350 351 352

6.1 EXPERIMENTS SETUP

353 **Datasets.** We evaluate our method in medical and legal domains. However, existing benchmarks, 354 such as DDXPlus (Tchango et al., 2022) and NLICE (Al-Ars et al., 2023) for medical diagnosis, and 355 SLJA (Deng et al., 2023b) for legal judgement, do not specify privacy spans in the input. To address 356 this gap, we leveraged GPT-3.5 (Ouyang et al., 2022) to classify symptoms in DDXPlus/NLICE 357 and case details in SLJA into sensitive and non-sensitive categories, treating the sensitive data as 358 privacy spans. The classification prompt is shown in Appendix O.1. Based on classification results, 359 we curated three privacy-preserving datasets: **Pri-DDXPlus**, **Pri-NLICE** and **Pri-SLJA**, with 149, 64 and 142 types of privacy spans, respectively. The process of dataset construction and statistical 360 361 information can be found in Appendix B.

Metrics. The evaluation access both performance and inference efficiency. For performance 362 evaluation, we employ MC1 and MC2 (Lin et al., 2021) to measure the model's accuracy in selecting 363 the correct answer among 4 options. Each sample in Pri-DDXPlus, Pri-NLICE, and Pri-SLJA is 364 assigned with 4 options, including one correct and three incorrect options. The detailed calculation 365 of MC1 and MC2 is outlined in Appendix J. We also evaluate the model's generation ability using 366 **ROUGE-L** and **LLM-Judge** (**LLM-J**)(Zheng et al., 2023a). For ROUGE-L, the reference text is the 367 initial output produced by the backbone LLM. As ROUGE-L primarily focuses on n-gram overlap 368 between generated text and reference texts, which may not fully capture the semantic meaning or 369 overall quality of the generated content, we further use a LLM (i.e., GPT-3.5) to assess the quality of 370 outputs considering relevance, clarity, and accuracy. The assessment prompt is shown in Appendix 371 O.3. The LLM-J score ranges from 1 to 10, with higher scores indicating better quality. For inference 372 efficiency, we use Throughput (TP), defined as the number of tokens generated per second, to 373 evaluate the inference efficiency.

Compared Methods. To demonstrate the effectiveness of our method, we compare our model with following baselines. d_{χ} -privacy. The client applies d_{χ} -privacy (Feyisetan et al., 2020) mechanism on the entire input, which injects noise into tokens' embedding and find the nearest tokens to replace the initial tokens. d_{χ} -privacy on privacy spans. The client employs d_{χ} -privacy (Feyisetan et al., 2020) mechanism on privacy spans in the input, rather than the entire input. **Paraphrase.** According to

Datasets	Methods	MC1 ↑	MC2 ↑	ROUGE-L \uparrow	LLM-J↑	$\mathrm{TP}\uparrow$
	d_{χ} -privacy	28.79±0.02	30.26±0.01	17.97±0.00	$1.17{\pm}0.00$	37.45±0.01
Pri-DDXPlus	d_{χ} -privacy on privacy spans	44.71±0.29	$42.36{\pm}0.00$	$29.17{\pm}0.04$	$3.31{\pm}0.00$	$33.21 {\pm} 0.00$
	Paraphrase	27.92 ± 0.56	$28.56{\pm}0.07$	$18.04{\pm}0.01$	$1.23{\pm}0.00$	$35.42 {\pm} 0.67$
	PrivacyRestore	62.97±0.00	$60.19{\pm}0.00$	$27.24{\pm}0.26$	$\textbf{4.47}{\pm 0.00}$	$26.09{\pm}0.08$
	d_{γ} -privacy	29.08±0.00	29.72±0.00	15.68±0.02	$1.41{\pm}0.00$	38.30±0.00
D . MI ICE	d_{χ} -privacy on privacy spans	30.00±0.09	$31.46 {\pm} 0.00$	$22.97 {\pm} 0.00$	$3.01 {\pm} 0.00$	35.73±0.57
Pri-NLICE	Paraphrase	$28.46 {\pm} 0.02$	$29.15 {\pm} 0.03$	$16.15 {\pm} 0.01$	$1.62{\pm}0.00$	$37.22 {\pm} 0.07$
	PrivacyRestore	62.23±1.70	$\textbf{57.94}{\pm}\textbf{0.09}$	$24.42{\pm}0.81$	$3.67{\pm}0.01$	$32.33{\pm}0.01$
	d_{γ} -privacy	16.66±0.37	17.57±0.04	23.35±0.00	$2.08 {\pm} 0.00$	36.83±0.03
D.: CLIA	d_{χ} -privacy on privacy spans	24.23±1.69	$26.63 {\pm} 0.67$	$40.10{\pm}0.00$	$4.54 {\pm} 0.00$	$36.16 {\pm} 0.00$
Pri-SLJA	Paraphrase	16.21±0.02	$17.52 {\pm} 0.02$	$24.90 {\pm} 0.01$	$2.07 {\pm} 0.01$	$31.31 {\pm} 0.05$
	PrivacyRestore	35.47±1.48	35.41±0.64	$37.56 {\pm} 0.06$	5.25±0.00	30.73±0.04

Table 1: Comparison of the performance and the inference efficiency between PrivacyRestore and other baselines across three privacy-preserving datasets. All experiments are conducted over 3 runs, with the average results and variances reported. The best results are highlighted in **bold**.

Mattern et al. (2022b); Utpala et al. (2023), clients can use generative models to paraphrase original inputs, achieving effects similar to DP.

397 **Implementation Details.** We use Llama2-chat-7b (Touvron et al., 2023) as the LLM backbone 398 on the server side, and BERT-base (Devlin et al., 2019) on the client side for weight estimation, as 399 described in Section 4.2. For fair comparison, we utilize flan-t5-base model (Chung et al., 2024) on 400 the client side for paraphrasing in the Paraphrase baseline as its model size is comparable to that of 401 BERT-base. During restoration vector training, the LLM parameters remain fixed, and we train the 402 restoration vectors for 5 epochs with a batch size of 1. The optimal number of edited heads K is 175 for Pri-DDXPlus/Pri-SLJA and 125 for Pri-NLICE. The search process is shown in Section I. To 403 evaluate the generation capabilities, we utilize GPT-3.5 to assess the generated outputs. The prompts 404 are detailed in Appendix O.2. To evaluate inference efficiency, we use the greedy search decoding 405 strategy and set the max generation length to 256 during generation. 406

407 Settings of Privacy Hyperparameters. The hyperparameters related to privacy protection strength 408 are ϵ for d_{χ} -privacy (on privacy spans) and PrivacyRestore, and τ for paraphrase. For fair comparison, 409 we ensure all methods under the same privacy budget. We show the calculation process of determining 410 values of ϵ and τ for different methods on different datasets in Appendix H. The values of ϵ and τ for 411 different privacy-preserving methods are shown in Table 5 in the Appendix.

412

392

393

394

413 6.2 MAIN RESULTS

As shown in Table 1, we evaluate the performance and inference efficiency of PrivacyRestore and other compared methods across three privacy-preserving datasets. Compared to d_{χ} -privacy and paraphrase, d_{χ} -privacy on privacy spans solely apply d_{χ} -privacy mechanism to those privacy spans and achieves higher scores in MC1/2, ROUGE-L and LLM-J. The possible reason for this is that both d_{χ} -privacy and paraphrase operate on the entire user input, instead of specific privacy spans. Injecting noise into the entire input creates larger disturbances during inference compared to only corrupting a limited number of privacy spans.

421 PrivacyRestore achieves best scores in MC1/2 and LLM-J compared to other privacy-preserving 422 methods. In terms of the ROUGE-L evaluation metric, PrivacyRestore achieve the best result in 423 Pri-NLICE while ranking second in the other two datasets. This discrepancy likely stems from ROUGE-L's dependence on n-gram overlap between the reference text and the generated output, 424 which does not fully reflect the quality of generated outputs. As demonstrated by the examples in 425 Figure 6 and Appendix N, PrivacyRestore often generates outputs with different sentence structures 426 while still providing accurate answers. Consequently, our method achieves slightly lower ROUGE-L 427 scores but significantly higher LLM-J scores compared to d_{χ} -privacy on privacy spans. Furthermore, 428 the ROUGE-L metric displays larger variance than the LLM-J metric, potentially due to its sensitivity 429 to expression rather than the underlying meaning of the generated output. 430

As for inference efficiency, d_{χ} -privacy achieves the highest throughput. In contrast, d_{χ} -privacy on privacy spans requires prior identification of privacy spans, while paraphrase necessitates rephrasing

432 the user input on the client side, leading to delays. PrivacyRestore also requires additional time for 433 the prior identification and removal of privacy spans, along with constructing the meta vector on the 434 client side. However, its throughput can reach nearly 70% on Pri-DDXPlus and 80% on Pri-NLICE 435 and Pri-SLJA, relative to the best results.

436 437

438 439

441

443

446

447

464

465 466

467

468

469

470

471

472

473

474

EMPIRICAL PRIVACY PROTECTION RESULTS 6.3

We not only provide a theoretical privacy proof of our method in Section 5, but also implement two attack methods to empirically evaluate the privacy protection capability of our approach. PrivacyRe-440 store sends only the meta vector and the incomplete input with privacy spans removed. It is less likely for adversaries to infer privacy spans from incomplete input. Adversaries can only attack by 442 intercepting the meta vector and inferring the corresponding privacy spans from it. Therefore, we implement the embedding inverse attack (Li et al., 2023b; Morris et al., 2023) and attribute inference 444 attack (Li et al., 2022), both commonly used methods for attacking embeddings. For d_{γ} -privacy, 445 d_{γ} -privacy and paraphrase methods, we obtain the hidden state of the last token in the last layer in Llama2-chat-7b as the embeddings. The embedding inverse attack utilizes ROUGE-L as the evaluation metric, while the attribute inference attack employs F1 to assess attack performance. 448 Detailed description about both attack methods and evaluation metrics can be found in Appendix M. 449

450 **Different Privacy Hyperparameter** ϵ . We 451 evaluate the attack performance of the embed-452 ding inverse attack and attribute inference attack 453 for PrivacyRestore and other baselines. Addi-454 tionally, we present the attack results without 455 any privacy protection, which serves as the upper bound of attack performance. The values of 456 ϵ on x-axis in Figure 2 represent the values used 457 in PrivacyRestore. Equivalent values of ϵ and τ 458 for other baselines are provided in the Appendix 459 H to ensure the same privacy budget. As shown 460 in Figure 2, although the privacy budget is the 461 same, PrivacyRestore demonstrates better pri-462 vacy protection performance, evidenced by its 463



(a) Embedding Inverse At-(b) Attribute Inference Attack tack

Figure 2: Results of all methods under embedding inverse attack and attribute inference attack under different privacy hyperparameters ϵ on three datasets.

lower ROUGE-L and F1 scores. All privacy-preserving methods effectively protect privacy compared to the upper bound of no protection.





475 Figure 3: (a) and (b) present the results of d_{χ} -privacy method under the prompt injection attack 476 and attribute inference attack under varying d_{χ} -privacy percentages across three privacy-preserving 477 datasets. (c) and (d) show the results of PrivacyRestore for the embedding inverse attack and attribute 478 inference attack under different privacy span ratios α on the same three datasets. 479

480 **Different** d_{χ} -privacy Percentage for d_{χ} -privacy. We randomly select a proportion of tokens in 481 user input to protect, denoted as the d_{χ} -privacy percentage. A larger d_{χ} -privacy percentage indicates 482 a larger number of tokens being protected and a longer protected text. As illustrated in Figure 3(a) and 3(b), both prompt injection attack and attribute inference attack exhibit increased attack performance 483 with larger d_{γ} -privacy percentages, as reflected by the increase of ROUGE-L and F1 scores. These 484 experimental results demonstrate the linear growth problem of the privacy budget in d_{γ} -privacy, as 485 proved in Section 5.1. Implementation details of attacks are presented in Appendix L.1.

486 **Different Privacy Span Ratio** α for **PrivacyRestore.** We randomly select a proportion of privacy 487 spans to protect, denoted as α , where a larger α indicates a larger number of privacy spans being 488 protected and longer protected text. As shown in Figure 3(c) and 3(d), the ROUGE-L scores of 489 embedding inverse attack are stable across different α values on three datasets. Similarly, the F1 490 scores of attribute inference attack also keep stable. The length of the protected text does not impact the privacy protection capability of PrivacyRestore. The stable performance in both attack settings 491 provides empirical evidence that our method inherently addresses the linear growth problem of 492 privacy budget in these DP variants, as proved in Section 5.2. The trend exhibits slight fluctuations on 493 both the Pri-NLICE and Pri-DDXPlus datasets. A detailed analysis of these fluctuations is provided 494 in Appendix L.2. 495

496 497 498

499 500 501

509

510

6.4 ABLATION STUDIES

Datasets	Methods	MC1↑	$\text{MC2} \uparrow$	ROUGE-L↑	LLM-J↑	$TP\uparrow$
Pri-DDXPlus	Equal Weighted Aggregation Attention-aware Weighted Aggregation	53.84 62.97	51.12 60.19	26.32 27.24	4.29 4.47	26.35 26.09
Pri-NLICE	Equal Weighted Aggregation	46.92	45.89	22.78	3.12	32.75
	Attention-aware Weighted Aggregation	62.23	57.94	24.42	3.67	32.33
Pri-SLJA	Equal Weighted Aggregation	30.88	30.70	30.96	4.10	31.00
	Attention-aware Weighted Aggregation	35.47	35.41	37.56	5.25	30.73

Table 2: Comparison of the performance and the inference efficiency between Equal Weighted Aggregation and Attention-aware Weighted Aggregation. The best results are highlighted in **bold**.

In order to verify the effectiveness of Attention-aware Weighted Aggregation (AWA) component,
 we compare the performance and the inference efficiency between equal weighted aggregation and
 attention-aware weighted aggregation. Different from attention-aware weighted aggregation, equal
 weighted aggregation computes the meta vector by simply summing up all restoration vectors.

As shown in Table 2, the MC1, MC2, ROUGE-L, and LLM-J scores of equal weighted aggregation are all lower than those of attention-aware weighted aggregation, indicating that simply summing all restoration vectors equally degrades performance. This degradation is primarily due to the equal weights diluting the influence of critical spans while amplifying the effect of irrelevant ones. In terms of inference efficiency, the throughput difference between Attention-Aware Weighted Aggregation and Equal Weighted Aggregation is negligible. This suggests that the weight computation, as defined in Eq 4, is efficient and does not significantly impact overall throughput.

522 523

524 525

526

6.5 ANALYSIS OF HYPERPARAMETER AND LLM BACKBONE

We analyze the performance of PrivacyRestore using different numbers of edited heads K. In addition, we analyze the performance of PrivacyRestore using a different LLM backbone (i.e., Llama-13b-chat). Due to space limitation, we put the analysis in the Appendix I and Appendix K.

527 528 529

7 CONCLUSION

530

531 We propose PrivacyRestore which protects the privacy within user inputs during inference in online 532 LLM inference services. PrivacyRestore achieves privacy protection by directly removing privacy 533 spans in the user input and then restoring these privacy spans via activation steering. PrivacyRestore 534 provides a practical and efficient solution for protecting privacy while maintaining satisfactory performance and inference efficiency. We demonstrate that PrivacyRestore inherently addresses 535 the linear growth problem of the privacy budget found in differential privacy variants. We curate 536 three privacy-preserving datasets covering medical and legal fields, and PrivacyRestore achieves 537 strong performance and inference efficiency across all datasets. Additionally, we implemented two 538 types of attacks, and the experimental results demonstrate PrivacyRestore's robust privacy protection capabilities.

540 REFERENCES

554

- Zaid Al-Ars, Obinna Agba, Zhuoran Guo, Christiaan Boerkamp, Ziyaad Jaber, and Tareq Jaber. Nlice:
 Synthetic medical record generation for effective primary healthcare differential diagnosis. In 2023 *IEEE 23rd International Conference on Bioinformatics and Bioengineering (BIBE)*, pp. 397–402.
 IEEE, 2023.
- Guillaume Alain and Yoshua Bengio. Understanding intermediate layers using linear classifier probes.
 arXiv preprint arXiv:1610.01644, 2016.
- Mário S. Alvim, Konstantinos Chatzikokolakis, Catuscia Palamidessi, and Anna Pazii. Invited paper: Local differential privacy on metric spaces: Optimizing the trade-off with utility. In *31st IEEE Computer Security Foundations Symposium, CSF 2018, Oxford, United Kingdom, July 9-12, 2018*, pp. 262–267, 2018. doi: 10.1109/CSF.2018.00026. URL https://doi.org/10.1109/ CSF.2018.00026.
 - Yonatan Belinkov. Probing classifiers: Promises, shortcomings, and advances. *Computational Linguistics*, 48(1):207–219, 2022.
- Haohan Bo, Steven H. H. Ding, Benjamin C. M. Fung, and Farkhund Iqbal. ER-AE: differentially private text generation for authorship anonymization. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2021, Online, June 6-11, 2021*, pp. 3997–4007, 2021. doi: 10.18653/V1/2021.NAACL-MAIN.314. URL https://doi.org/10.18653/v1/2021.
 naacl-main.314.
- Konstantinos Chatzikokolakis, Miguel E. Andrés, Nicolás Emilio Bordenabe, and Catuscia Palamidessi. Broadening the scope of differential privacy using metrics. In *Privacy Enhancing Technologies - 13th International Symposium, PETS 2013, Bloomington, IN, USA, July 10-12, 2013. Proceedings*, pp. 82–102, 2013. doi: 10.1007/978-3-642-39077-7_5. URL https://doi.org/10.1007/978-3-642-39077-7_5.
- Zeming Chen, Alejandro Hernández Cano, Angelika Romanou, Antoine Bonnet, Kyle Matoba,
 Francesco Salvi, Matteo Pagliardini, Simin Fan, Andreas Köpf, Amirkeivan Mohtashami,
 et al. Meditron-70b: Scaling medical pretraining for large language models. *arXiv preprint arXiv:2311.16079*, 2023.
- Zhongzhi Chen, Xingwu Sun, Xianfeng Jiao, Fengzong Lian, Zhanhui Kang, Di Wang, and
 Chengzhong Xu. Truth forest: Toward multi-scale truthfulness in large language models through
 intervention without tuning. In *Proceedings of the AAAI Conference on Artificial Intelligence*,
 volume 38, pp. 20967–20974, 2024.
- 576 Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan 577 Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Alex Castro-Ros, Marie Pellat, 578 Kevin Robinson, Dasha Valter, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Y. Zhao, 579 Yanping Huang, Andrew M. Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, 580 Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. Scaling instruction-finetuned language 581 models. J. Mach. Learn. Res., 25:70:1-70:53, 2024. URL https://jmlr.org/papers/ 582 v25/23-0870.html. 583
- Wentao Deng, Jiahuan Pei, Keyi Kong, Zhe Chen, Furu Wei, Yujun Li, Zhaochun Ren, Zhumin
 Chen, and Pengjie Ren. Syllogistic reasoning for legal judgment analysis. In Houda Bouamor,
 Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Conference on Empirical Meth- ods in Natural Language Processing*, pp. 13997–14009, Singapore, December 2023a. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.864. URL https:
 //aclanthology.org/2023.emnlp-main.864.
- Wentao Deng, Jiahuan Pei, Keyi Kong, Zhe Chen, Furu Wei, Yujun Li, Zhaochun Ren, Zhumin Chen, and Pengjie Ren. Syllogistic reasoning for legal judgment analysis. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pp. 13997–14009, 2023b. doi: 10.18653/V1/2023.EMNLP-MAIN.864. URL https://doi.org/10.18653/v1/2023.emnlp-main.864.

594	Jacob Devlin Ming-Wei Chang Kenton Lee and Kristina Toutanova BERT. Pre-training of
595	deep bidirectional transformers for language understanding. In Jill Burstein, Christy Doran, and
596	Thamar Solorio (eds.), Proceedings of the 2019 Conference of the North American Chapter of the
597	Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and
598	Short Papers), pp. 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational
599	Linguistics. doi: 10.18653/v1/N19-1423. URL https://aclanthology.org/N19-1423.
600	
601	John C. Duchi, Michael I. Jordan, and Martin J. Wainwright. Local privacy and statistical minimax
602	rates. In 51st Annual Allerton Conference on Communication, Control, and Computing, Allerton
603	2013, Allerton Park & Retreat Center, Monticello, IL, USA, October 2-4, 2013, pp. 1592, 2013.
604	dol: 10.1109/ALLEKION.2013.0/30/18. UKL https://dol.org/10.1109/Allerton.
605	2013.0730710.
606	Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam D. Smith, Calibrating noise to sensitivity
607	in private data analysis. J. Priv. Confidentiality, 7(3):17–51, 2016. doi: 10.29012/JPC.V7I3.405.
608	URL https://doi.org/10.29012/jpc.v7i3.405.
609	
610	Oluwaseyi Feyisetan, Tom Diethe, and Thomas Drake. Leveraging hierarchical representations for
611	preserving privacy and utility in text. In 2019 IEEE International Conference on Data Mining,
612	<i>ICDM 2019, Beijing, China, November 8-11, 2019</i> , pp. 210–219, 2019. doi: 10.1109/ICDM.2019.
613	00031. URL https://doi.org/10.1109/ICDM.2019.00031.
614	
615	Oluwaseyi Feyisetan, Borja Balle, Thomas Drake, and Tom Dietne. Privacy-and utility-preserving
616	conference on web search and data mining pp. 178–186–2020
617	conjerence on web search and data mining, pp. 178–180, 2020.
618	Kanay Gupta, Neha Jawalkar, Ananta Mukheriee, Nishanth Chandran, Divya Gupta, Ashish Panwar,
619	and Rahul Sharma. Sigma: Secure gpt inference with function secret sharing. Cryptology
620	ePrint Archive, Paper 2023/1269, 2023. URL https://eprint.iacr.org/2023/1269.
621	https://eprint.iacr.org/2023/1269.
622	
623	Meng Hao, Hongwei Li, Hanxiao Chen, Pengzhi Xing, Guowen Xu, and Tianwei Zhang. Iron: Private
624	inference on transformers. Advances in neural information processing systems, 35:15718–15731,
625	2022a.
626	Mang Hao Hongwai Li Hanvigo Chan Dangzhi Ving Guowan Vu and Tianwai
627	Zhang Iron: Private inference on transformers. In S Koveio S Mohamed
628	A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.). Advances in Neural Infor-
620	mation Processing Systems, volume 35, pp. 15718–15731, Curran Associates, Inc.
620	2022b. URL https://proceedings.neurips.cc/paper files/paper/2022/
621	file/64e2449d74f84e5b1a5c96ba7b3d308e-Paper-Conference.pdf.
031	
622	Evan Hernandez, Belinda Z. Li, and Jacob Andreas. Inspecting and editing knowledge representations
033	in language models, 2023.
034	lines Hang Nach Les and James Thomas OPPO: manalithic anti-
635	Jiwoo Hong, Noan Lee, and James Thorne. ORPO: monolithic preference optimization without
636	reference model. $COKK$, abs/2403.0/091, 2024. doi: 10.48550/AKAIV.2403.0/091. UKL
637	https://doi.org/10.40550/arxiv.2405.07091.
638	Dacheng Li, Hongyi Wang, Rulin Shao, Han Guo, Eric Xing, and Hao Zhang, MPCFORMER: FAST,
039	PERFORMANT AND PRIVATE TRANSFORMER INFERENCE WITH MPC. In <i>The Eleventh</i>
040	International Conference on Learning Representations, 2023a. URL https://openreview.
641	net/forum?id=CWmvjOEhgH
642	
643	Haoran Li, Yangqiu Song, and Lixin Fan. You don't know my favorite color: Preventing dialogue
644	representations from revealing speakers' private personas. In <i>Proceedings of the 2022 Conference</i>
645	of the North American Chapter of the Association for Computational Linguistics: Human Language
646	<i>Iechnologies, NAACL 2022, Seattle, WA, United States, July 10-15, 2022, pp. 5858–5870, 2022.</i> doi: 10.18652/01/2022. NAACL MAIN 420. LIPL https://doi.org/10.10652/01/2020
647	10.1003/01/2022.NAACL-101AIN.429. UKL NULPS://doi.org/10.10053/V1/2022.

640

653

667

669

673

040	Haoran Li, Mingshi Xu, and Yangqiu Song. Sentence embedding leaks more information than you
649	expect: Generative embedding inversion attack to recover the whole sentence. In Findings of
650	the Association for Computational Linguistics: ACL 2023, Toronto, Canada, July 9-14, 2023,
651	pp. 14022-14040, 2023b. doi: 10.18653/V1/2023.FINDINGS-ACL.881. URL https://doi.
652	org/10.18653/v1/2023.findings-acl.881.

- Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. Inference-654 time intervention: Eliciting truthful answers from a language model. In A. Oh, T. Nau-655 mann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), Advances in Neural 656 Information Processing Systems, volume 36, pp. 41451-41530. Curran Associates, Inc., 657 2023c. URL https://proceedings.neurips.cc/paper_files/paper/2023/ 658 file/81b8390039b7302c909cb769f8b6cd93-Paper-Conference.pdf. 659
- Yansong Li, Zhixing Tan, and Yang Liu. Privacy-preserving prompt tuning for large language 660 model services. CoRR, abs/2305.06212, 2023d. doi: 10.48550/ARXIV.2305.06212. URL 661 https://doi.org/10.48550/arXiv.2305.06212. 662
- 663 Yu Li, Zhihua Wei, Han Jiang, and Chuanyang Gong. DESTEIN: navigating detoxification of lan-664 guage models via universal steering pairs and head-wise activation fusion. CoRR, abs/2404.10464, 665 2024. doi: 10.48550/ARXIV.2404.10464. URL https://doi.org/10.48550/arXiv. 666 2404.10464.
- Zi Liang, Pinghui Wang, Ruofei Zhang, Nuo Xu, Shuo Zhang, Lifeng Xing, Haitao Bai, and Ziyang 668 Zhou. Merge: Fast private text generation. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 38, pp. 19884–19892, 2024. 670
- 671 Stephanie Lin, Jacob Hilton, and Owain Evans. Truthfulqa: Measuring how models mimic human 672 falsehoods. arXiv preprint arXiv:2109.07958, 2021.
- Xuanqi Liu and Zhuotao Liu. Llms can understand encrypted prompt: Towards privacy-computing 674 friendly transformers. arXiv preprint arXiv:2305.18396, 2023. 675
- 676 Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike 677 Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining 678 approach. arXiv preprint arXiv:1907.11692, 2019.
- Justus Mattern, Zhijing Jin, Benjamin Weggenmann, Bernhard Schoelkopf, and Mrinmaya Sachan. 680 Differentially private language models for secure data sharing. In Yoav Goldberg, Zornitsa 681 Kozareva, and Yue Zhang (eds.), Proceedings of the 2022 Conference on Empirical Meth-682 ods in Natural Language Processing, Abu Dhabi, United Arab Emirates, December 2022a. 683 Association for Computational Linguistics. doi: 10.18653/v1/2022.emnlp-main.323. URL 684 https://aclanthology.org/2022.emnlp-main.323. 685
- Justus Mattern, Benjamin Weggenmann, and Florian Kerschbaum. The limits of word level differential 686 privacy. In Findings of the Association for Computational Linguistics: NAACL 2022, Seattle, WA, 687 United States, July 10-15, 2022, pp. 867–881, 2022b. doi: 10.18653/V1/2022.FINDINGS-NAACL. 688 65. URL https://doi.org/10.18653/v1/2022.findings-naacl.65. 689
- 690 John X. Morris, Volodymyr Kuleshov, Vitaly Shmatikov, and Alexander M. Rush. Text embeddings 691 reveal (almost) as much as text. In Proceedings of the 2023 Conference on Empirical Methods in 692 Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023, pp. 12448–12460, 693 2023. doi: 10.18653/V1/2023.EMNLP-MAIN.765. URL https://doi.org/10.18653/ v1/2023.emnlp-main.765. 694
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong 696 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow 697 instructions with human feedback. Advances in neural information processing systems, 35:27730– 27744, 2022. 699
- Fábio Perez and Ian Ribeiro. Ignore previous prompt: Attack techniques for language models. CoRR, 700 abs/2211.09527, 2022. doi: 10.48550/ARXIV.2211.09527. URL https://doi.org/10. 48550/arXiv.2211.09527.

702	Chen Qu, Weize Kong, Liu Yang, Mingyang Zhang, Michael Bendersky, and Marc Najork. Natural
703	language understanding with privacy-preserving BERT. In CIKM '21: The 30th ACM International
704	Conference on Information and Knowledge Management, Virtual Event, Queensland, Australia,
705	November 1 - 5, 2021, pp. 1488-1497, 2021. doi: 10.1145/3459637.3482281. URL https:
706	//doi.org/10.1145/3459637.3482281.
707	
708	Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language
709	models are unsupervised multitask learners. 2019.
710	Vishen See, Signal answert, A new anneath to answert answert injustice attacks and interview the
711	integrated applications. CoPP, abs/2401.07612, 2024. doi: 10.48550/APXIV.2401.07612. LIPL
712 713	https://doi.org/10.48550/arXiv.2401.07612.
714	$\mathbf{A} = \mathbf{A} + $
715	plus: A new dataset for automatic medical diagnosis. In Sanmi Koyejo, S. Mohamed,
716	A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh (eds.), Advances in Neural Infor-
717	mation Processing Systems 35: Annual Conference on Neural Information Process-
718	ing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December
719	9, 2022, 2022. UKL http://papers.nips.cc/paper_files/paper/2022/
720	Hash/Cae/Say/4590C0edd95ae/aeae09159C-Abstract-Datasets_and_ Benchmarks_html
721	Denchmarks.nemt.
722	Ian Tenney, Dipanian Das, and Ellie Paylick. Bert rediscovers the classical nlp pipeline. arXiv
723	preprint arXiv:1905.05950, 2019.
724	
725	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay
726	Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation
727	and fine-tuned chat models. arXiv preprint arXiv:2307.09288, 2023.
728	
729	Alexander Matt Turner, Lisa Thiergart, David Udell, Gavin Leech, Ulisse Mini, and Monte MacDi-
730	armid. Activation addition: Steering language models without optimization, 2023.
731	Saitaia Utnala, Sara Hooker, and Din Vu Chen. Locally differentially private document gen
732	eration using zero shot prompting. In Findings of the Association for Computational Lin-
733	<i>guistics: EMNLP 2023, Singapore, December 6-10, 2023, pp. 8442–8457, 2023, doi: 10.</i>
734	18653/V1/2023.FINDINGS-EMNLP.566. URL https://doi.org/10.18653/v1/2023.
735	findings-emnlp.566.
736	
737	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
738	Kaiser, and Illia Polosukhin. Attention is all you need. Advances in neural information processing
739	systems, 30, 2017.
740	Chille We Open Incore Otenen Les Verlies Debrauel 11 Med Decides Other Colors De 11
741	ian Kambadur, David Rosenberg, and Cideon Mann. Bloomberganti, A large language model for
742	finance arXiv preprint arXiv:2303 17564 2023
743	тинов. и ли роргии и ли. 2505.17507, 2025.
744	Yiquan Wu, Yuhang Liu, Yifei Liu, Ang Li, Siying Zhou, and Kun Kuang, wisdominterrogatory URL
745	https://github.com/zhihaiLLM/wisdomInterrogatory. Available at GitHub.
746	
/4/	Qianqian Xie, Weiguang Han, Xiao Zhang, Yanzhao Lai, Min Peng, Alejandro Lopez-Lira, and Jimin
748	Huang. Pixiu: A large language model, instruction data and evaluation benchmark for finance.
749	arXiv preprint arXiv:2306.05443, 2023.
/50	
/51	Canwen Xu, Daya Guo, Nan Duan, and Julian J. McAuley. Baize: An open-source chat model with
/52	parameter-efficient tuning on self-chat data. In <i>Proceedings of EMNLP</i> , pp. 6268–6278, 2023.
/53	Zakun Xu, Abbinay Aggamual Olympaayi Equicaton and Nathangal Triggian A differentially arisest
754 755	text perturbation method using a regularized mahalanobis metric. <i>CoRR</i> , abs/2010.11947, 2020. URL https://arxiv.org/abs/2010.11947.

756 757 758 759 760 761 762 763	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. Judging llm-as-a-judge with mt-bench and chatbot arena. In Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023, 2023a. URL http://papers.nips.cc/paper_files/paper/2023/hash/91f18a1287b398d378ef22505bf41832-Abstract-Datasets_and_Benchmarks.html.
764 765	Mengxin Zheng, Qian Lou, and Lei Jiang. Primer: Fast private transformer inference on encrypted data. In 2023 60th ACM/IEEE Design Automation Conference (DAC), pp. 1–6. IEEE, 2023b.
766	
767	
768	
769	
770	
771	
772	
774	
775	
776	
777	
778	
779	
780	
781	
782	
783	
784	
785	
786	
787	
788	
789	
790	
791	
792	
793	
795	
796	
797	
798	
799	
800	
801	
802	
803	
804	
805	
806	
807	
808	
009	

⁸¹⁰ A NOTATIONS

Here we present all notations used in our paper in Table 3.

813 814 815

816 817

818

812

B DATASETS

B.1 CONSTRUCTION PROCESS

We used GPT-3.5 (Ouyang et al., 2022) to classify symptoms in DDXPlus and NLICE, as well as case details in SLJA, into five levels ranging from non-sensitive to highly sensitive. The assessment prompt template is shown in Appendix O.1. A higher level indicates that the symptom or case detail is more sensitive. We define all symptoms and case details with a sensitivity level greater than 3 as privacy spans.

We assign each sample a correct answer along with three randomly selected incorrect options. For DDXPlus and NLICE, we randomly select three diagnosis results to combine with the correct diagnosis as the choices. In the SLJA dataset, we randomly select three legal judgments to pair with the correct one as the options.

The initial dataset is extensive, and we observed that for most samples, removing all privacy spans 829 often yields outputs similar to those obtained when privacy spans are provided. Privacy preserving 830 for these samples is meaningless because users can directly hide those privacy spans and obtain 831 approximate result outputs. In real-world scenarios, sensitive privacy spans often play a crucial role in 832 medical diagnoses and legal judgments, making privacy preservation highly valuable. Our dataset is 833 designed to benchmark various privacy-preserving methods and must include samples where privacy 834 spans are crucial for generating outputs. We utilize the KL divergences to measure the importance 835 scores of samples. We calculate the KL divergence between the model output distributions with 836 and without the privacy symptoms included. A higher KL divergence indicates that the absence of sensitive privacy spans may lead to different or incorrect outputs. We selected only samples 837 with high KL divergence to construct the privacy-preserving datasets. As a result, we curated three 838 privacy-preserving datasets: Pri-DDXPlus and Pri-NLICE for medical diagnosis, and Pri-SLJA for 839 legal judgment. 840

841 842

843

B.2 STATISTICAL INFORMATION

We show the statistics of the obtained Pri-DDXPlus, Pri-NLICE and Pri-SLJA datasets in Table 4. We tally the number of user inputs, privacy span types, and answer types. We also compute the average occurrence of privacy spans per instance. In Pri-DDXPlus and Pri-NLICE, the privacy spans are the symptoms, and the answers are the diagnoses. In Pri-SLJA, the privacy spans are the case details, and the answers are the legal judgments.

Pri-DDXPlus commonly contains more instances and more privacy symptoms types compared to
Pri-NLICE and Pri-SLJA. Each sample in Pri-DDXPlus contains an average of six privacy symptoms,
while samples in Pri-NLICE have an average of four privacy spans, and samples in Pri-SLJA have an
average of three privacy spans.

853 854

855 856

857

C PRELIMINARIES FOR METHODOLOGY

The d_{χ} -privacy protection and activation steering technique are two crucial components of our method. Here, we provide an illustration of these techniques for better understanding of our method.

858 859

C.1 d_{χ} -privacy protection

860 861

 d_{χ} -privacy (Feyisetan et al., 2019) is a variant of the differential privacy mechanism designed to protect privacy by incorporating a distance measure into the privacy budget. Its detailed definition is provided in Definition 2.1. Typically, to implement the d_{χ} -privacy mechanism, noise is added to the

Notations	Definitions
s	A single privacy span.
S	All possible privacy spans.
\mathcal{S}_I	All privacy spans in user input <i>I</i> .
h	A single edited head.
\mathcal{H}_{c}	The common top-K heads set.
\mathcal{H}_{a}	The set of all heads.
\mathcal{H}_{1}^{s}	The top-K heads set of the privacy span s.
L_{h}	The score list of the head h across all privacy spans.
K	The number of selected edited heads
\mathcal{F}^s	The probe of privacy span s on head h
\mathcal{J}_{h} \mathcal{A}^{s}	The prope of privacy span s on read n .
σ_h	The sigmoid function
0	The signified function.
\mathbf{u}_h	The output hidden state offer restoration on head h
\mathbf{u}_h	The output induction state after restoration on head n .
r_s^n	The restoration vector for privacy span s on head h .
Θ	All restoration vectors for all privacy spans on all edited heads.
λ	The tradeoff hyperparameter of ORPO loss.
w_s	The weight of privacy span s.
n	The number of tokens in the user input.
n_h	The number of heads in the lightweight model.
$\operatorname{Attn}_h(x,y)$	The attention score of y attending to x on head h .
Z_h	The normalized weighted sum of restoration vectors on head h.
Z_h, Z'_h	Any two normalized weighted sums.
\mathcal{R}_h	The meta vector on head h.
\mathcal{N}^{n}	The added noise on the normalized weighted sums for meta vector construction
$ x _{2}$	The l_2 -normalization of x
$I_{all} = \{I_1,, \}$	I_m All user inputs in the training set.
$Y_{2} = \{y_{1}, y_{2}\}$	The labels indicating whether the corresponding input contains s
$m = (g_1,, g_m)$	The size of training set
T T'	Any two user inputs
$I = \int i_{-} i_{-}$	The tokens of the input I
$I = \iota_1,, \iota_n$	The tokens of the user input I .
l_t	The token on the ddings of the input I .
$\{e_1,, e_n\}$	The rescible system tests for <i>I</i> .
	The possible output sets for I .
$O = \{o_1,, o_n\}$	I he possible output sets for tokens of the input I , also represented as O .
1	The user input with all privacy spans removed.
$\hat{I}_{all} = \{\hat{I}_1,, \}$	\hat{I}_m All user inputs with privacy spans removed in the training set.
a	The initial output given the complete input <i>I</i> .
\hat{a}	The output given the incomplete input with privacy spans removed \hat{I}
\overline{O} O'	Any two queries to the database
$O = \{a_1, \dots, a_n\}$	The sub-queries of the query O
$\mathcal{G} = \{q_1, \dots, q_n\}$	The possible query result sets for O
0	The possible result sets for sub-queries of the query Ω also represented
$G = \{g_1,, g_n\}$	$\frac{1}{1}$ $\frac{1}{100}$ $\frac{1}{1$
	as tr.
e	The privacy hyperparameter.
τ	i ne generation temperature.
0	The privacy hyperparameter.
n_{ps}	The number of tokens associated with the privacy spans in the input.
α	The proportion of privacy spans selected for protection.
d_{χ}	Any distance function used by d_{χ} -privacy.
d_e	The prior value of distance between token embeddings.
d_z	The prior value of distance between normalized weighted sums.
	Table 3: Definitions of all notations used in our paper.

Datasets	Dataset Split	User inputs	Privacy Span Type	Avg. Privacy Spans
	All	7759	149	5.95
	Train	5901	149	6.03
FII-DDAFius	Dev	309	60	5.37
	Test	1549	78	5.77
	All	4062	64	4.49
D. NI ICE	Train	3282	64	4.55
PII-NLICE	Dev	130	58	4.25
	Test	650	64	4.24
	All	3901	142	2.67
	Train	3117	142	2.56
PII-SLJA	Dev	130	95	3.21
	Test	654	142	3.09

Table 4: The statistics of Pri-DDXPlus and Pri-NLICE. Average privacy symptoms indicate the average privacy spans occur in one query.

initial embedding or vector for privacy protection, as follows:

$$\mathcal{R} = Z + \mathcal{N},\tag{8}$$

$$\mathbb{P}(\mathcal{N}) \propto \exp(-\epsilon ||\mathcal{N}||),$$
(9)

where Z is the protected embedding/vector, N is the added noise, R is the protected results and ϵ is the privacy parameter of the mechanism. According to Feyisetan et al. (2019), in order to sampling the noise N from the distribution, we can compute as the following:

$$\mathbf{v} \in \{ v \in \mathbb{R}^n : ||v|| = 1 \}$$

$$\tag{10}$$

$$(\mathbf{l}) \propto \frac{\mathbf{l}^{n-1}e^{-\epsilon \mathbf{l}}}{\Gamma(n)\epsilon^{-n}},\tag{11}$$

$$\mathcal{N} = \mathbf{l} \cdot \mathbf{v}, \tag{12}$$

where n is the size of the embedding/vector and ϵ is the privacy parameter.

 \mathbb{P}

C.2 ACTIVATION STEERING TECHNIQUE

Activation steering methods (Li et al., 2023c; Turner et al., 2023; Hernandez et al., 2023) control the behavior of LLM by modifying their activations during the inference stage, without incurring training costs. It serves as a crucial part of our methodology to restore privacy information during LLM inference. Typically, the attention mechanism (Vaswani et al., 2017) in LLM is responsible for capturing contextual information, and it can be expressed as:

$$q = W_q \cdot \mathbf{i}, \tag{13}$$

$$\mathbf{u} = \operatorname{Softmax}(\frac{q \cdot K^T}{\sqrt{d_k}}) \cdot V, \tag{14}$$

where **i** is the input hidden state, **u** is the output hidden state, W_q is the query weight matrix, K is the key of the context and V is the value of the context and d_k is the dimension of the key. Activation steering methods add some steering vectors into the output hidden state and in our methos we add the restoration meta vector into the output hidden state to restore privacy information, which can be expressed as:

$$\mathbf{u} = \mathbf{u} + \mathcal{R},\tag{15}$$

where \mathcal{R} is the steering/restoration meta vector.

D FORMAL DEFINITIONS OF CDP AND LDP

D.1 CENTRALIZED DIFFERENTIAL PRIVACY (CDP)

971 CDP protects individual privacy only **after data has been aggregated in a central repository**, which is defined as,

972 **Definition D.1.** (*CDP*). A randomized mechanism $\mathcal{M} : \mathcal{Q} \to \mathcal{G}$ fulfills (ϵ, δ) -differential privacy if 973 for all adjacent queries $Q, Q' \in \mathcal{Q}$ and all possible query results $G \subset \mathcal{G}$, 974

$$\mathbb{P}\left(\mathcal{M}(Q)\in G\right)\leq \exp(\epsilon)\mathbb{P}\left(\mathcal{M}(Q')\in G\right)+\delta.$$

 $\begin{array}{ll} \textbf{CDP} (Dwork et al., 2016) ensures that adversaries cannot distinguish between <math>Q$ and Q' based on Gdue to the similar probabilities, meaning the query results are probabilistically indistinguishable. This prevents adversaries from inferring characteristics about the repository based on multiple queries and their results. \\ \end{array}

981 982

987

988

989 990

991

992

993 994 995

996 997

998

999

975

D.2 LOCAL DIFFERENTIAL PRIVACY (LDP)

However, privacy risks can also emerge during the data collection process itself, as attackers may
 intercept user inputs while they are being transmitted to the central repository. LDP (Duchi et al., 2013) protect user inputs during transmission process by ensuring that attackers cannot distinguish
 between any two adjacent inputs, which is defined as,

Definition D.2. (*LDP*). A randomized mechanism $\mathcal{M} : \mathcal{I} \to \mathcal{O}$ fulfills (ϵ, δ) -differential privacy if for all adjacent inputs $I, I' \in \mathcal{I}$ and all possible outputs $O \subset \mathcal{O}$,

$$\mathbb{P}\left(\mathcal{M}(I)\in O\right)\leq \exp(\epsilon)\mathbb{P}\left(\mathcal{M}(I')\in O\right)+\delta.$$

The mechanism A processes the user input before transmitting it. LDP ensures that, even if attackers intercept O, they cannot distinguish between the initial user input I and the adjacent one I'.

E PROOF OF THEOREM 5.1

Proof of d_{χ} **-privacy**. As shown in Definition 2.1, if the input length is 1, indicating a single token i_1 , d_{χ} can be:

$$\mathbb{P}\left(\mathcal{M}(i_1) \in o_1\right) \le \exp(\epsilon d_{\chi}(i_1, i'_1)) \mathbb{P}\left(\mathcal{M}(i'_1) \in o_1\right),$$

where o_1 is the possible output set for $\mathcal{M}(i_1)$ and the privacy budget is $\epsilon d_{\chi}(i_1, i'_1)$. When the input becomes the sequential tokens $I = \{i_1, i_2, ..., i_n\}$ with corresponding output sets $O = \{o_1, ..., o_n\}$, the LCP for the sequence of length n is:

$$\mathbb{P}\left(\mathcal{M}(I)\in O\right) = \mathbb{P}\left(\mathcal{M}(i_1)\in o_1\right)\cdot\mathbb{P}\left(\mathcal{M}(i_2)\in o_2\right)\cdot\ldots\cdot\mathbb{P}\left(\mathcal{M}(i_n)\in o_n\right)$$
$$\leq \left[\exp(\epsilon d_{\chi}(i_1,i'_1))\mathbb{P}\left(\mathcal{M}(i'_1)\in o_1\right)\right]\cdot\ldots\cdot\left[\exp(\epsilon d_{\chi}(i_n,i'_n))\mathbb{P}\left(\mathcal{M}(i'_n)\in o_n\right)\right]$$

1003 1004

1006 1007 1008

1009

1017

1021

1023

$$\leq [\exp(\epsilon d_{\chi}(i_1, i'_1)) \mathbb{P} \left(\mathcal{M}(i'_1) \in o_1\right)] \cdot \ldots \cdot [\exp(\epsilon d_{\chi}(i_n, i'_n))]$$

$$= \exp[\epsilon (d_{\chi}(i_1, i'_1) + \ldots + d_{\chi}(i_n, i'_n))] \mathbb{P} \left(\mathcal{M}(I') \in O\right)$$

$$= \exp[\epsilon \sum_{j=0}^n d_{\chi}(i_j, i'_j)] \mathbb{P} \left(\mathcal{M}(I') \in O\right),$$

where the privacy budget is $\epsilon \sum_{j=1}^{n} d_{\chi}(i_j, i'_j)$. Commonly, we use the Euclidean distance as the d_{χ} function and obviously $\sum_{j=1}^{n} d_{\chi}(i_j, i'_j) \propto n$. Therefore, the privacy budget of LDP grows linearly with the length n.

Proof of LDP. As stated in Definition D.2, if the length of input is 1, corresponding to a single token i_1 , LDP can be expressed as:

$$\mathbb{P}\left(\mathcal{M}(i_1) \in o_1\right) \le \exp(\epsilon) \mathbb{P}\left(\mathcal{M}(i_1') \in o_1\right) + \delta,$$

where o_1 is the possible output set of $\mathcal{M}(i_1)$ and the privacy budget is controlled by (ϵ, δ) . Considering the sequence tokens $I = \{i_1, i_2, ..., i_n\}$ and corresponding output sets $O = \{o_1, o_2, ..., o_n\}$, the CDP for the sequence of length n can be written as:

$$\mathbb{P}\left(\mathcal{M}(I)\in O\right) = \mathbb{P}\left(\mathcal{M}(i_{1})\in o_{1}\right)\cdot\mathbb{P}\left(\mathcal{M}(i_{2})\in o_{2}\right)\cdot\ldots\cdot\mathbb{P}\left(\mathcal{M}(i_{n})\in o_{n}\right)$$
$$\leq \left[\exp(\epsilon)\mathbb{P}\left(\mathcal{M}(i_{1}')\in o_{1}\right)+\delta\right]\cdot\ldots\cdot\left[\exp(\epsilon)\mathbb{P}\left(\mathcal{M}(i_{n}')\in o_{n}\right)+\delta\right]$$

1024
1025
$$= \exp(n\epsilon)\mathbb{P}\left(\mathcal{M}(I') \in O\right) + \delta \cdot \sum_{i=1} \prod_{j=i} \mathbb{P}\left(\mathcal{M}(i'_j) \in o_j\right) + \delta^2 \cdot ...,$$

1026 where δ is typically considered a very small value. When δ approaches 0, we consider only the first 1027 two terms and then,

1028 1029 1030

1042 1043 1044

1047

1061

1062 1063

$$\mathbb{P}\left(\mathcal{M}(I)\in O\right) \leq \exp(n\epsilon)\mathbb{P}\left(\mathcal{M}(I')\in O\right) + \delta \cdot \sum_{i=1}^{n} \prod_{j!=i} \mathbb{P}\left(\mathcal{M}(i'_{j})\in o_{j}\right),$$

1031 which indicating the privacy budget becomes $(n\epsilon, \delta \cdot \sum_{i=1}^{n} \prod_{j!=i} \mathbb{P}(\mathcal{M}(i'_{j}) \in o_{j}))$, according to 1032 the Definition of LDP in Section D.2. The first term $n\epsilon$ obviously grows linearly with the length 1033 n. The second term can be view as δ multiplied by $\sum_{i=1}^{n} \prod_{j=i} \mathbb{P}\left(\mathcal{M}(i'_{j}) \in o_{j}\right)$. The second 1034 term summarizes n multiplicative terms, each bounded within (0,1). The second term can be 1035 approximately considered to grow linearly with n. Therefore, the privacy budget of CDP also grows 1036 linearly with the length n. 1037

Proof of CDP. The definition of CDP, as shown in Section D.1, is similar to LDP, with the only 1038 difference being that CDP applies to the user query $Q = \{q_1, q_2, ..., q_n\}$ rather than the text input 1039 $I = \{i_1, i_2, ..., i_n\}$. If we consider the token i_n as a sub-query q_n , then similarly the definition of 1040 CDP for sequential sub-queries $Q = \{q_1, q_2, ..., q_n\}$ can be: 1041

$$\mathbb{P}\left(\mathcal{M}(Q)\in G\right)\leq \exp(n\epsilon)\mathbb{P}\left(\mathcal{M}(Q')\in G\right)+\delta\cdot\sum_{i=1}^{n}\prod_{j!=i}\mathbb{P}\left(\mathcal{M}(q'_{j})\in g_{j}\right),$$

where $G = \{g_1, g_2, ..., g_n\}$ are the possible output sets for sequence inputs Q. The privacy budget is 1045 $(n\epsilon, \delta \cdot \sum_{i=1}^{n} \prod_{j \neq i} \mathbb{P}(\mathcal{M}(q'_{j}) \in g_{j}))$ and also grows linearly with the length n. 1046

1048 F **PROOF OF THEOREM 5.2** 1049

1050 As shown in Figure 1, during the inference stage, only the meta vector and the query with privacy 1051 spans removed are transmitted from the client to the server. The meta vector holds information 1052 about all privacy spans and could be vulnerable to interception by adversaries who may attempt to 1053 reverse-engineer these spans.

1054 PrivacyRestore protects the meta vector by adding noise \mathcal{N} which is sampling from the distribution 1055 $p(\mathcal{N}) \propto \exp(-\epsilon \|\mathcal{N}\|)$, before transmission, as shown in Eq 6. Firstly, to simply the question, we 1056 only consider the situation of only one edited head h. Assume Z_h represents the normalized weighted 1057 sum of restoration vectors of all privacy spans on head h without adding noise, \mathcal{R}_h denotes the vector 1058 on head h after adding noise, as shown in Eq 5 and 6. The process of adding noise can be represented 1059 by \mathcal{M} . Then, the possibility that Z_h becomes \mathcal{R}_h after adding noise \mathcal{N} is

$$\mathbb{P}(\mathcal{M}(Z_h) = \mathcal{R}_h) =$$

$$\mathbb{P}(\mathcal{M}(Z_h) = \mathcal{R}_h) = \mathbb{P}(Z_h + \mathcal{N} = \mathcal{R}_h)$$
$$= \mathbb{P}(\mathcal{N} = \mathcal{R}_h - Z_h)$$

$$= \exp(-\epsilon ||\mathcal{R}_h - Z_h||).$$

Then for any normalized weighted sums on head h, Z_h and Z'_h , we have 1064

1065
1066
1067
1069

$$\frac{\mathbb{P}[\mathcal{M}(Z_h) = \mathcal{R}_h]}{\mathbb{P}[\mathcal{M}(Z'_h) = \mathcal{R}_h]} = \frac{\exp(-\epsilon ||\mathcal{R}_h - Z_h||)}{\exp(-\epsilon ||\mathcal{R}_h - Z'_h||)}$$

$$= \exp(\epsilon(||\mathcal{R}_h - Z'_h|| - ||\mathcal{R}_h - Z_h||))$$

1070 According to the definition of d_{χ} -privacy in Section 2.1, the mechanism \mathcal{M} satisfies d_{χ} -privacy. In other words, by adding noise \mathcal{N} , adversaries cannot infer Z_h from \mathcal{R}_h even if \mathcal{R}_h is intercepted. 1071 Moreover, the privacy budget of our methods is $\epsilon ||Z'_h - Z_h||$. 1072

 $\leq \exp(\epsilon ||Z_h' - Z_h||).$

Then for multiple edited heads, the only difference is to concatenated all restoration vectors \mathcal{R}_h to 1074 form a single vector, as shown in Section 4.2. Then, the concatenated vector is also added with 1075 noise to protect privacy and the privacy budget of our methods become $\epsilon ||Z' - Z||$, where Z' and Z represent any pair of normalized weighted sums of restoration vectors concatenated across all edited heads. It is independent of the input length n and depends on the hyperparameter ϵ the term 1077 ||Z' - Z||. Therefore PrivacyRestore fulfills d_{χ} -privacy and provides a privacy budget $\epsilon ||Z' - Z||$ 1078 which is independent of the input length and inherently addresses the problem of the linear growth of 1079 privacy budget.

1080 **TOP-K HEADS SELECTOR ALGORITHM** G

1082

In the section, we present the detail implementation of Top-K Heads Selector, as shown in Algorithm 1083 1. Firstly, we initialize an empty score list L_h for each head. Secondly, each privacy span s has its 1084 corresponding top-K heads set \mathcal{H}_k^s . For each head h in \mathcal{H}_k^s , we append Score (h, \mathcal{H}_k^s) into its score 1085 list L_h . Score (h, \mathcal{H}_k^s) is defined as the rank of head h among \mathcal{H}_k^s in ascending order based on the 1086 accuracy of the probe associated with the head h and privacy span s. Thirdly, we calculate the average value of each score list L_h as the score of the corresponding head h. Finally, we sort all heads in the 1087 LLM by the scores and pick up top-K heads as the common top-K heads set \mathcal{H}_c . 1088

1089 Algorithm 1 Top-K Heads Selector

1090 **Input:** S is the set of privacy spans; \mathcal{H}_a is the set of all heads; \mathcal{H}_k^s is the top-K heads set of the 1091 privacy span s; Score (h, \mathcal{H}_k^s) return the rank of head h among \mathcal{H}_k^s in ascending order based on the 1092 accuracy of the probe associated with the head h and privacy span s. The score of the head with 1093 lowest accuracy is 1. The score of the head with highest accuracy is K. 1094 1: Initialize an empty score list $L_h = []$ for each head h in \mathcal{H}_a . 1095 2: for s in S do 3: for h in \mathcal{H}_{k}^{s} do Append Score (h, \mathcal{H}_k^s) into L_h . 4: end for 5: 1099 6: end for 1100 7: for h in \mathcal{H}_a do 1101 7: $score_h = average(L_h)$ 1102 8: end for 9: Sort \mathcal{H}_a according to score_h and select top K heads to obtain **common top-K heads set** \mathcal{H}_c . 1103 1104 **Output:** \mathcal{H}_c is the common top-K heads set. 1105

Settings of Privacy Hyperparameter ϵ 1107 Η

1108

1106

1109 As demonstrated in Appendix E, the privacy budget of d_{γ} -privacy is $\epsilon n d_e$ and the privacy budget 1110 of $d\chi$ -privacy on privacy spans is $\epsilon n_{sp} d_e$, where n represents the length of user inputs, n_{sp} denotes the length of privacy spans, d_e denotes the average distance between word embeddings and ϵ is the 1111 privacy hyperparameter. In addition, as pointed by Mattern et al. (2022b); Utpala et al. (2023), the 1112 privacy budget of paraphrase method is $2n/\tau$, where τ is the generation temperature used during 1113 paraphrasing, and n represents the average length of user inputs. To ensure same privacy budget for 1114 fair comparison, we need to determine the values of different hyperparameters for different methods 1115 on different datasets. 1116

The privacy budget of PrivacyRestore is ϵd_z , according to Appendix F, where d_z denotes the 1117 average distance between normalized weighted sums, computed as $||Z'_h - Z_h||$. We set the privacy 1118 hyperparameter ϵ to 75.00. 1119

1120 To achieve the same privacy budget when using other privacy-preserving methods, we first analyze the 1121 distribution of users inputs lengths n, privacy spans lengths n_{sp} , distances between word embeddings 1122 d_e and distances between normalized weighted sums d_z across three privacy-preserving datasets in Figure 4. The average values of these distributions are used as prior estimates for the corresponding 1123 parameters. We then compute the corresponding ϵ for d_{χ} -privacy (on privacy spans) and τ for 1124 paraphrase across different datasets. The privacy hyperparameter ϵ and τ for different baselines 1125 across three privacy-preserving datasets are shown in Table 5. 1126

1127

Ι 1128 HYPERPARAMETER ANALYSIS

1129

We evaluate the performance of our methods using different numbers of edited heads, K, across 1130 the development sets of three privacy-preserving datasets. For simplicity, we compute MC2 to 1131 represent classification performance, LLM-J to measure generation performance, and TP to indicate 1132 inference efficiency. As shown in Table 6, according to the MC2 score, the optimal value of K is 175 1133 for the Pri-DDXPlus and Pri-SLJA datasets, and 125 for the Pri-NLICE dataset. The performance



Figure 4: The distributions of privacy span lengths n_{sp} , total input lengths n, distances between word embeddings d_e , and distances between normalized weighted sums d_z on the dev set across all three privacy-preserving datasets are analyzed. The average values from these distributions are used as prior estimates for the corresponding parameters.

	a	l_{χ} -privacy	y	d_{χ} -priva	icy on priv	acy spans	Parapl	ırase	Privacy	Restore	Total Privacy
Datasets	n	d_e	ϵ	n_{sp}	d_e	ϵ	n	τ	d_z	ϵ	Budget
Pri-DDXPlus	69.63	1.12	0.87	20.28	1.12	3.00	69.63	2.04	0.91	75.00	<u>68.25</u>
Pri-NLICE	40.14	1.09	1.06	14.08	1.09	3.03	40.14	1.72	0.62	75.00	<u>46.50</u>
Pri-SLJA	129.76	1.15	0.54	16.33	1.15	4.35	129.76	3.17	1.09	75.00	<u>81.75</u>

Table 5: The settings of privacy hyperparameters for different baselines across all privacy-preserving datasets.

1188 degradation as K increases can be attributed to the cumulative effect of multiple edited heads. As 1189 more heads are modified, the activations progressively deviate from their initial values, potentially 1190 compromising the LLM's general capabilities. Moreover, throughput increases with larger K because 1191 we need to inject the meta vector for each head in \mathcal{H}_c using Eq 7 on the server. Consequently, more 1192 heads indicate more injections, which increases the inference time on the server.

Datasets	Metrics	K = 75	K = 100	K = 125	K = 150	K = 175	K = 200
	MC2↑	52.20	56.17	59.39	58.96	62.95	62.64
Pri-DDXplus	LLM-J↑	4.51	4.38	4.45	4.33	4.71	4.55
1	TP↑	24.31	21.51	19.72	20.07	22.68	21.91
Pri-NLICE	MC2↑	37.15	51.01	58.97	51.89	58.11	58.45
	LLM-J↑	3.27	3.66	3.80	3.44	3.40	3.62
	TP↑	20.05	19.14	18.23	16.08	15.89	15.48
Pri-SLJA	MC2↑	28.75	30.65	35.07	32.41	35.13	32.08
	LLM-J↑	5.21	5.41	5.00	5.33	5.15	5.28
	TP↑	36.28	35.25	34.62	32.97	30.51	29.87

Table 6: The performance of PrivacyRestore on the development set using various numbers of edited 1204 heads K. MC2 reflects classification capability, while LLM-J indicates generation performance. The 1205 TP assesses inference efficiency. We report results across three datasets to identify the optimal K for 1206 each datasets. The best results are highlighted in bold. 1207

1208

1201 1202 1203

1209 J CALCULATION PROCESS OF MC1 AND MC2 1210

1211 We evaluate the model's classification ability using two metrics: MC1 and MC2 (Lin et al., 2021). 1212 We assign each sample in Pri-DDXPlus, Pri-NLICE and Pri-SLJA with four options, including one 1213 correct answer and three incorrect ones. The details of calculation process is as follows: 1214

Calculation of MC1. For each user input, we select the option with the highest probability as the 1215 model's choice. MC1 is defined as the model's accuracy, which is calculated as the proportion of 1216 correctly answered inputs. 1217

1218 Calculation of MC2. For each user input, we compute the normalized probability of the correct 1219 answer among the four options. The average of these normalized probabilities across all inputs is calculated as the MC2 score. 1220

1221

1222 Κ VARYING LLM BACKBONE 1223

1224 We evaluate the performance of PrivacyRestore and other privacy-preserving baselines on a larger 1225 model, Llama-13b-chat. As shown in Figure 5, PrivacyRestore outperforms the other baselines in 1226 terms of both MC2 and LLM-J values across all three privacy-preserving datasets. Notably, the 1227 performance of all privacy-preserving methods on the larger model, Llama-13b-chat, is worse than on the smaller model, Llama-7b-chat. This suggests that as model size increases, the model becomes 1228 more sensitive to the injected disturbances introduced by these privacy-preserving methods, leading 1229 to performance degradation. 1230

1231

1232 SUPPLEMENTS OF EMPIRICAL PRIVACY PROTECTION L

1233

1234 We present empirical evidence of the privacy protection capabilities of d_{γ} -privacy and PrivacyRestore by implementing various attacks on these privacy-preserving methods. Lower attack performance 1236 indicates stronger privacy protection provided by these methods.

1237

1238 L.1 EMPIRICAL PRIVACY PROTECTION OF d_{χ} -privacy 1239

As presented by Feyisetan et al. (2019; 2020), the d_{γ} -privacy mechanism protects input by injecting 1240 noise into the token embeddings and replacing the original tokens with their nearest neighbors. To 1241 attack the garbled text input, we implement two types of attacks: prompt injection attack (Suo, 2024)



Figure 5: The MC2 and LLM-J results of PrivacyRestore and other privacy-preserving baselines on larger model, Llama-13b-chat, across three datasets.

and attribute inference attack (Li et al., 2022), both commonly used for attacking text inputs. In
 prompt injection attack, additional instructions are added before and after the garbled text input,
 prompting the model to output the original text. The attack's performance is measured by calculating
 the ROUGE-L score between the generated text and the original input. Attribute inference attack
 performs classification on the garbled text, where the target labels are the token IDs of the original
 input. The attack's performance is evaluated using the classification F1 score. The details of these
 attack methods are presented in Appendix M.

1266 1267 L.2 EMPIRICAL PRIVACY PROTECTION OF PRIVACYRESTORE

1268 As stated in Section 6.3, the attack object of PrivacyRestore is the meta vector and we implement 1269 two types of attack: the embedding inverse attack (Li et al., 2023b; Morris et al., 2023) and attribute 1270 inference attack (Li et al., 2022), both commonly used methods for attacking embeddings. Embedding inverse attack utilizes the generative model to generate the privacy spans from the meta vector. 1271 Embedding inverse attack measure the attack performance by computing the ROUGE-L between 1272 the generate output and the privacy spans. Attribute inference attack performs classification on the 1273 meta vector and the target labels are the token IDS of the privacy spans. Attribute inference attack 1274 compute the F1 score to evaluate the attack performance. The details of these attack methods are 1275 presented in Appendix M. 1276

As shown in Figure 3(c) and 3(d), the ROUGE-L score for the embedding inverse attack remains nearly stable across different α values in the Pri-SLJA and Pri-DDXPlus datasets. What's a little strange is the ROUGE-L score in the Pri-NLICE dataset shows a slight increase. The possible reason is that higher ratio indicating more privacy spans consider and resulting longer reference string when compute the ROUGE-L score. Since ROUGE-L measures the overlap between the generated output and the reference string, a longer reference string may slightly boost the score. The F1 score for the attribute inference attack remains stable across all three datasets. The stable performance in both attack scenarios provides empirical support for Theorem 5.2.

1284 1285

1287

1286 M

DETAILED IMPLEMENTATION OF ATTACK METHODS

1288 M.1 PROMPT INJECTION ATTACK

1289 1290 d_{χ} -privacy injects noise into the original user inputs and transmits the garbled inputs to the server to 1291 protect the privacy spans. We employ prompt injection attack to recover the initial question from 1292 the garbled inputs. Following Perez & Ribeiro (2022); Suo (2024), we insert additional instructions 1293 before and after the user inputs to prompt the model to output the original user input rather than the 1294 normal response. The template for the additional instructions is provided in Appendix O.4.

1295 We set the maximum generation length for the prompt injection attack to 256 tokens. To evaluate the attack's performance, we calculate the ROUGE-L score between the generated output and the

original user input. A higher ROUGE-L score indicates greater overlap between the recovered text
 and the original input, signifying more successful attack results.

1299 M.2 ATTRIBUTE INFERENCE ATTACK

1301 Attribute inference attack attempts to steal user inputs by performing classification on the garbled 1302 inputs, where the target labels correspond to the token IDs of the original inputs. Since each input contains multiple tokens, this classification task is naturally a multi-label classification problem. 1303 Following Li et al. (2022), we utilize a multi-layer perceptron (MLP) model as the classifier. The 1304 input dimension is 4096 and the output dimension is the size of the whole vocabulary size. To 1305 evaluate the attack's performance, we calculate the F1 score of the classification, where a higher F1 1306 score indicates a more successful attack. The attack targets can include garbled text from d_{γ} -privacy, 1307 paraphrased text, or the meta vector from PrivacyRestore. The implementation details for text and 1308 vectors may vary slightly. 1309

1310Attribute inference attack on meta vector. To attack the meta vector from PrivacyRestore, we can
directly use a fully-connected layer to transform the meta vector's dimension from $128 \times K$ to the
classifier's input dimension 4096. We then perform classification on the transformed meta vector.

1313 Attribute inference attack on text. For garbled text from d_{χ} -privacy or paraphrased text, we first 1314 transform the text into a vector representation. We utilize Llama2-chat-7b to process the text input 1315 and obtain the last token's hidden state as the vector representation. Classification is then performed 1316 on this hidden state.

- 1317
- 1318 M.3 EMBEDDING INVERSE ATTACK

1319 Different from attribute inference attack, embedding inverse attack steal the user inputs through 1320 the generative model to generate the original user inputs. We utilize the GPT-2 model (Radford 1321 et al., 2019) as the generative model and set the maximum generation length to 256. We finetune the 1322 GPT-2 model on the training set for 20 epoch using the learning rate of 1e-5. To evaluate the attack's 1323 performance, we compute the ROUGE-L score between the generated output of the GPT-2 model 1324 and the original user input, where higher scores indicate better attack effectiveness. Similar to the 1325 attribute inference attack, the implementation of the embedding inverse attack differs between text 1326 and meta vectors.

Embedding inverse attack on meta vector. We use a fully-connected layer to transform the meta vector's dimension to the dimension of hidden state of GPT-2 model. Then we directly input the transformed meta vector as the input embedding.

Embedding inverse attack on text. We use Llama-2-chat-7b to process the text input and extract the last token's hidden state as the vector representation. This vector is then transformed by a fully connected layer to match the hidden state dimension of the GPT-2 model. Finally, it is fed into the GPT-2 model as the input embedding for subsequent generation.

1334

N EXAMPLE OUTPUTS OF PRIVACYRESTORE

1336 1337

We provide some example outputs of our method in Figure 6. As shown in these examples, applying 1338 d_{γ} -privacy to privacy spans results in outputs with higher ROUGE-L scores but lower LLM-J scores 1339 compared to our method. After analyzing these outputs in detail, the high ROUGE-L scores from 1340 d_{γ} -privacy on privacy spans likely result from a greater overlap with the initial output. However, 1341 the overlapping sections consist mainly of meaningless sentence structures and lack diagnostic 1342 information. Moreover, the final diagnosis is incorrect, leading to lower LLM-J scores. In contrast, PrivacyRestore generates outputs with a different structure but provides the same, correct diagnosis. 1344 As a result, our method achieves slightly lower ROUGE-L scores but significantly higher LLM-J 1345 scores compared to d_{γ} -privacy on privacy spans.

- 1346
- 1347

1348



Figure 6: Example Outputs of PrivacyRestore and d_{χ} on privacy spans in Pri-DDXPlus dataset. Text highlighted with a yellow background represents the privacy spans in user inputs. Text highlighted with a green background indicates the correct diagnosis. Text highlighted with a red background denotes the incorrect diagnosis. Underscored text marks sections that overlap with the initial output.

PROMPT TEMPLATE DETAILS

O.1 CLASSIFICATION OF PRIVACY SPANS.

O.1.1 MEDICAL DATASETS (PRI-DDXPLUS/PRI-NLICE).

The prompt template is for GPT and used to classify symptoms in Pri-DDXPlus/Pri-NLICE dataset into sensitive and non-sensitive categories. GPT grades the symptoms on a scale of one to five based on sensitivity, with levels greater than three considered private spans in the Pri-DDXPlus/Pri-NLICE dataset.

1415	Y
1416	sy
417	
418	



O.1.2 LEGAL DATASET (PRI-SLJA).

The prompt template is for GPT and used to classify the case details in Pri-SLJA dataset into sensitive and non-sensitive categories. GPT grades the symptoms on a scale of one to five based on sensitivity, with levels greater than three considered private spans in the Pri-SLJA dataset.

1440	You are a privacy expert with extensive legal knowledge, please categorize the given case element into five levels based on the
1441	degree of privacy involved.
1442	
1443	Utilize the classification guidelines provided below to categorize the case element:
1440	Degrade 1: public legal information: information that is publicly available and does not reveal personal or sensitive legal
1444	details. Examples include general legal concepts, publicly available case laws, or minor regulatory matters.
1445	Degrade 2: non-sensitive personal legal information: Legal information that may involve a specific individual but is not
1446	particularly sensitive or revealing. Examples include minor legal infractions like parking tickets, common civil disputes, or
1440	routine contracts such as lease agreements.
1447	Degrade 3: potentially sensitive legal information: Legal information that could imply more personal aspects but is not
1448	immediately sensitive or damaging. Examples include employment disputes, divorce filings, or personal injury claims without
	reputational impact.
1449	Degrade 4: sensitive personal legal information: Information that could be stigmatizing or have significant social, financial, or
1450	reputational consequences. Examples include legal matters involving bankruptcy, criminal records (misdemeanors), or family
4.454	law issues like child custody disputes.
1451	Degrade 5: highly sensitive legal information: Legal details that are highly personal, stigmatizing, or potentially life-altering
1452	if disclosed. This includes serious criminal cases (e.g., felonies, sexual offenses), immigration status disputes, or legal
1453	proceedings involving whistleblowing or confidential corporate matters.
4 4 5 4	Output your evaluation in the following format:
1454	#thereason: your analysis here [note: this should be step-by-step analysis following the steps: (a) Give an overview of the given
1455	case element. (b) Evaluate the degree of the given case element according to classification guidelines.]
1456	#thescore: your degrade here. [note: provide your degrade in based on your analysis and the classification criteria. Please note
1457	that only the degrade is required to be written in here] Case element: [INSERT CASE ELEMENTS HERE].

1458 O.2 GENERATION PROMPTS.

1460 O.2.1 MEDICAL DATASETS (PRI-DDXPLUS/PRI-NLICE).

The prompt template is for Llama model and used during model generation for the Pri-DDXPlus/Pri-NLICE datasets.

[INST] <<SYS>>

You are a medical doctor answering real-world medical entrance exam questions. Based on your understanding of basic and clinical science, medical knowledge, and mechanisms underlying health, disease, patient care, and modes of therapy, answer the following diagnosis question. Base your answer on the current and standard practices referenced in medical guidelines. <</SYS>>

[INSERT USER INPUT HERE] [/INST]

O.2.2 LEGAL DATASET (PRI-SLJA).

¹⁴⁷⁴ The prompt template is for Llama model and used during model generation for the Pri-SLJA datasets.

1/177	[INST] < <sys>></sys>
1477	You are the world's best legal judgment machine, tasked with rendering a verdict in a legal case. Below are four paragraphs
1478	describing the case, though some information may be missing.
479	Subject refers to the individual with criminal responsibility who committed the crime;
480	Subjective Element consists of intention and negligence, representing the psychological state of the subject concerning the
1400	criminal behavior and its outcomes;
1481	Object refers to the social relationships protected by criminal law that have been infringed upon by the criminal acts;
482	Objective Elements include the specific details of the crime, such as time, location, actions, and consequences.
483	Using these four components, generate a judgment to sentence the subject of the crime.
403	<>
1484	
1485	[INSERT USER INPUT HERE] [/INST]

1488 O.3 EVALUATION OF GENERATED OUTPUT.

14891490O.3.1MEDICAL DATASETS (PRI-DDXPLUS/PRI-NLICE).

The prompt template is for GPT and evaluates the quality of generated output based on relevance, clarity, and accuracy, for Pri-DDXPlus/Pri-NLICE dataset. Scores range from 1 to 10, with higher values indicating better output.

1495	
1496	You are a neutral legal evaluator. You will be provided with a legal question containing the description of a case and the AI assistant's indement of the case
1497	Your task is to assess whether the AI assistant's judgment is useful and accurately points out the crime. Consider factors such as the relevance,
1/108	clarity, and accuracy of the judgment.
1450	Rate the response on a scale from 1 to 10, where 1 indicates poor quality and 10 indicates excellent quality.
1499	Start your answer with the rating score by strictly following the format: `[[rating]]`. After assigning a rating, provide a brief explanation for
1500	your score, justifying your evaluation.
1000	[User Question]
1501	NSERT QUESTION HERE].
1502	[The Start of Assistant Response]
1503	INSERT GENERATED OUTPUT HERE]. The End of Assistant Response]
1504	
1505	Rating:
1506	

1508 O.3.2 LEGAL DATASET (PRI-SLJA).

The prompt template is for GPT and evaluates the quality of generated output based on relevance,
 clarity, and accuracy, for Pri-SLJA dataset. Scores range from 1 to 10, with higher values indicating better output.

You are a neutral medical evaluator. You will be provided with a medical question from a patient and the AI assistant's response to the question. Your task is to assess whether the AI assistant's response is useful and accurately addresses the diagnosis. Consider factors such as the relevance, clarity, and accuracy of the information provided. Rate the response on a scale from 1 to 10, where 1 indicates poor quality and 10 indicates excellent quality. Start your answer with the rating score by strictly following the format: `[[rating]]`. After assigning a rating, provide a brief explanation for your score, justifying your evaluation." [User Question] [INSERT QUESTION HERE]. [The Start of Assistant Response] [INSERT GENERATED OUTPUT HERE]. [The End of Assistant Response] Rating:

O.4 PROMPT INJECTION ATTACK.

The prompt template is for Llama model and used to carry out a prompt injection attack, translating the garbled text back into the original text.

 1530
 INST]<<SYS>

 1531
 You are a text generating machine. Given input with garbled characters, output normal text without garbled characters. ONLY

 1532
 generate translated normal text.

 1533
 [INSERT USER INPUT HERE]

 1534
 Translated normal text:[/INST]

 1536