

# Simple Yet Effective: Extracting Private Data Across Clients in Federated Fine-Tuning of Large Language Models

Yingqi Hu<sup>1</sup>, Zhuo Zhang<sup>1</sup>, Jingyuan Zhang<sup>2</sup>, Jinghua Wang<sup>1</sup>,  
Qifan Wang<sup>3</sup>, Lizhen Qu<sup>4\*</sup>, Zenglin Xu<sup>5,6</sup>

<sup>1</sup>Harbin Institute of Technology, Shenzhen, China

<sup>2</sup>Kuaishou Technology, China    <sup>3</sup>Meta AI, USA

<sup>4</sup>Monash University, Australia    <sup>5</sup>Fudan University, China

<sup>6</sup>Shanghai Academy of Artificial Intelligence for Science, China

## Abstract

Federated large language models (FedLLMs) enable cross-silo collaborative training among institutions while preserving data locality, making them appealing for privacy-sensitive domains such as law, finance, and healthcare. However, the memorization behavior of LLMs can lead to privacy risks that may cause cross-client data leakage. In this work, we study the threat of *cross-client data extraction*, where a semi-honest participant attempts to recover personally identifiable information (PII) memorized from other clients' data. We propose three simple yet effective extraction strategies that leverage contextual prefixes from the attacker's local data, including frequency-based prefix sampling and local fine-tuning to amplify memorization. To evaluate these attacks, we construct a Chinese legal-domain dataset with fine-grained PII annotations consistent with CPIS, GDPR, and CCPA standards, and assess extraction performance using two metrics: *coverage* and *efficiency*. Experimental results show that our methods can recover up to 56.6% of victim-exclusive PII, where names, addresses, and birthdays are particularly vulnerable. These findings highlight concrete privacy risks in FedLLMs and establish a benchmark and evaluation framework for future research on privacy-preserving federated learning. Code and data are available at <https://github.com/SMILELab-FL/FedPII>.

## 1 Introduction

Federated large language models (FedLLMs) (Ye et al., 2024a,b; Zhang et al., 2023; Chen et al., 2024; Yao et al., 2024) have recently emerged as a promising approach for cross-silo federated learning (FL) (Li et al., 2024c), enabling collaborative model training while maintaining data locality and institutional privacy. In cross-silo FL, organizations, such as courts, banks, and hospitals,

collaboratively fine-tune<sup>1</sup> a shared model without exchanging private records. Prior studies mainly focused on improving algorithmic efficiency and convergence (Bai et al., 2024; Wu et al., 2025; Li et al., 2020), while the privacy vulnerabilities of FedLLMs remain largely unexplored.

A growing body of work has shown that large language models (LLMs) tend to memorize and reproduce fragments of their training data, including PII such as names, addresses, and dates of birth (Carlini et al., 2021, 2023; Shao et al., 2024; Kim et al., 2023; Nakka et al., 2024). Although FL mitigates privacy risks by exchanging model updates instead of raw data, our preliminary experiments (Appendix B.4) indicate that FedLLMs remain vulnerable to *verbatim data extraction* (VDE)—where adversaries can recover verbatim text sequences from the aggregated global model. However, most existing VDE studies assume that attackers have privileged access or significant knowledge of the victim's data (Yu et al., 2023; Huang et al., 2022), which is unrealistic in practical cross-silo deployments.

We instead consider a more realistic *semi-honest* threat model, where each participant follows the FL protocol but may attempt to infer private information from the global model. For instance, in a federation of courts, a participant could exploit its own local case records as contextual prefixes to elicit sensitive information memorized from other courts' data (as illustrated in Figure 1). To investigate this, we propose three extraction strategies: (1) PII Contextual Prefix Sampling, which queries the global model using local contextual prefixes; (2) Frequency-prioritized (FP) Sampling, which focuses on high-frequency prefixes to improve extraction efficiency; and (3) Latent Association Fine-tuning (LAFt), which locally fine-tunes the global

<sup>1</sup>In current practice, *FedLLM* typically refers to the federated fine-tuning of large language models rather than federated pre-training. See Appendix B.3 for details.

\*Corresponding author.

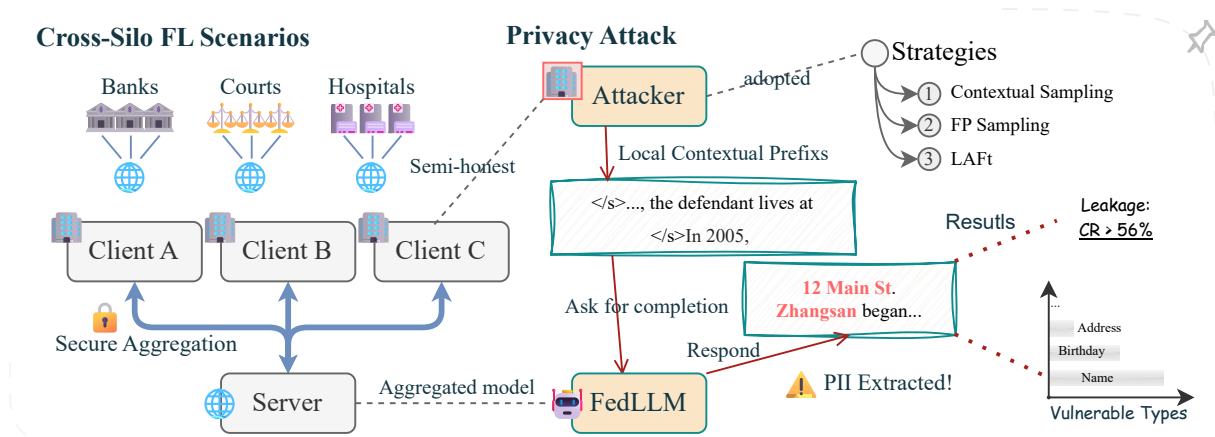


Figure 1: Overview of cross-silo FedLLMs and the proposed privacy attack. In cross-silo FL, institutions such as banks, courts, and hospitals collaboratively fine-tune a shared model under the coordination of a central server, keeping data local. A semi-honest client leverages its local data to construct PII-related prefixes and queries the aggregated global FedLLM, leading to cross-client data leakage. The proposed strategies—Contextual Prefix Sampling, Frequency-prioritized (FP) Sampling, and Latent Association Fine-tuning (LAFt)—achieve up to 56.6% recovery of victim-exclusive PII, with names, addresses, and birthdays being the most vulnerable categories.

model using stored prefix-PII pairs, thereby improving its ability to extract PII implicitly memorized within the model.

To evaluate these attacks, we build a benchmark dataset by annotating a real-world Chinese legal corpus with fine-grained PII labels aligned with privacy regulations such as CPIS, GDPR, and CCPA (see Acronyms List A). We assess the attacks using two metrics: *coverage* (the proportion of target PII successfully extracted) and *efficiency* (the amount of PII recovered under a limited query budget). Our experiments reveal that the proposed attacks can extract up to 56.6% of victim-exclusive PII, with *Name*, *Address*, and *Birthday* being the most vulnerable types. Moreover, we observe diminishing returns as query budgets grow, while FP Sampling and LAFt enhance diversity under tighter budgets. These findings expose concrete privacy risks in FedLLMs and underscore the need for stronger privacy-preserving mechanisms.

In summary, our main contributions are:

1. We propose three novel extraction strategies for FedLLMs, independent of existing gradient-based or membership inference attacks, and evaluate them using two rigorous metrics: coverage and efficiency.
2. Extensive experiments show that our methods can recover up to 56.6% of cross-client unique PII, with larger prefix sets yielding diminishing returns in efficiency, revealing a trade-off between coverage and computational cost.
3. We construct a real-world benchmark dataset

by augmenting a legal-domain corpus with fine-grained PII annotations aligned with CPIS, GDPR, and CCPA standards, filling the gap in public resources for privacy research in federated learning.

## 2 Related Work

This study is related to the fields of data extraction attacks and federated learning. For the reader’s convenience, a brief introduction to these concepts is provided in Appendix B. In this section, we review only the work directly related to our method.

### 2.1 PII Extraction Attacks in LLM

Large language models, due to their massive parameter scale, are capable of memorizing exact training data samples, making them vulnerable to data extraction attacks. These attacks can target different granularities of information: sample-level and entity-level.

At the sample level, an attacker with access to the full prefix of a training sample can query the LLM to regenerate the exact suffix (Yu et al., 2023; Shi et al., 2024; Zhang et al., 2024). This technique, known as *verbatim training data extraction* (Carlini et al., 2021, 2023; Schwarzschild et al., 2024), is widely used to detect data contamination and copyright violations (Dong et al., 2024).

At the entity level, attackers may know a subset of PII entities—such as names or affiliations—about a particular subject. By combining these known details with prompt templates

(either manually crafted or automatically generated (Kassem et al., 2025)), they can elicit the model to produce additional PII records about the same subject. This is known as an associative data extraction attack (Shao et al., 2024; Kim et al., 2023; Zhou et al., 2024).

Broadly, PII extraction attacks refer to any attack that aims at eliciting outputs from the model that contain PII (Lukas et al., 2023; Nakka et al., 2024; Huang et al., 2022). Both verbatim and associative techniques can be used to conduct such attacks.

While most prior work assumes centralized training with full data access, we investigate PII extraction under federated fine-tuning, where the attacker has limited observability and control. We elaborate on this in Section 4.1.

## 2.2 Privacy Threats in Federated Learning

Threats in Federated Learning can be categorized into two main areas: security and privacy (Wang et al., 2024a; Xie et al., 2024; Li et al., 2024b). Security threats typically aim to disrupt the entire FL system by invalidating model training (Shejwalkar and Houmansadr, 2021) and introducing backdoors (Bagdasaryan et al., 2020; Chang et al., 2024). In contrast, privacy threats have attracted more attention from researchers and focus on stealing confidential information from the FL system, such as inferring sensitive properties (Melis et al., 2019), reconstructing clients’ private datasets (Zhu et al., 2019; Geiping et al., 2020), and determining the membership and source of training data (Rashid et al., 2025; Vu et al., 2024; Hu et al., 2024). To achieve these attacks, researchers often make different assumptions regarding the attacker’s knowledge. Common assumptions typically fall into two dimensions: whether the attacker is a client or a server (Chu et al., 2023), and whether the attacker is semi-honest (Applebaum, 2017; Hu et al., 2024) or malicious. These assumptions determine whether the attacker has access to gradients, local datasets, model parameters, and the ability to manipulate them.

## 3 Dataset

### 3.1 Data Sources and Preprocessing

The majority of our dataset is sourced from the Challenge of AI in Law (CAIL) (Li et al., 2024a), supplemented by smaller portions from CJRC (Duan et al., 2019) and JEC-QA (Zhong et al., 2020). CAIL is a renowned annual competi-

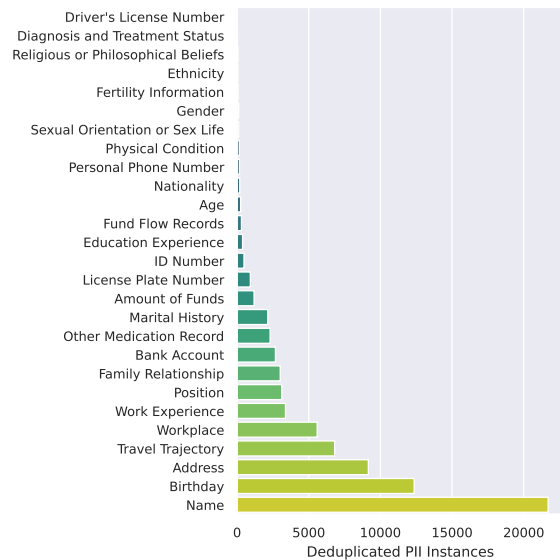


Figure 2: Distribution of de-duplicated PII instances by label category.

tion featuring a variety of legal NLP tasks. In this study, we focus on two natural language generation tasks: *Judicial Summary* and *Judicial Reading Comprehension*, and three natural language understanding tasks: *Similar Case Matching*, *Judicial Exam Question Answering*, and *Legal Case Classification*. Detailed task descriptions are provided in Appendix D, with representative examples shown in Table 6.

Following prior work (Zhang et al., 2023; Yue et al., 2024), we further preprocess and curate the dataset to fit our setting. The complete preprocessing pipeline is described in Appendix E, where Table 7 reports the dataset statistics.<sup>2</sup>

### 3.2 PII Labeling

We reviewed the definitions and examples of PII in various legal provisions, including CPIS, GDPR, CCPA, and Singapore PDPC (see Acronyms List in Appendix A), and used them as references to establish a systematic PII labeling standard. We selected PII types relevant to the text modality and removed types that are unlikely to appear in legal texts (e.g., browser history, SMS content, IP & MAC addresses), as well as those that are difficult to describe or evaluate (e.g., medical examination

<sup>2</sup>The datasets contain PII from publicly available government-published legal documents. They were de-identified and used in prior work, e.g., as in Yue et al. (2024). We use curated versions from these papers. Since our study concerns privacy risks in FedLLMs, real-world PII is necessary to evaluate model vulnerabilities.

reports, psychological trends). Ultimately, we defined labeling guidelines encompassing 7 major categories and 36 subcategories. The distribution of labeled PII types is shown in Figure 2 in the main text, whereas a complete summary of these standards is provided in Table 9 in Appendix I.

We employed a combination of machine-assisted annotation and manual verification to label the data. For each major PII category, we designed a dedicated prompt depicted in Figure 12 and employed GPT-4o (OpenAI et al., 2024) to generate annotations. We then recruited students to verify a subset of those annotations with the help of Label Studio (Tkachenko et al., 2020-2025). The agreement between human evaluations and GPT-4o annotations achieves an F1 score of 89.9%. Owing to space limitations, the full evaluation results are shown in Table 10 in Appendix J, together with further details of the annotation process, including annotator backgrounds, annotation instructions, and user interfaces.

## 4 Method

### 4.1 Problem Definition

We study a novel extraction attack tailored to FedLLMs, which differs from traditional verbatim data extraction in three key aspects:

**Assumptions.** Unlike VDE that assumes the attacker has access to most or all of the training data, our setting limits the attacker to a small, isolated subset of the overall training corpus.

**Setup.** In our formulation, the prefix and its corresponding target suffix are not drawn from a contiguous span of training data. Instead, extraction prefixes are sampled from the attacker’s local dataset  $D_a$ , while the target suffixes reside exclusively in other clients’ private data and are absent from  $D_a$ . Thus, each prefix must generalize beyond local context to trigger the generation of unseen suffixes.

**Goals.** An attacker does not aim to recover all training completions, but instead focuses on extracting specific, high-value information—most notably, PII—from the global model.

Formally, we consider a **cross-silo** FL system comprising  $c$  clients  $\mathcal{C} = \{C_1, C_2, \dots, C_c\}$ , where each client  $C_i$  holds a local dataset  $D_i$ . Among them, one client—denoted as the attacker  $C_a \in \mathcal{C}$ —is assumed to be **semi-honest** (Applebaum, 2017; Hu et al., 2024). That is,  $C_a$  faithfully follows the FL protocol (e.g., does not poison data or

manipulate model weights), but acts adversarially in a passive manner, attempting to infer PII contained in other clients’ datasets by analyzing the global model  $\theta$ .

In this setting, the attacker issues queries to the model  $\theta$  to extract data without knowing which client any particular output originates from. However, for evaluation purposes, we designate one client as the reference victim to measuring the attack’s effectiveness. Let  $S_a$  and  $S_v$  denote the sets of PII instances held by the attacker and the victim client, respectively. The attacker constructs a prompt set  $P$  and queries the FedLLM  $\theta$ , obtaining a corresponding output set  $Y$ . We formalize key definitions and evaluation metrics in the following subsections.

**Definition 1** (Extracted). *A PII instance  $s \in S_v$  is considered successfully extracted if there exists a prompt  $p \in P$  and a corresponding model output  $y \in Y$  such that:*

$$\exists u \in \Sigma^* \text{ such that } y = s \oplus u, \quad (1)$$

where  $\Sigma^*$  is the set of all finite-length strings over the vocabulary, and  $\oplus$  denotes string concatenation. In other words, the model output  $y$  begins with  $s$ .

**Definition 2** (Coverage Rate). *The coverage rate measures how thoroughly the attacker recovers the PII unique to the victim client. It is defined as:*

$$S_E = \{s_i \mid \exists y \in Y \text{ such that } s_i \text{ is extracted by } y\},$$

$$CR = \frac{|(S_v \setminus S_a) \cap S_E|}{|S_v \setminus S_a|}. \quad (2)$$

A higher CR indicates that a larger fraction of the victim’s unique PII has been successfully extracted.

**Definition 3** (Efficiency). *Efficiency quantifies the precision of extraction with respect to the number of queries. Let  $Q$  denote the number of queries, the efficiency is defined as:*

$$EF = \frac{|(S_v \setminus S_a) \cap S_E|}{Q}. \quad (3)$$

A higher EF indicates that more PII is extracted with fewer queries.

Building upon these definitions, the central challenge is to design algorithms that enable the attacker to extract PII both comprehensively and efficiently—that is, achieving high coverage and high efficiency.



## 4.2 Attacking Algorithms

### 4.2.1 PII-contextual Prefix Sampling

We start with a simple method for constructing query prompts using contextual prefixes of PII, which are word sequences immediately preceding PII instances in the attacker’s dataset  $D_i$ . This mitigates reliance on manually crafted prompts (e.g., ‘my phone number is’), which often deviate from the model’s training distribution and exhibit limited empirical applicability.

Let the attacker’s training corpus be  $U_a = \{t_0, t_1, \dots, t_{|U_a|}\}$ , formed by concatenating samples in  $D_i$ , with  $\mathcal{S}$  the multiset of labeled PII. For a PII instance  $s \in \mathcal{S}$ , let  $\text{Loc}(s)$  be the index of its first token in  $U_a$ . We define a  $\lambda$ -length contextual prefix function:

$$\mathcal{T}_\lambda(U, s) = t_{\text{Loc}(s)-\lambda} \cdots t_{\text{Loc}(s)-1}.$$

The contextual prefix set of a PII set  $\mathcal{S}$  is given by:

$$P_c = \{\mathcal{T}_\lambda(U_a, s) \mid s \in \mathcal{S}\}. \quad (4)$$

For each  $p \in P_c$ , the attacker  $\mathcal{C}_a$  queries the global model  $\theta$  to generate a suffix  $y$  of up to  $m$  tokens:

$$y = \{x_1, \dots, x_m\} \sim \mathbf{P}(y \mid p; \theta).$$

To enhance diversity,  $n$  independent suffixes may be sampled per prefix:

$$Y_p = \{y_1, \dots, y_n\}, \quad Y = \bigcup_{p \in P_c} Y_p, \quad Q = n \cdot |P_c|.$$

A generalized version extends  $P_c$  by including all substrings ending before each PII:

$$\text{SUP}(P_c) = \{t_i \cdots t_{\text{Loc}(s)-1} \mid (\text{Loc}(s)-i) \in [1, \lambda]\}.$$

This yields broader coverage but incurs high query cost due to the large prefix set.

### 4.2.2 Frequency-Prioritized Prefix Sampling

Following prior work (Shao et al., 2024), which associates extraction effectiveness with co-occurrence frequency, we posit that prefixes frequently occurring before PII entities tend to capture stronger and more diverse associations. Therefore, we emphasize high-frequency prefixes to construct a compact, information-rich set.

Formally, we partition  $\text{SUP}(P_c)$  by prefix frequency. For each  $\sigma \geq 1$ ,

$$P_\sigma = \{p \in \text{SUP}(P_c) \mid \text{Count}_{\text{SUP}(P_c)}(p) = \sigma\}.$$

This yields

$$\text{Set}(\text{SUP}(P_c)) = \bigcup_{\sigma \geq 1} P_\sigma.$$

Given a threshold  $\sigma_a$ , the frequent prefix set is

$$P_{f \geq \sigma_a} = \bigcup_{\sigma \geq \sigma_a} P_\sigma,$$

sorted by frequency. Setting  $\sigma_a = 1$  recovers the full contextual prefix set. With a budget  $B$ , we select the top- $B$  prefixes from  $P_{f \geq \sigma_a}$ , thereby emphasizing frequent contexts.

### 4.2.3 Latent Association Fine-tuning

We conjecture that a model’s vulnerability to PII extraction stems from its capacity to capture the conditional probability  $\mathbf{P}(\mathcal{B} \mid \mathcal{A}; \theta)$  under model parameters  $\theta$ , where  $\mathcal{A}$  denotes prefixes that typically precede PIIs and  $\mathcal{B}$  represents the corresponding PII instances.

Since the association between PII and their prefixes is implicitly encoded in the model’s representations, we propose Latent Association Fine-tuning (LAFt), which updates parameters to maximize  $\mathbf{P}(\mathcal{B} \mid \mathcal{A}; \theta)$ . The goal is to strengthen the mapping between indicative prefixes and PII, thereby improving extraction.

As the first step, we build a fine-tuning dataset  $D_{\text{ft}}$  by pairing frequent prefixes with known PII:

$$D_{\text{ft}} = \{(p, s) \mid p \in P_f, s \in S_a\}, \quad (5)$$

where  $P_f$  is the frequent-prefix set from  $D_a$ , and  $S_a$  the attacker’s PII set. The model is then fine-tuned with the standard causal LM objective:

$$\theta' = \arg \min_{\theta} \sum_{(p,s) \in D_{\text{ft}}} \sum_{t=1}^{|s|} -\log \mathbf{P}(s_t \mid p, s_{<t}; \theta).$$

The updated model  $\theta'$  is subsequently used for extraction with prefixes from  $P_f$  or  $P_c$ .

To note that, LAFt follows the **semi-honest** FL setting such that the fine-tuned model remains local and is not uploaded to the server.

## 5 Experiment

### 5.1 Experimental Setup

**Federated Setup.** Our federated setup consists of two main components: the data partitioning across clients and the federated fine-tuning procedure.

For the data partitioning, we simulate a system with 5 clients using a label-skewed non-IID partitioning based on clustering of language embeddings (Li et al., 2023), and ensured that each client receives a comparable number of samples.

For the federated fine-tuning, we perform training on legal tasks using the OpenFedLLM framework (Ye et al., 2024b), with FedAVG (McMahan et al., 2023) as the aggregation method over 10 communication rounds. All clients adopt parameter-efficient fine-tuning (LoRA) and a shared prompt template. Hyperparameter settings and implementation details are provided in Appendix K.

After the federated fine-tuning, we evaluate the utility of the final global model on a held-out global test set. Following common practice, we compare it to a centrally trained (non-FL) baseline evaluated on the same test set. The results are reported in Table 8 in the Appendix.

**Models and Metrics.** Since our data and tasks are derived from Chinese legal documents, we focus primarily on LLMs with proficiency in Chinese. Specifically, we evaluate Qwen1-8B (Bai et al., 2023), Baichuan2-7B (Yang et al., 2023), Qwen3-8B (Yang et al., 2025), GLM4-8B (GLM et al., 2024), and Llama3-Chinese (Raytrfr, 2024).<sup>3</sup>

We evaluate model performance using two primary metrics: Coverage Rate (CR), Efficiency (EF). In addition, we introduce an intermediate metric, Victim-exclusive Extracted PII (VxPII), defined as  $|(S_v \setminus S_a) \cap S_E|$ , which directly measures the amount of extracted information.

**Attack Strategies.** We designate client 0 as the attacker and client 1 as the victim, and evaluate three strategies: (1) PII-contextual prefix sampling. The attacker builds a prefix set  $P_c$  from its local dataset  $D_0$  with prefix length  $\lambda = 50$ . Each prefix queries the global model 15 times, generating up to  $m = 10$  tokens per query—sufficient to recover most PII with manageable cost. (2) Frequency-prioritized sampling. Prefixes in  $\text{Set}(\text{SUP}(P_c))$  are ranked by frequency to form  $P_{f \geq 1}$  and used in descending order. Sweeping the prefix budget  $B$  varies the frequency threshold  $\sigma_a$ , enabling analysis of coverage–efficiency trade-offs. (3) Latent association fine-tuning. The attacker fine-tunes the global model (1 epoch, LR =  $5e-5$ , LoRA:  $r = 16$ ,  $\alpha = 32$ ) using 10k frequent prefixes and 10k randomly sampled PII from its own data to reinforce

<sup>3</sup>All models are publicly available on HuggingFace: Qwen1, Baichuan2, Qwen3, GLM4, Llama3-Chinese.

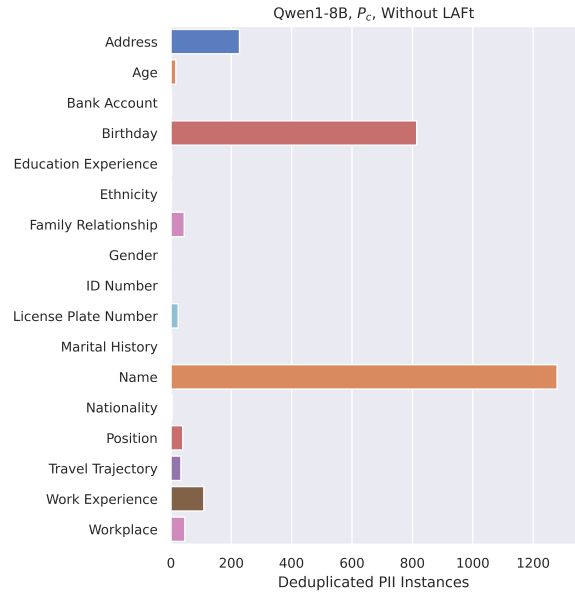


Figure 3: Label distribution of deduplicated victim-exclusive PII extracted by Qwen1-8B (without LAFt, using prefix set  $P_c$ ). Results for Baichuan2-7B are shown in Appendix Figure 11.

prefix–PII associations. Further implementation details are provided in Appendix K.3.2.

**Evaluation Protocol.** To ensure a fair evaluation, the set of victim-exclusive PII ( $S_v \setminus S_a$ ) is obtained by applying two filters: (1) retain only those victim PII that do not appear in the attacker’s training corpus (i.e.,  $s_i \in S_v$  but  $s_i \notin U_a$ ); and (2) remove PII that share a common prefix to avoid ambiguity in identifying which PII was extracted (see Equation (1)). This is enforced by constraining the length of the longest common prefix (LCP) between any two PII:

$$\text{LCP}(s_i, s_j) = 0, \quad \forall s_i \neq s_j \in S_v$$

Metrics defined in Equations (2) and (3) are then computed on this filtered, prefix-disjoint set.

## 5.2 Results and Discussions

**RQ1: How effective is the PII extraction attack using contextual prefixes?** We first evaluate the coverage rate (CR) and efficiency (EF) of our extraction attacks by querying federated fine-tuned LLMs using the PII-contextual prefix set  $P_c$ . Table 1 presents the results. With  $P_c$ , our attack achieves a considerable CR of 22.93% on Qwen1-8B and 28.95% on Baichuan2-7B.

To understand what types of PII are most vulnerable, we analyze the extracted instances. Figure 3 shows the label distribution of deduplicated

Table 1: Summary of attack results using the PII-contextual prefix sampling method (with and without LAFt), where client 0 (attacker) targets client 1 (victim). The victim-exclusive set ( $S_v \setminus S_a$ ) includes 8,870 unique PII items.

Model	Prefix Set	CR	EF	VxPII Count	Prefix Set Size
<i>wo LAFt</i>					
Qwen1-8B	$P_c$	22.93%	0.1910%	2034	71006
Baichuan2-7B	$P_c$	28.95%	0.2411%	2568	71006
Qwen3-8B	$P_c$	30.69%	0.2556%	2722	71006
GLM4-9B	$P_c$	28.20%	0.2348%	2501	71006
Llama3-Chinese-8B	$P_c$	19.73%	0.1643%	1750	71006
Qwen1-8B	Set(SUP( $P_c$ ))	56.20%	0.0110%	4985	3013161
Baichuan2-7B	Set(SUP( $P_c$ ))	53.56%	0.0105%	4751	3013161
<i>w LAFt</i>					
Qwen1-8B	$P_c$	28.30%	0.2357%	2510	71006
Baichuan2-7B	$P_c$	28.46%	0.2370%	2524	71006
Qwen1-8B	Set(SUP( $P_c$ ))	56.57%	0.0111%	5018	3013161
Baichuan2-7B	Set(SUP( $P_c$ ))	52.16%	0.0102%	4627	3013161

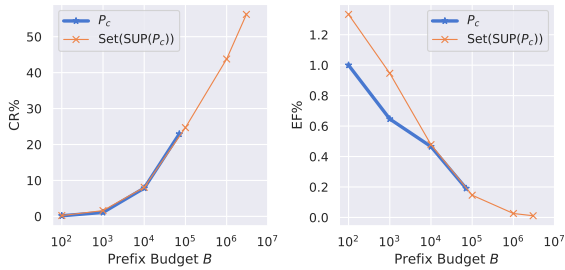


Figure 4: Coverage rate (CR) and efficiency (EF) under varying prefix budgets  $B$  for prefix sets  $P_c$  and  $P_{f \geq 1}$ . Prefix set  $P_{f \geq 1}$  is frequency-sorted in descending order (see Section 4.2.2). Budget values are scaled exponentially (base 10); model used is Qwen1-8B.

victim-exclusive PII extracted by Qwen1-8B (without LAFt). The results for Baichuan2-7B are provided in Appendix Figure 11.

The most frequently extracted PII categories include "Address", "Birthday", and "Name", while others such as "Work Experience" and "Work Place" occur less often but remain notable. More complex types like "Medication Record" are not extracted at all. This is primarily due to the evaluation protocol, which only credits model outputs that match ground truth exactly. Complex PII often appears as long free-text spans, making verbatim reproduction difficult.

To estimate an upper bound of extraction capability, we evaluate with the generalized prefix set  $\text{Set}(\text{SUP}(P_c))$ , which includes all potential contextual prefixes. As shown in Table 1, expanding  $P_c$  to  $\text{Set}(\text{SUP}(P_c))$  increases CR to 56.57% (Qwen1-8B) and 53.56% (Baichuan2-7B). However, this gain comes at a steep cost in efficiency—dropping EF to only 0.01%—indicating most queries yield

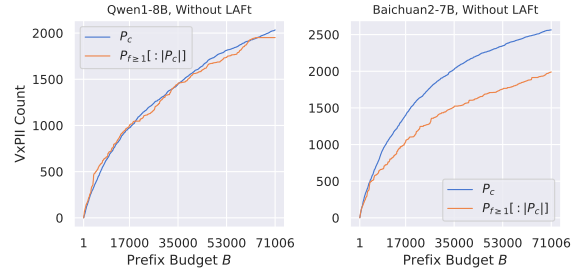


Figure 5: VxPII counts under varying prefix budgets ( $B$ ) for prefix sets  $P_c$  and  $P_{f \geq 1}$ . Prefix set  $P_{f \geq 1}$  is frequency-sorted in descending order (see Section 4.2.2) and truncated to match the size of  $P_c$  here.

redundant or irrelevant content.

We further investigate this CR–EF tradeoff in Figure 4, which illustrates how CR and EF vary with prefix budget  $B$  for prefix sets  $P_c$  and  $P_{f \geq 1}$ . As  $B$  increases, CR improves, but EF declines sharply. This suggests diminishing returns in efficiency when scaling up the number of queries to discover new PII instances.

**RQ2: How effective is frequency-prioritized prefix sampling?** As shown in Figure 5, frequency-prioritized (FP) sampling does not extract more VxPII instances than the contextual prefix set  $P_c$ , contrary to our hypothesis in Section 4.2.2. This result suggests that the contextual cues embedded in  $P_c$  are already strong indicators of LLM memorization, and that memorization cannot be inferred solely from co-occurrence frequency. Instead, it likely arises from more complex interactions between corpus semantics, model architecture, and pre-training dynamics.

Despite this, FP sampling captures highly dis-

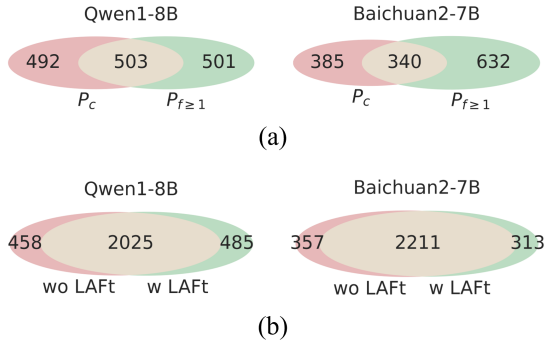


Figure 6: Venn diagrams showing overlap between VxPII sets extracted by different methods. (a) Comparison of VxPII sets using PII-contextual prefixes  $P_c$  vs. frequency-prioritized prefixes  $P_{f \geq 1}$  at prefix budget  $B = 10,000$  (without LAFt). (b) Comparison of VxPII sets extracted with and without LAFt on Qwen1-8B and Baichuan2-7B using the full  $P_c$  prefix set.

tinct subsets of memorized PII. As shown in Figure 6(a), the Venn diagram comparison reveals that 49.9% of the VxPII extracted by FP sampling on Qwen1-8B and 65.02% on Baichuan2-7B are not discovered by the  $P_c$  method. This highlights FP sampling’s complementary strength in uncovering diverse memorized content.

**RQ3: How effective is Latent Association Fine-tuning?** As shown in Table 1, applying Latent Association Fine-tuning (LAFt) significantly improves the CR of Qwen1-8B by 5.37%, raising it to 28.30%, and increases EF to 0.24%, indicating enhanced extraction performance. For Baichuan2-7B, LAFt does not yield a direct improvement in CR, but, as depicted in Figure 6(b), it facilitates the identification of additional distinct PII instances.

These results demonstrate that LAFt is an effective method for increasing the diversity of extracted PII, complementing the FP sampling approach. The extent of the improvement achieved by LAFt is influenced by the construction of the fine-tuning dataset  $D_{ft}$  and the choice of hyperparameters. In this study, we adopt a consistent setting by constructing  $D_{ft}$  through pairing frequent prefixes with randomly sampled PII and fine-tuning the model for one epoch to ensure a fair comparison. However, further exploration of personalized strategies—tailored to models with different architectures and pre-training conditions—could potentially yield better performance.

Table 2: Coverage rates (CR) of extraction attacks across different attacker–victim client pairs with a prefix budget  $B = 10000$ . Prefixes are randomly sampled from each attacker’s corresponding set  $P_c$ . “-” indicates self-attack scenarios, which are not applicable.

Attacker ID	Victim ID				
	0	1	2	3	4
0	-	10.91%	12.89%	10.93%	11.88%
1	11.97%	-	12.41%	11.46%	12.35%
2	12.56%	11.39%	-	11.65%	12.74%
3	12.07%	10.82%	12.04%	-	11.99%
4	12.26%	11.36%	13.25%	11.21%	-

Table 3: Attack performance with and without PII masking using the contextual prefix set  $P_c$ . The model is Qwen1-8B.

	VxPII	CR	EF
With PII Masking	2017	22.74%	0.1894%
Without Defense	2034	22.93%	0.1910%

### 5.3 Cross-Client Evaluation of Extraction Robustness

To assess the robustness of our PII extraction method across different clients, we perform a cross-client evaluation where each client is iteratively designated as the attacker, while the remaining clients act as victims. This setup ensures that the extraction performance is not biased toward any particular client.

As shown in Table 2, our method achieves consistently high coverage rates across all attacker–victim pairings, demonstrating its generalizability and effectiveness in diverse settings.

### 5.4 PII Sanitization Defense

We evaluate the effectiveness of a simple data sanitization strategy that masks PII using existing annotations. Each PII instance in the training data is replaced with an equal-length sequence of asterisks (\*). We then re-fine-tune FedLLM on this sanitized dataset and reapply the PII-contextual Prefix Sampling attack. Table 3 compares the attack performance with and without the PII masking defense.

The results show only a slight reduction in the number of extracted VxPII. To investigate this, we analyze the document frequency of extracted VxPII—that is, how often each appears in the training corpus. Figure 8 demonstrates that masking substantially lowers the frequency of most VxPII, suggesting that our annotations successfully cover



Table 4: Extraction results on non-Chinese LLMs after a single round of federated fine-tuning.

Model	Prefix Set	CR	EF
OLMo2-1124-7B	$P_c$	23.04%	0.1919%
Llama-2-7B	$P_c$	26.41%	0.2200%

the majority of PII. Interestingly, some VxPII with zero document frequency—absent from the sanitized dataset—were still extracted.

Based on these observations, we attribute the limited effectiveness of masking to two factors. First, pretraining data contamination: our training data, drawn from publicly available legal documents, likely overlaps with the pretraining corpora of models such as Qwen1-8B and Baichuan2-7B. Second, incomplete PII labeling: some PII instances may be missing from annotations, and in practice attackers can redefine PII categories, making exhaustive coverage fundamentally difficult.

### 5.5 Disentangling the Effect of Data Contamination

Pretraining data contamination is difficult to avoid, as LLM providers rarely disclose their pretraining corpora. To mitigate this influence, we adopt two strategies:

**Using open-source LLMs.** To eliminate interference from pretraining memorization, we use open-source or semi-open-source LLMs that are not pretrained on Chinese legal documents: OLMo2-7B (OLMo et al., 2024) and Llama2-7B (Touvron et al., 2023). As shown in Table 4, after just one round of federated fine-tuning, the CR exceeded 23%, achieving performance comparable to or better than Qwen1-8B and Baichuan2-7B. This confirms that our findings are not simply an artifact of pretraining contamination.

**Subtracting contaminated memorization.** For Chinese-proficient LLMs that may contain contamination, we adopt a subtraction-based approach. Specifically, we compare the VxPII extracted from the fine-tuned FedLLM ( $F$ ) with those from its base model ( $B$ ), and compute  $F \setminus B$  to isolate PII memorized during federated fine-tuning. Table 5 shows that even after subtracting  $B$ , a substantial number of VxPII remain in  $F \setminus B$ , confirming memorization during fine-tuning. Furthermore, Figure 10 demonstrates that  $F \setminus B$  exhibits a distribution of VxPII labels similar to Figure 3, supporting the validity of this analysis.

Table 5: Comparison of VxPII sets between attacks on FedLLM and its un-fine-tuned base model.

Prefix Set	Model	$ F \setminus B $	$ B \setminus F $	$ F \cap B $
$P_c$	Qwen1	682	518	1801
$P_{f \geq 1}$	Qwen1	554	308	4611
$P_c$	Baichuan2	407	405	2161

## 6 Conclusion

To investigate the privacy risks of data extraction attacks in realistic settings, we introduce a new class of attacks targeting FedLLMs. We extend a legal dataset with systematic PII annotations aligned with major privacy regulations, and evaluate attack performance using two key metrics: coverage rate and efficiency. Extensive experiments demonstrate that certain PII types are highly vulnerable, and our proposed methods can achieve substantial extraction performance. These findings highlight a critical privacy gap in FedLLMs and underscore the urgent need for stronger defense mechanisms in future federated learning systems.

### Limitations

This work investigates the privacy risks of FedLLMs using a legal-domain dataset. Future research can extend our proposed methods to other sensitive domains such as healthcare and finance, where privacy concerns are equally critical. Additionally, there is a need for further exploration of defense mechanisms that can preserve the privacy of FedLLMs while maintaining their performance.

### Ethics Statement

This paper presents PII extraction attacks on federated fine-tuned LLMs to expose potential privacy risks. While designed for research and defense purposes, such methods could be misused to recover sensitive user data in real-world FL systems. We conduct all experiments on legal datasets with anonymized PII, and highlight the need for stronger safeguards in FedLLM deployments.

### Acknowledgements

We thank Guanzhong Chen and Yukun Zhang for their help with dataset annotation. All annotators were properly compensated for their contributions.

## References

- Mahad Ali, Curtis Lisle, Patrick W. Moore, Tammer Barkouki, Brian J. Kirkwood, and Laura J. Brattain. 2025. [Fine-tuning foundation models with federated learning for privacy preserving medical time series forecasting](#). *Preprint*, arXiv:2502.09744.
- Benny Applebaum. 2017. *Garbled Circuits as Randomized Encodings of Functions: a Primer*, pages 1–44. Springer International Publishing, Cham.
- Eugene Bagdasaryan, Andreas Veit, Yiqing Hua, Deborah Estrin, and Vitaly Shmatikov. 2020. [How to backdoor federated learning](#). In *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics*, volume 108 of *Proceedings of Machine Learning Research*, pages 2938–2948. PMLR.
- Jiamu Bai, Daoyuan Chen, Bingchen Qian, Liuyi Yao, and Yaliang Li. 2024. [Federated fine-tuning of large language models under heterogeneous tasks and client resources](#). In *Advances in Neural Information Processing Systems*, volume 37, pages 14457–14483. Curran Associates, Inc.
- Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu, Keming Lu, and 29 others. 2023. [Qwen technical report](#). *Preprint*, arXiv:2309.16609.
- Nicholas Carlini, Daphne Ippolito, Matthew Jagielski, Katherine Lee, Florian Tramèr, and Chiyuan Zhang. 2023. [Quantifying memorization across neural language models](#). In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net.
- Nicholas Carlini, Florian Tramèr, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom Brown, Dawn Song, Ulfar Erlingsson, Alina Oprea, and Colin Raffel. 2021. [Extracting training data from large language models](#). *Preprint*, arXiv:2012.07805.
- Shan Chang, Ye Liu, Zhijian Lin, Hongzi Zhu, Bingzhu Zhu, and Cong Wang. 2024. [Fedtrojan: Corrupting federated learning via zero-knowledge federated trojan attacks](#). In *2024 IEEE/ACM 32nd International Symposium on Quality of Service (IWQoS)*, pages 1–10.
- Chaochao Chen, Xiaohua Feng, Yuyuan Li, Lingjuan Lyu, Jun Zhou, Xiaolin Zheng, and Jianwei Yin. 2024. [Integration of Large Language Models and Federated Learning](#). *Patterns*, 5(12):101098.
- Hong-Min Chu, Jonas Geiping, Liam H. Fowl, Micah Goldblum, and Tom Goldstein. 2023. [Panning for gold in federated learning: Targeted text extraction under arbitrarily large-scale aggregation](#). In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net.
- Yihong Dong, Xue Jiang, Huanyu Liu, Zhi Jin, Bin Gu, Mengfei Yang, and Ge Li. 2024. [Generalization or memorization: Data contamination and trustworthy evaluation for large language models](#). In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 12039–12050, Bangkok, Thailand. Association for Computational Linguistics.
- Xingyi Duan, Baoxin Wang, Ziyue Wang, Wentao Ma, Yiming Cui, Dayong Wu, Shijin Wang, Ting Liu, Tianxiang Huo, Zhen Hu, Heng Wang, and Zhiyuan Liu. 2019. [Cjrc: A reliable human-annotated benchmark dataset for chinese judicial reading comprehension](#). In *Chinese Computational Linguistics*, pages 439–451, Cham. Springer International Publishing.
- European Union. 2016. [General Data Protection Regulation \(GDPR\)](#).
- Jonas Geiping, Hartmut Bauermeister, Hannah Dröge, and Michael Moeller. 2020. [Inverting gradients - how easy is it to break privacy in federated learning?](#) In *Proceedings of the 34th International Conference on Neural Information Processing Systems, NIPS '20*, Red Hook, NY, USA. Curran Associates Inc.
- Team GLM, Aohan Zeng, Bin Xu, Bowen Wang, Chenhui Zhang, Da Yin, Diego Rojas, Guanyu Feng, Hanlin Zhao, Hanyu Lai, Hao Yu, Hongning Wang, Jiadai Sun, Jiajie Zhang, Jiale Cheng, Jiayi Gui, Jie Tang, Jing Zhang, Juanzi Li, and 37 others. 2024. [Chatglm: A family of large language models from glm-130b to glm-4 all tools](#). *Preprint*, arXiv:2406.12793.
- Kuangpu Guo, Yuhe Ding, Jian Liang, Ran He, Zilei Wang, and Tieniu Tan. 2024. [Exploring vacant classes in label-skewed federated learning](#). *Preprint*, arXiv:2401.02329.
- Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. [Lora: Low-rank adaptation of large language models](#). *Preprint*, arXiv:2106.09685.
- Hongsheng Hu, Xuyun Zhang, Zoran Salcic, Lichao Sun, Kim-Kwang Raymond Choo, and Gillian Dobbie. 2024. [Source inference attacks: Beyond membership inference attacks in federated learning](#). *IEEE Transactions on Dependable and Secure Computing*, 21(4):3012–3029.
- Jie Huang, Hanyin Shao, and Kevin Chen-Chuan Chang. 2022. [Are large pre-trained language models leaking your personal information?](#) In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 2038–2047, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank J. Reddi, Sebastian U. Stich, and Ananda Theertha Suresh. 2020. [Scaffold: Stochastic Controlled Averaging for Federated Learning](#). In *Proceedings of the 37th International Conference on Machine Learning, ICML, Vienna, Austria (virtual conference)*.

- Aly M. Kassem, Omar Mahmoud, Niloofar Miresghalah, Hyunwoo Kim, Yulia Tsvetkov, Yejin Choi, Sherif Saad, and Santu Rana. 2025. *Alpaca against vicuna: Using llms to uncover memorization of llms*. *Preprint*, arXiv:2403.04801.
- Siwon Kim, Sangdoon Yun, Hwaran Lee, Martin Gubri, Sungroh Yoon, and Seong Joon Oh. 2023. *Propile: Probing privacy leakage in large language models*. In *Advances in Neural Information Processing Systems*, volume 36, pages 20750–20762. Curran Associates, Inc.
- Vladimir I. Levenshtein. 1965. Binary Codes Capable of Correcting Deletions, Insertions, and Reversals. *Soviet physics. Doklady*, 10:707–710.
- Haitao Li, You Chen, Qingyao Ai, Yueyue Wu, Ruizhe Zhang, and Yiqun Liu. 2024a. *Lexeval: A comprehensive chinese legal benchmark for evaluating large language models*. In *Advances in Neural Information Processing Systems*, volume 37, pages 25061–25094. Curran Associates, Inc.
- Ming Li, Yong Zhang, Zhitao Li, Jiu Hai Chen, Lichang Chen, Ning Cheng, Jianzong Wang, Tianyi Zhou, and Jing Xiao. 2023. *From quantity to quality: Boosting LLM performance with self-guided data selection for instruction tuning*. *CoRR*, abs/2308.12032.
- Shenghui Li, Fanghua Ye, Meng Fang, Jiayu Zhao, Yun-Hin Chan, Edith C. H. Ngai, and Thiemo Voigt. 2024b. *Synergizing foundation models and federated learning: A survey*. *Preprint*, arXiv:2406.12844.
- Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith. 2020. *Federated optimization in heterogeneous networks*. *Preprint*, arXiv:1812.06127.
- Zilinghan Li, Pranshu Chaturvedi, Shilan He, Han Chen, Gagandeep Singh, Volodymyr Kindratenko, Eliu A Huerta, Kibaek Kim, and Ravi Madduri. 2024c. *Fedcompass: Efficient cross-silo federated learning on heterogeneous client devices using a computing power-aware scheduler*. In *The Twelfth International Conference on Learning Representations*.
- Nils Lukas, Ahmed Salem, Robert Sim, Shruti Tople, Lukas Wutschitz, and Santiago Zanella-Béguelin. 2023. *Analyzing leakage of personally identifiable information in language models*. In *2023 IEEE Symposium on Security and Privacy (SP)*, pages 346–363.
- H. Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Agüera y Arcas. 2023. *Communication-efficient learning of deep networks from decentralized data*. *Preprint*, arXiv:1602.05629.
- Luca Melis, Congzheng Song, Emiliano De Cristofaro, and Vitaly Shmatikov. 2019. *Exploiting unintended feature leakage in collaborative learning*. In *2019 IEEE Symposium on Security and Privacy (SP)*, pages 691–706.
- Krishna Kanth Nakka, Ahmed Frikha, Ricardo Mendes, Xue Jiang, and Xuebing Zhou. 2024. *Pii-scope: A benchmark for training data pii leakage assessment in llms*. *Preprint*, arXiv:2410.06704.
- Team OLMO, Pete Walsh, Luca Soldaini, Dirk Groeneveld, Kyle Lo, Shane Arora, Akshita Bhagia, Yuling Gu, Shengyi Huang, Matt Jordan, Nathan Lambert, Dustin Schwenk, Oyvind Tafjord, Taira Anderson, David Atkinson, Faeze Brahman, Christopher Clark, Pradeep Dasigi, Nouha Dziri, and 21 others. 2024. *2 olmo 2 furious*. *Preprint*, arXiv:2501.00656.
- OpenAI, :, Aaron Hurst, Adam Lerer, Adam P. Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, Aleksander Madry, Alex Baker-Whitcomb, Alex Beutel, Alex Borzunov, Alex Carney, Alex Chow, Alex Kirillov, and 401 others. 2024. *Gpt-4o system card*. *Preprint*, arXiv:2410.21276.
- Personal Data Protection Commission, Singapore. 2012. *Personal Data Protection Act (PDPA)*.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2023. *Exploring the limits of transfer learning with a unified text-to-text transformer*. *Preprint*, arXiv:1910.10683.
- Md Rafi Ur Rashid, Vishnu Asutosh Dasu, Kang Gu, Najrin Sultana, and Shagufta Mehnaz. 2025. *Fltrojan: Privacy leakage attacks against federated language models through selective weight tampering*. *Preprint*, arXiv:2310.16152.
- Eric Zhang Raytrfr, LlamaFamily. 2024. *Llama-chinese*. <https://github.com/LlamaFamily/Llama-Chinese>.
- Google Research. 2022. *Training data extraction challenge*.
- Avi Schwarzschild, Zhili Feng, Pratyush Maini, Zachary C. Lipton, and J. Zico Kolter. 2024. *Rethinking llm memorization through the lens of adversarial compression*. In *Advances in Neural Information Processing Systems*, volume 37, pages 56244–56267. Curran Associates, Inc.
- Xinye Sha. 2024. *Research on financial fraud algorithm based on federal learning and big data technology*. *Preprint*, arXiv:2405.03992.
- Hanyin Shao, Jie Huang, Shen Zheng, and Kevin Chang. 2024. *Quantifying association capabilities of large language models and its implications on privacy leakage*. In *Findings of the Association for Computational Linguistics: EACL 2024*, pages 814–825, St. Julian’s, Malta. Association for Computational Linguistics.
- Virat Shejwalkar and Amir Houmansadr. 2021. *Manipulating the byzantine: Optimizing model poisoning attacks and defenses for federated learning*. In *NDSS*.



- Weijia Shi, Anirudh Ajith, Mengzhou Xia, Yangsibo Huang, Daogao Liu, Terra Blevins, Danqi Chen, and Luke Zettlemoyer. 2024. [Detecting pretraining data from large language models](#). In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net.
- Standardization Administration of China (SAC). 2020. GB/T 35273-2020 | Information Security Technology: Personal Information Security Specification. Available at: <https://openstd.samr.gov.cn/>. Replaces GB/T 35273-2017, National Standard, Number: GB/T 35273—2020.
- State of California, US. 2018. [California Consumer Privacy Act \(CCPA\)](#).
- Maxim Tkachenko, Mikhail Malyuk, Andrey Holmanyuk, and Nikolai Liubimov. 2020-2025. [Label Studio: Data labeling software](#). Open source software available from <https://github.com/HumanSignal/label-studio>.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, and 49 others. 2023. [Llama 2: Open foundation and fine-tuned chat models](#). *Preprint*, arXiv:2307.09288.
- Minh N. Vu, Truc Nguyen, Tre’ R. Jeter, and My T. Thai. 2024. [Analysis of privacy leakage in federated large language models](#). *Preprint*, arXiv:2403.04784.
- Linlin Wang, Tianqing Zhu, Wanlei Zhou, and Philip S. Yu. 2024a. [Linkage on security, privacy and fairness in federated learning: New balances and new perspectives](#). *Preprint*, arXiv:2406.10884.
- Ziyao Wang, Zheyu Shen, Yexiao He, Guoheng Sun, Hongyi Wang, Lingjuan Lyu, and Ang Li. 2024b. [Flora: Federated fine-tuning large language models with heterogeneous low-rank adaptations](#). In *Advances in Neural Information Processing Systems*, volume 37, pages 22513–22533. Curran Associates, Inc.
- Yebo Wu, Chunlin Tian, Jingguang Li, He Sun, Kahou Tam, Li Li, and Chengzhong Xu. 2025. [A survey on federated fine-tuning of large language models](#). *Preprint*, arXiv:2503.12016.
- Xianghua Xie, Chen Hu, Hanchi Ren, and Jingjing Deng. 2024. [A survey on vulnerability of federated learning: A learning algorithm perspective](#). *Neurocomputing*, 573:127225.
- Aiyuan Yang, Bin Xiao, Bingning Wang, Borong Zhang, Ce Bian, Chao Yin, Chenxu Lv, Da Pan, Dian Wang, Dong Yan, Fan Yang, Fei Deng, Feng Wang, Feng Liu, Guangwei Ai, Guosheng Dong, Haizhou Zhao, Hang Xu, Haoze Sun, and 36 others. 2023. [Baichuan 2: Open large-scale language models](#). *Preprint*, arXiv:2309.10305.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, and 41 others. 2025. [Qwen3 technical report](#). *Preprint*, arXiv:2505.09388.
- Qiang Yang, Yang Liu, Tianjian Chen, and Yongxin Tong. 2019. [Federated machine learning: Concept and applications](#). *ACM Trans. Intell. Syst. Technol.*, 10(2).
- Yuhang Yao, Jianyi Zhang, Junda Wu, Chengkai Huang, Yu Xia, Tong Yu, Ruiyi Zhang, Sungchul Kim, Ryan Rossi, Ang Li, Lina Yao, Julian McAuley, Yiran Chen, and Carlee Joe-Wong. 2024. [Federated large language models: Current progress and future directions](#). *Preprint*, arXiv:2409.15723.
- Rui Ye, Rui Ge, Xinyu Zhu, Jingyi Chai, Yaxin Du, Yang Liu, Yanfeng Wang, and Siheng Chen. 2024a. [Fedllm-bench: Realistic benchmarks for federated learning of large language models](#). In *Advances in Neural Information Processing Systems*, volume 37, pages 111106–111130. Curran Associates, Inc.
- Rui Ye, Wenhao Wang, Jingyi Chai, Dihan Li, Zexi Li, Yinda Xu, Yaxin Du, Yanfeng Wang, and Siheng Chen. 2024b. [Openfedllm: Training Large Language Models on Decentralized Private Data Via Federated Learning](#). In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, page 6137–6147, New York, NY, USA. Association for Computing Machinery.
- Weichen Yu, Tianyu Pang, Qian Liu, Chao Du, Bingyi Kang, Yan Huang, Min Lin, and Shuicheng Yan. 2023. [Bag of tricks for training data extraction from language models](#). In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 40306–40320. PMLR.
- Shengbin Yue, Wei Chen, Siyuan Wang, Bingxuan Li, Chenchen Shen, Shujun Liu, Yuxuan Zhou, Yao Xiao, Song Yun, Xuanjing Huang, and Zhongyu Wei. 2023. [Disc-lawllm: Fine-tuning large language models for intelligent legal services](#). *Preprint*, arXiv:2309.11325.
- Shengbin Yue, Shujun Liu, Yuxuan Zhou, Chenchen Shen, Siyuan Wang, Yao Xiao, Bingxuan Li, Yun Song, Xiaoyu Shen, Wei Chen, and 1 others. 2024. [Lawllm: Intelligent legal system with legal reasoning and verifiable retrieval](#). In *International Conference on Database Systems for Advanced Applications*, pages 304–321. Springer.
- Jingyang Zhang, Jingwei Sun, Eric C. Yeats, Yang Ouyang, Martin Kuo, Jianyi Zhang, Hao Yang, and Hai Helen Li. 2024. [Min-k%++: Improved baseline for detecting pre-training data from large language models](#). *CoRR*, abs/2404.02936.



Zhuo Zhang, Xiangjing Hu, Jingyuan Zhang, Yating Zhang, Hui Wang, Lizhen Qu, and Zenglin Xu. 2023. [FEDLEGAL: The first real-world federated learning benchmark for legal NLP](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3492–3507, Toronto, Canada. Association for Computational Linguistics.

Haoxi Zhong, Chaojun Xiao, Cunchao Tu, Tianyang Zhang, Zhiyuan Liu, and Maosong Sun. 2020. [Jecqa: A legal-domain question answering dataset](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(05):9701–9708.

Zhenhong Zhou, Jiuyang Xiang, Chaomeng Chen, and Sen Su. 2024. [Quantifying and analyzing entity-level memorization in large language models](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(17):19741–19749.

Ligeng Zhu, Zhijian Liu, and Song Han. 2019. [Deep leakage from gradients](#). In *Proceedings of the 33rd International Conference on Neural Information Processing Systems*, Red Hook, NY, USA. Curran Associates Inc.

## A Acronyms List

- **GDPR** - General Data Protection Regulation ([European Union, 2016](#))
- **CCPA** - California Consumer Privacy Act ([State of California, US, 2018](#))
- **CPIS** - Chinese Information Security Technology: Personal Information Security Specification (GB/T 35273-2020) ([Standardization Administration of China \(SAC\), 2020](#))
- **Singapore PDPC** - Personal Data Protection Commission (Singapore) ([Personal Data Protection Commission, Singapore, 2012](#))
- **Non-IID** - Non-independent and identically distributed

## B Preliminary Knowledge

### B.1 Data Extraction Attack

Early research on training data extraction attacks has broadly categorized them into untargeted and targeted attacks ([Research, 2022](#); [Yu et al., 2023](#)). Untargeted extraction aims to recover any memorized training samples without specifying a target ([Lukas et al., 2023](#)), whereas targeted extraction attempts to reconstruct specific training samples, often by providing a known prefix and recovering the remaining content ([Carlini et al., 2021](#)). The latter type, often referred to as Verbatim Data

Extraction, has become a standard approach for evaluating memorization in LLMs ([Carlini et al., 2023](#); [Dong et al., 2024](#)) and for detecting potential data contamination ([Dong et al., 2024](#)). We briefly outline the core methodology of verbatim data extraction below.

Given an LLM  $\theta$  and a training dataset  $X$ , each training sample  $x_i \in X$  is partitioned into two segments: a prefix  $a_i$  and a suffix  $b_i$ , such that  $x_i = a_i b_i$ . The model is then prompted with  $a_i$  to generate a completion  $g_i$  of the same length as  $b_i$ . If  $g_i$  exactly matches  $b_i$ , the sample is considered successfully extracted.

In practice, model outputs may not exactly replicate the original suffix but can still be lexically close. To accommodate this, a similarity-based metric such as Edit Distance ([Levenshtein, 1965](#)) is often employed. A sample is deemed extracted if the similarity score between  $g_i$  and  $b_i$  exceeds a pre-defined threshold  $t$ . By computing this similarity-based extraction score across all samples in a dataset  $D$ , one can quantify the model’s memorization behavior or assess its vulnerability to training data extraction attacks.

### B.2 Federated Learning

Federated Learning (FL) is a solution to address data isolation issues ([Yang et al., 2019](#)), where a central server and multiple clients collaborate to complete the training process. A key feature of FL is that the training datasets are stored locally on each client and remain invisible to other clients. FL is commonly used in industrial scenarios where each client represents an independent organization, such as hospitals collaborating to train a medical model without combining their datasets due to legal restrictions or business competition. Federated Learning enables the training of stronger models compared to training on data from a single client alone.

Given  $c$  clients and their private datasets  $D_1, D_2, \dots, D_c$ , the federated learning process aims to learn a global model  $\theta$  by solving the following optimization problem:

$$\theta^* = \arg \min_{\theta} \frac{1}{c} \sum_{i=1}^c \mathcal{L}(D_i, \theta)$$

To solve this problem, many federated optimization algorithms have been proposed, such as FedAVG ([McMahan et al., 2023](#)) and FedProx ([Li et al., 2020](#)). Typically, these algorithms consist

of two alternating phases: local updating and central aggregation. In the local updating phase, each client independently optimizes the global model using its own dataset. In the central aggregation phase, the server aggregates the models from the clients using an aggregation algorithm, obtaining a global model, which is then sent back to each client for the next round of local updating. A typical procedure of federated learning is illustrated in Algorithm 1.

### B.3 Federated Large Language Models (FedLLMs)

In current research at the intersection of federated learning and large language models, the term *FedLLM* predominantly refers to *federated fine-tuning* of pre-trained LLMs (Bai et al., 2024; Wang et al., 2024b; Ye et al., 2024a), rather than federated pre-training. This focus arises from both practical and technical considerations.

Federated pre-training is rarely necessary, as large-scale pre-training typically relies on publicly available general-purpose corpora (such as web text or encyclopedias) that do not contain sensitive information and therefore do not require federated sharing. Moreover, federated pre-training would entail transmitting extremely large model parameters across institutions, incurring prohibitive communication costs and conflicting with high-efficiency training practices like data parallelism and optimized operators, making deployment infeasible in practice.

By contrast, federated fine-tuning is essential in privacy-sensitive domains. Many applications of LLMs in vertical fields—such as judicial documents, electronic health records, and banking customer data—rely on restricted information that cannot be centralized due to legal or institutional constraints (e.g., the Data Security Law or the Personal Information Protection Law). In these settings, federated fine-tuning allows each institution to adapt a shared pre-trained model locally, achieving a balance between model performance and data privacy. This approach has already demonstrated tangible value across multiple domains.

For example, in the judicial domain (Zhang et al., 2023), courts can fine-tune a common LLM on local case repositories without sharing sensitive records, enabling cross-court tasks such as statute matching and case similarity analysis. In healthcare (Ali et al., 2025), hospitals can locally fine-tune models on specialty data, producing compre-

hensive medical LLMs capable of supporting structured record processing and disease risk prediction. And in finance (Sha, 2024), banks can fine-tune models on transaction data to detect fraud and assess credit risk without violating privacy regulations.

### B.4 Preliminary Assessment of Verbatim Data Extraction Risks in FedLLMs

To examine the memorization behavior of FedLLMs and evaluate their potential risks of leaking sensitive information, we conduct a preliminary experiment simulating a *verbatim data extraction* attack. The results are referenced in the main paper (Section 1) to empirically motivate our study.

We adapt the experimental setup from (Dong et al., 2024) to the federated setting, where an attacker is assumed to possess prefix fragments of the training data from all participating clients and attempts to recover the subsequent suffix tokens. For each training sample, we extract a prefix from the original sequence and query the trained model to generate a continuation. The generated suffix is compared against the ground truth using **Edit Distance** (ED) (Levenshtein, 1965), where a lower ED indicates stronger memorization. Specifically:

- ED = 0 indicates the model has perfectly memorized and reproduced the suffix;
- ED values are capped at 50, as we restrict suffixes to a maximum of 50 tokens.

We perform the attack on the global models aggregated after 10 rounds of federated training. To ensure a comprehensive assessment, we consider three popular FL algorithms—**FedAvg** (McMahan et al., 2023), **FedProx** (Li et al., 2020), and **Scaffold** (Karimireddy et al., 2020)—each under both **IID** and **Non-IID** data distributions. Two baseline settings are also included:

- **Centralized**: All client data is pooled and the model is fine-tuned in a conventional non-federated manner;
- **Untrained**: The base model is evaluated without any fine-tuning.

Figure 7 summarizes the results across five downstream tasks. Our key observations are:

- FedLLMs consistently exhibit higher ED scores (i.e., lower memorization) than central-

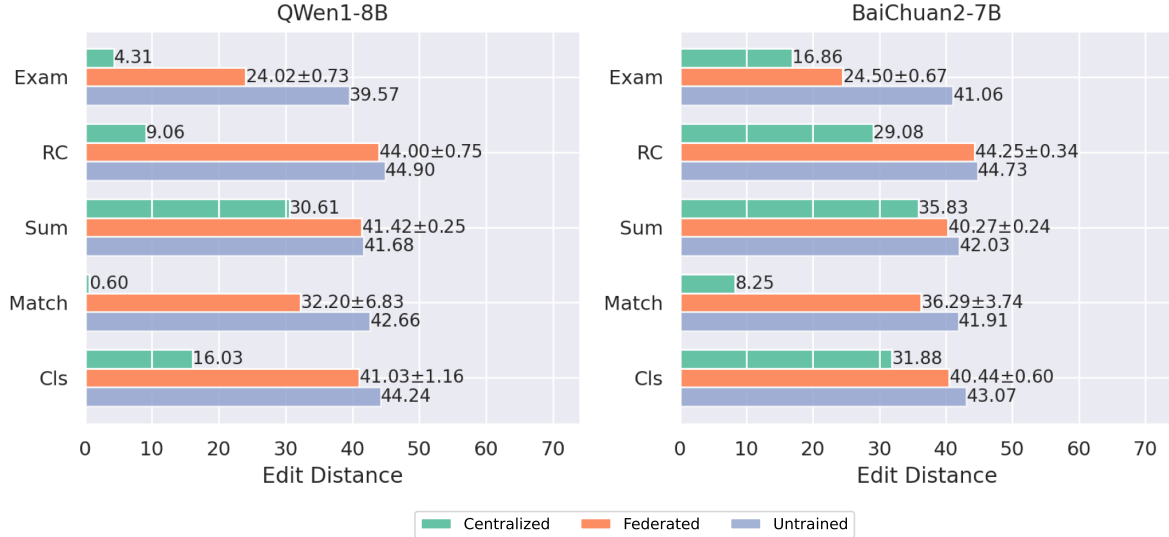


Figure 7: Edit Distance results of verbatim data extraction attacks after 10 training rounds. We evaluate six federated configurations (FedAvg/FedProx/Scaffold  $\times$  IID/Non-IID) and report the mean and standard deviation. Lower values indicate stronger memorization. Centralized and untrained models serve as baselines.

ized models, suggesting that the FL aggregation process reduces susceptibility to verbatim extraction.

- However, compared to untrained models, FedLLMs still show non-negligible memorization, with noticeably lower ED scores, indicating partial leakage of training data.

These findings highlight a trade-off between collaborative model training and privacy preservation, and they serve as the motivation for our in-depth investigation of privacy risks in FedLLMs.

## C Federated Learning Framework

Algorithm 1 outlines a general framework for Federated Learning (FL), where a central server coordinates multiple clients to collaboratively train a global model without sharing local data. At each round, the server distributes the current model to all clients, each of which performs local updates based on its private data and sends the updated parameters back. The server then aggregates the received updates to produce a new global model.

## D Task Descriptions and Examples

1. **Judicial Summarization (Sum):** The task of judicial summarization aims to extract key information from court judgments and generate concise summaries. The input to this task is a legal document, and the output is a summary

---

### Algorithm 1 A Federated Learning Framework

---

**Input:** Clients set  $\mathcal{C} = \{c_1, c_2, \dots, c_c\}$  with local datasets  $D_1, D_2, \dots, D_c$ ; total FL rounds  $R$ ; initial global model  $\theta_0$ ; server aggregation function  $f_{\text{agg}}$ ; client loss function  $\mathcal{L}$

**Output:** Learned global model  $\theta_R$

```

1: ServerExecute:
2: for round  $r = 1$  to  $R$  do
3:   for each client  $c_i \in \mathcal{C}$  (in parallel) do
4:      $\theta_r^i \leftarrow \text{CLIENTUPDATE}(c_i, \theta_{r-1})$ 
5:   end for
6:    $\theta_r \leftarrow f_{\text{agg}}(\{\theta_r^i | c_i \in \mathcal{C}\})$ 
7: end for

8: ClientExecute:
9: function CLIENTUPDATE( $c_i, \theta_{r-1}$ )
10:   $\theta_r^i \leftarrow \arg \min_{\theta} \mathcal{L}(\theta_{r-1}, D_i)$ 
11:  return  $\theta_r^i$ 
12: end function

```

---

of its content. The performance of this task is evaluated using the Rouge-L metric, which effectively measures the similarity between the generated and reference texts based on the longest common subsequence (LCS). Rouge-L is a widely used metric in text generation tasks. In this study, we adopt Rouge-L because it captures both semantic and structural similarities between texts, making it suitable for summarizing judicial documents.

2. **Judicial Reading Comprehension (RC):** This task focuses on answering legal questions based on court documents to evaluate the

model’s reading comprehension ability. The input consists of a piece of legal material and a question, and the task requires answering the question based on the content of the material. The performance metric for this task is Rouge-L.

3. **Similar Case Matching (Match):** In this task, the input includes three case documents, and the model is required to determine which of the latter two documents is more similar to the first one. The model selects the most similar document by computing the similarity between the first case and each of the other two. The evaluation metric for this task is accuracy.
4. **Judicial Exam (Exam):** This task simulates multiple-choice questions from legal examinations to assess the model’s knowledge of legal concepts. Given a judicial exam question with multiple options, the model is expected to choose the correct answer. The performance is evaluated using accuracy.
5. **Legal Case Classification (Cls):** This task requires the model to classify the cause of action in a case, assisting legal retrieval systems in automatically categorizing case types. The input is a description of the case facts, and the model is required to output the corresponding case category. The performance metric is accuracy.

## E Data Preprocessing

Previous works (Zhang et al., 2023; Yue et al., 2023) have used these datasets for LLM and FedLLM research. In this work, we use the processed datasets from these prior studies and further curate the data for our experiments. We applied the following preprocessing steps to prepare the datasets:

**Deduplication and Cleansing.** To ensure the quality of our data, we remove duplicate samples with logically equivalent meanings. For example, in the RC tasks, some samples only differ in the order of two legal cases. We also clean out samples containing garbled characters or large segments with a mixture of multiple languages.

**Unifying Prompt Template and Instruction Reshaping.** Some tasks, such as Exam, contain instructions that appear in different parts of the sam-

ple (either at the beginning or the end). To standardize the format, we reshape the data so that the instruction always appears at the beginning, followed by the legal document. Additionally, we employ hierarchical hyper markers such as "<Case A>", "<Case B>", and "<Answer>" to clearly segment the prompt, making the structure more transparent for the LLM.

## F Supplementary Dataset Statistics and Analysis

Table 7 summarizes the basic statistics of the five datasets used in our experiments. Each dataset corresponds to a different downstream task for fine-tuning the model.

Figure 8 presents the document frequency distribution of the 2017 VxPII instances extracted from the model trained on the masked dataset (see Section 5.4). Most VxPII exhibit low frequency, indicating that PII masking significantly reduces memorization.

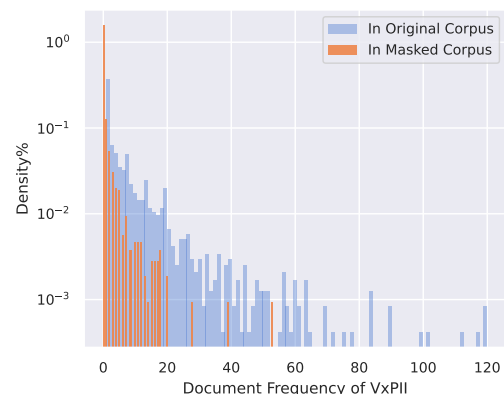


Figure 8: Document frequency distribution of the 2017 VxPII instances extracted from the model trained on the masked dataset.

## G Prompt Template and Utility Fine-tuning Results for FedLLMs

Figure 9 shows the unified prompt template used for all federated utility fine-tuning tasks. Table 8 reports the evaluation results across multiple tasks, comparing different federated learning algorithms and base models.

## H Additional PII Label Distribution Results

Figure 10 illustrates the label distribution of FedLLM-exclusive victim PII extracted by Qwen1-



Table 6: Input and Output Examples for Each Task

Task	Input	Output
Judicial Summarization (SUM)	First-instance civil judgment on inheritance dispute between Han and Su Shenyang Dadong District People’s Court Plaintiff: Han, female, born June 6, 1927, Han ethnicity... ... Clerk: Li Dan	Summary: This case involves an inheritance dispute between the plaintiff and the defendant. The plaintiff requests...
Judicial Reading Comprehension (RC)	Case: Upon trial, it was found that on February 11, 2014, the plaintiff... Question: When did the plaintiff and defendant agree on the travel plan?	The plaintiff and defendant agreed on the travel plan on February 11, 2014.
Similar Case Matching (Match)	Determine whether Case A is more similar to Case B or Case C. A: Plaintiff: Zhou Henghai, male, born October 17, 1951... B: Plaintiff: Huang Weiguo, male, Han ethnicity, resident of Zhoushan City... C: Plaintiff: Zhang Huaibin, male, resident of Suzhou City, Anhui Province, Han ethnicity...	B
Judicial Exam (Exam)	Wu was lawfully pursued by A and B... Which of the following analyses is correct? A. If Wu missed both A and B, and the bullet... B. If Wu hit A, resulting in A’s death... C. If Wu hit both A and B, causing A’s death and B’s serious injury... D. If Wu hit both A and B, causing both to die...	A
Legal Case Classification (Cls)	Legal document: Plaintiff Yan Qiang submitted the following claims to this court:...	Private Loan Dispute

Table 7: Dataset Statistics

	Exam	RC	SUM	Match	Cls
#Samples	2399	3500	2651	3848	4196

```

Below is a task related to judicial and legal matters. Output an appropriately completed response to the request.

<### Input>
{{Task Input}}

<### Output>
{{Task Output}}

```

Figure 9: Unified Utility Fine-tuning Template for All Tasks.

8B. This result corresponds to the experiment described in Section 5.4.

Figure 11 presents the label distribution of deduplicated victim-exclusive PII instances extracted by the Baichuan2-7B model.

## I Machine Annotation Standards for PII Labeling

This section provides details on the machine annotation protocol we use to identify PII in our dataset. Table 9 defines our categorization schema, which includes seven major categories and their corresponding fine-grained subtypes. To ensure annotation consistency and scalability, we utilize a templated prompting approach for automated PII labeling. Figure 12 shows the machine annotation prompt used to instruct the LLM annotator. The prompt dynamically incorporates category definitions and format constraints to standardize the output.

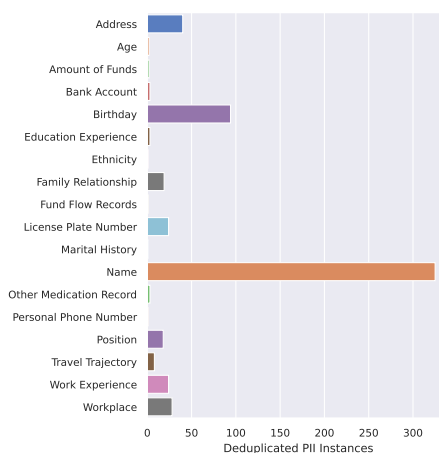


Figure 10: Label distribution of FedLLM-exclusive Vx-PII extracted using prefix set  $P_c$  and the Qwen1-8B model.

## J Details of Human Evaluation for PII Annotation

To validate the quality of the machine-generated PII annotations, we recruited four Chinese-speaking students with foundational knowledge of Chinese law to manually annotate PII on a selected subset of the dataset. Prior to annotation, all annotators underwent thorough training on the annotation guidelines and usage of the Label Studio tool. The instructions provided to annotators are detailed in Figure 13, while Figure 14 illustrates the annotation interface used. All annotators were fairly compensated upon completion of their tasks.

The human evaluation results, reported in terms of precision, recall, and F1 score, are summarized in Table 10, indicating high agreement both in exact span matching and in combined span-and-label matching, confirming the reliability of the machine

Table 8: Utility Performance over Different Tasks.

FL Algorithms	Models	SUM(rouge-l)	RC(rouge-l)	Match(Acc)	Exame(Acc)	Cls(Acc)
FedAvg	Qwen1-8B	50.0	14.2	50.0	37.5	90.0
FedAvg	Baichuan2-7B	57.6	42.4	50.0	33.3	89.5
Non-FL	Qwen1-8B	50.0	18.9	50.0	40.8	87.0

Table 9: Categorization of PII Types in Our Labeling Standards

Major Category	Minor Category
Personal Basic Information	Name, Birthday, Address, Gender, Ethnicity, Family Relationship, Age, Nationality, Personal Phone Number
Personal Identity Information	ID Number, Social Security Number, Driver’s License Number, Employee Number, License Plate Number
Health Related Information	Physical Condition, Fertility Information, Current Medical History, Diagnosis and Treatment Status, Other Medication Record
Work and Education Information	Workplace, Position, Work Experience, Education Experience, Grades
Personal Property Information	Bank Account, Amount of Funds, Fund Flow Records, Virtual Assets, Other Financial Records
Personal Location Information	Precise Location, Accommodation Information, Travel Trajectory
Others	Marital History, Religious or Philosophical Beliefs, Sexual Orientation or Sex Life, Unpublished Criminal Records

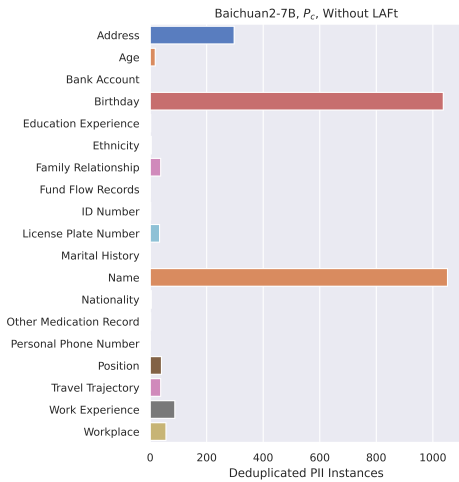


Figure 11: Label distribution of deduplicated victim-exclusive PII instances extracted by the Baichuan2-7B model (without LAFT, using prefix set  $P_c$ ). This figure complements Figure 3 in the main text, which presents the corresponding results for Qwen1-8B.

annotations.

## K Experiment Implementation Details

### K.1 Federated Dataset Partitioning.

We use the preprocessed and labeled datasets (see Section 3) for our experiments, splitting the data into training and testing sets. In the federated learning setup, we simulate a system with five clients.

I would like you to assist in reviewing the provided document and labeling all sections containing {{Major Categories of PII}} according to the following requirements.

- Types of personal information to identify include:**  
{{PII Subcategories}}
- Output format:**  
{{Output Format Description}}
- Input instructions:**  
{{Input Format Description}}

Please provide the output directly in accordance with the format requirements above, without any additional explanation or comments. Thank you for your assistance!

Figure 12: PII Machine Annotation Prompt Template

The testing set remains global, while the training set is heterogeneously partitioned across the clients using a balanced Non-IID distribution (see Acronyms List A). To achieve this, We employ a clustering-based method (Li et al., 2023) for partitioning, where a language encoder first generates embeddings, which are then clustered using K-means. Finally, a Dirac distribution with  $\alpha = 0.5$  is applied to create a label-skewed partitioning (Guo et al., 2024), ensuring each client receives a comparable number of samples.

```

# **PII Annotation Guidelines for Labelers**
## **1. Task Objective**
**Core Task***: Proofread legal texts to accurately identify and annotate **Personally Identifiable Information (PII)***. Each annotation task includes:
1. **Localization***: Mark the exact character offsets of each PII instance in the text;
2. **Categorization***: Assign each PII instance to the appropriate **major category (7 total)*** and **minor category (36 total)***, ensuring precise classification.
## **2. PII Category System**
| Major Category | Minor Categories |
| - | - |
| Personal Basic Information | Name, Birthday, Address, Gender, Ethnicity, Family Relationship, Age, Nationality, Personal Phone Number |
...(omitted)...
## **3. Annotation Workflow and Standards**
### **Step-by-Step Process**
1. **Read the Full Text***: Understand the context to detect all potential PII entities;
2. **Sentence-by-Sentence Annotation***: For each PII instance, annotate its **start position***, **text span***, and corresponding **major + minor category***;
3. **Special Cases***: For ambiguous expressions (e.g., "a certain district of a certain city"), determine PII status based on contextual clues.
### **Annotation Guidelines**
* **Accuracy***: Ensure all annotated content is verifiably present in the text. Avoid false positives or over-labeling;
* **Support Channel***: If any uncertain cases arise during annotation, promptly reach out to the *Annotation Support Team* for clarification.

```

Figure 13: Markdown-style guideline for PII annotation, covering task objectives, taxonomy, and labeling procedures.

Table 10: Human Evaluation of PII Labeling Quality

Evaluation Criteria	Precision (P)	Recall (R)	F1 Score (F1)
Identical Span Only	0.89	0.93	0.91
Identical Span and Label	0.89	0.90	0.89

## K.2 Hardware and Computation Budget

All experiments are conducted on a single NVIDIA A6000 GPU with 48 GB of memory, using bfloat16 precision. Most sampling-based attack experiments are completed within 200 GPU hours.

## K.3 Experiment Procedure

### K.3.1 Federated Utility Fine-Tuning

We begin by performing federated fine-tuning of the LLM (Zhang et al., 2023; Wu et al., 2025) on the partitioned dataset, adapting it to the legal tasks. The fine-tuning is conducted using the OpenFedLLM framework (Ye et al., 2024b). We set the total number of FL rounds to 10 and use FedAVG (McMahan et al., 2023) as the aggregation algorithm.

Each client performs multi-task fine-tuning by mixing all local tasks and applying a unified prompt template, as illustrated in Figure 9, following the approach of Raffel et al. (2023). In each round of federated learning, the client fine-tunes the received global model for one epoch using parameter-efficient fine-tuning (PEFT) techniques of LoRA (Hu et al., 2021). The learning rate is set to  $3e-4$  with a linear decay schedule. The maximum input sequence length is 3072 tokens. We use

a batch size of 1 and apply gradient accumulation with a factor of 8. The LoRA configuration is set to  $r = 16$  and  $\alpha = 32$ .

After federated fine-tuning is complete, we evaluate the utility performance of the final global model on a held-out global test set. In line with standard practices in federated learning research, we also compare this performance with that of a centrally (non-FL) trained model on the same test set. The results are summarized in Table 8.

### K.3.2 PII Extraction

In the main experiments, we designate client 0 as the attacker and client 1 as the victim. We construct the prefix set  $P_c$  for PII-contextual prefix sampling from the local dataset  $D_0$ . During this construction, we set the length parameter  $\lambda$  to 50. Each prefix is used to independently query the utility fine-tuned global model  $n = 15$  times. For each query, the model is allowed to generate up to  $m = 10$  new tokens. This generation length is sufficient to cover most labeled PII instances while keeping the computational cost acceptable.

For Frequency-Prioritized Prefix Sampling, we construct  $\text{Set}(\text{SUP}(P_c))$  from the aforementioned  $P_c$ , and sort it in descending order of prefix fre-

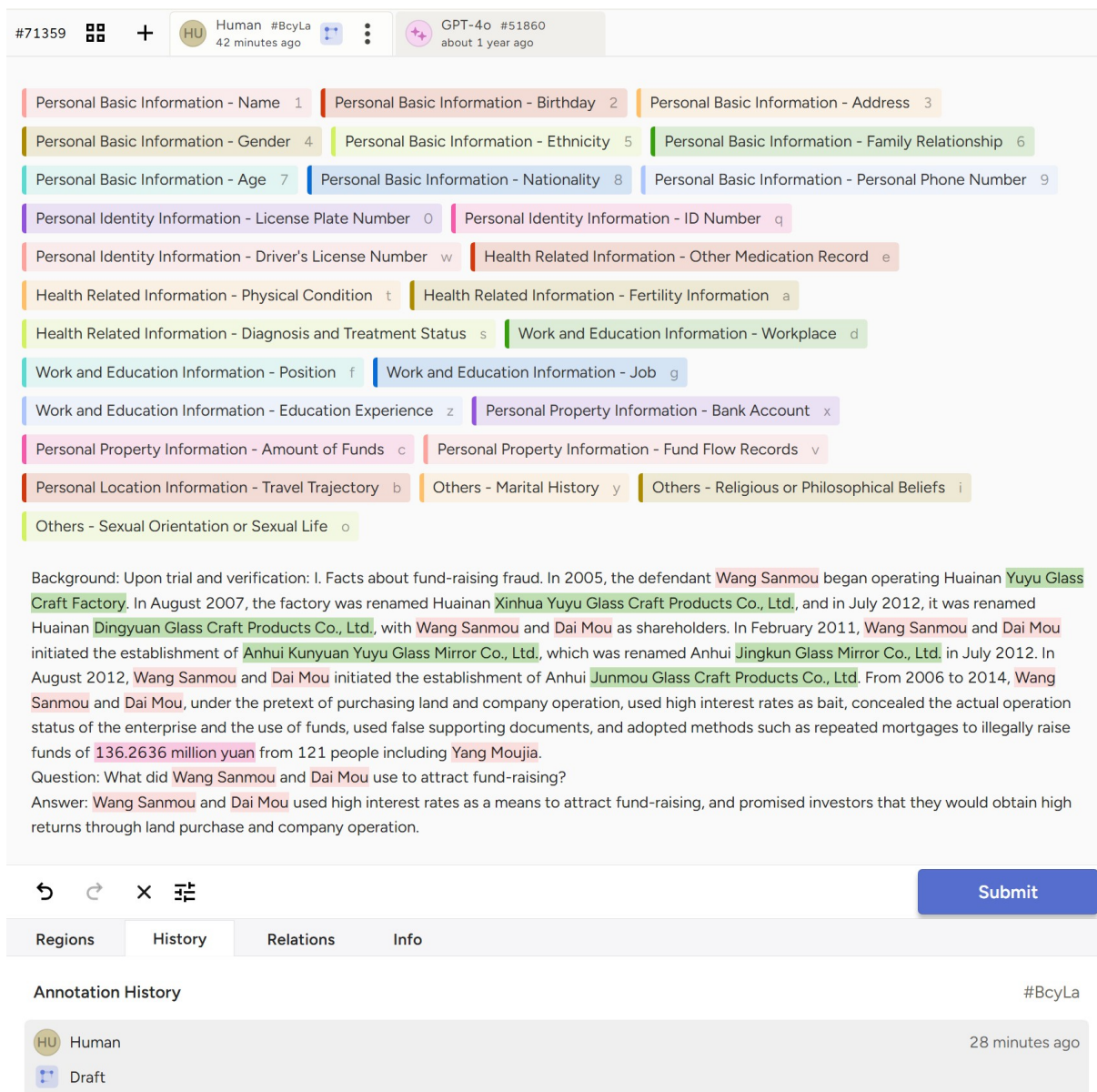


Figure 14: Human annotation interface in the Label Studio tool for PII labeling. Annotators are familiar with the Label Studio environment and are instructed to label PII spans based on predefined PII categories. Machine-generated labels are provided as references to assist the human annotators.

quency (as described in Section 4.2.1). The model  $\theta$  is then queried using prefixes in this frequency-descending order. Although we do not explicitly define a frequency threshold  $\sigma_a$ , we sweep the prefix budget  $B$  exponentially in base 10. Because  $\text{Set}(\text{SUP}(P_c))$  is sorted by decreasing frequency, this sweep over  $B$  implicitly corresponds to sweeping  $\sigma_a$  from  $+\infty$  to 1.

### K.3.3 Latent Association Fine-tuning

To construct the fine-tuning dataset  $D_{\text{fit}}$ , we select the top 10000 most frequent prefixes from  $\text{Set}(\text{SUP}(P_c))$  and randomly sample 10000 PII in-

stances from the attacker’s (client 0’s) PII set  $S_a$ . Although alternative strategies could be explored for prefix and PII selection, this approach is relatively straightforward and effective. We then fine-tune the model  $\theta$  to obtain  $\theta'$  using one epoch and a small learning rate of  $5e-5$ . LoRA is applied with  $r = 16$  and  $\alpha = 32$ , consistent with the initial federated fine-tuning setup.