

PII-Bench: Evaluating Query-Aware Privacy Protection Systems

Hao Shen^{1,2}, Zhouhong Gu³, HaoKai Hong¹, Weili Han^{1,2*}, Hongfeng Chai¹

¹Institute of Fintech, Fudan University

²Laboratory of Data Analytics and Security, Fudan University

³Fudan University

{hshen22, zhgu22, hkhong23}@m.fudan.edu.cn {wlhan, hfchai}@fudan.edu.cn

Abstract

The widespread adoption of Large Language Models (LLMs) has raised significant privacy concerns regarding the exposure of personally identifiable information (PII) in user prompts. To address this challenge, we propose a query-unrelated PII masking strategy and introduce PII-Bench, the first comprehensive evaluation framework for assessing privacy protection systems. PII-Bench comprises 2,842 test samples across 7 PII types with 55 fine-grained subcategories, featuring diverse scenarios from single-subject descriptions to complex multi-party interactions. Each sample is carefully crafted with a user query, context description, and standard answer indicating query-relevant PII. Our empirical evaluation reveals that while current models perform adequately in basic PII detection, they show significant limitations in determining PII query relevance. Even advanced LLMs struggle with this task, particularly in handling complex multi-subject scenarios, indicating substantial room for improvement in achieving intelligent PII masking.

1 Introduction

Recent years have witnessed the widespread adoption of Large Language Models (LLMs), with an increasing number of users directly interacting with these models through APIs for various tasks, ranging from daily conversations to complex analytical work (Sun et al., 2023; Yang et al., 2024b; Wong et al., 2023). Despite the convenience these services offer, users often overlook a significant privacy risk: the prompts submitted to LLMs frequently contain substantial personally identifiable information (PII) (Achiam et al., 2023). Such information is vulnerable not only to interception by malicious actors during transmission (Parast et al., 2022) but also to potential misuse by unethical service providers who might collect and incorporate it

*Corresponding authors: Weili Han
wlhan@fudan.edu.cn

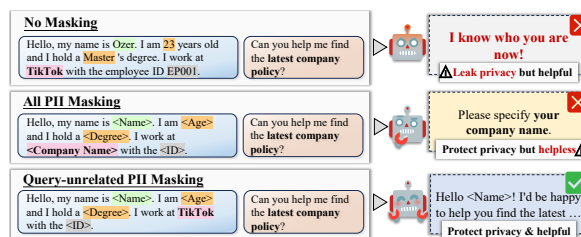


Figure 1: The overall performance of three PII Masking strategies: No Masking, All PII Masking, and Query-unrelated PII Masking. Effective Privacy Protection Systems are required to maintain LLMs' functionality while protect user's privacy as much as possible.

into subsequent model training, leading to permanent privacy breaches (Liu et al., 2023).

Current practices reveal that the vast majority of users adopt a zero-protection approach when utilizing LLM services, submitting original prompts containing PII directly to the LLMs. While an obvious protection strategy would be to mask all PII (Nakamura et al., 2020; Biesner et al., 2022; Lukas et al., 2023), as shown in Figure 1, this approach significantly compromises service quality. An ideal Privacy Protection System should maintain LLMs' functionality while maximizing user privacy protection. For instance, when a user inquires about a candidate's suitability for a senior researcher position, masking their educational background and work experience would render the LLM incapable of making an effective assessment.

This observation motivates our proposal of a query-unrelated PII masking strategy: Masking only the PII irrelevant to user queries while retaining essential information. In the aforementioned example, this approach would preserve the candidate's educational and professional information while masking unrelated personal details such as contact information.

The implementation of query-unrelated PII masking strategy faces two-tier challenges. The first involves accurate identification of all PII within the prompt, serving as foundational work.

The second requires determining the relevance of identified PII to user queries. While existing research has made progress in basic PII detection, systematic studies considering query relevance remain scarce.

To advance the field of privacy-preserving language models, we present PII-Bench, a comprehensive evaluation framework designed to assess Privacy Protection Systems' efficacy in preserving LLMs' core functionalities while optimizing user privacy safeguards. PII-Bench comprising 2,842 carefully designed test samples across 7 PII types with 55 fine-grained subcategories, ranging from basic personal information to complex social relationship data. Each sample consists of three key components: (1) A user query simulating real-world information needs. (2) A context description containing diverse PII. (3) A standard answer indicating query-relevant PII and masking requirements.

Our experimental analysis reveals that while existing models, including Bidirectional Long Short-Term Memory with Conditional Random Fields (BiLSTM-CRF) (Chen et al., 2017), perform adequately in basic PII detection, they demonstrate notable limitations in determining PII query relevance. Even advanced LLMs face challenges in this task, indicating substantial room for improvement in achieving intelligent PII masking. Despite the recent advances in model architecture and training techniques, Small Language Models (SLM) still show considerable performance gaps compared to larger LLMs, particularly in determining PII query relevance.

The primary contributions of this work include: 1. The first proposal of query-unrelated PII masking strategy, offering novel approaches to maintain LLM service quality while protecting privacy. 2. Development of PII-Bench evaluation framework, enabling systematic assessment of models' capabilities in PII identification and query relevance determination. 3. Experimental revelation of current model limitations in this task, providing direction for future research.

2 Related Work

2.1 Privacy-Preserving Text Processing

Text privacy protection has emerged as a critical challenge in natural language processing applications. Papadopoulou et al. (2022) proposed text sanitization that combines entity detection with privacy

risk assessment to guide masking decisions. Shen et al. (2024) extended this approach with an end-to-end framework that preserves task utility during privacy protection. Exploring information preservation, Meisenbacher and Matthes (2024) introduced differential privacy techniques for text modification, demonstrating improved semantic retention over traditional masking methods. While these approaches have advanced privacy protection techniques, they primarily focus on document-level sanitization without considering the dynamic nature of user interactions. Our work introduces query-aware privacy protection that adaptively balances information utility with privacy requirements.

2.2 Query-Aware PII Detection

Traditional PII detection has evolved from rule-based systems (Ruch et al., 2000; Douglass et al., 2005) to neural architectures (Deleger et al., 2013; Deroncourt et al., 2017; Johnson et al., 2020), with recent work demonstrating the effectiveness of transformer-based models in identifying sensitive information (Asimopoulos et al., 2024). Large language models have shown promising results in recognizing diverse PII types (Singhal et al., 2024; Bubeck et al., 2023), yet they treat all sensitive information with uniform importance. Our framework introduces a novel dimension to PII detection by incorporating query relevance assessment. Rather than applying uniform protection measures, we focus on identifying which PII elements are essential for addressing user queries. This approach enables more nuanced privacy protection by distinguishing between query-related and query-unrelated sensitive information, though the actual masking or protection mechanisms are left to downstream applications.

2.3 Privacy Protection Benchmarks

Existing benchmarks for evaluating privacy protection methods have primarily focused on general PII detection capabilities. Pilán et al. (2022) introduced TAB, a benchmark based on legal court cases, which evaluates text anonymization performance. However, it does not assess the model's ability to distinguish query-related information. The recent work by Sun et al. (2024) proposed evaluation metrics for privacy-preserving prompts, but their focus remains limited to general desensitization effectiveness. Li et al. (2024) developed LLM-PBE to assess privacy risks in language models, though their emphasis is on model-side privacy

Symbol	Description
p	A prompt consisting of a user description and a query
d	User description containing personal information
q	User query specifying the information need
d'	Modified description with masked PII
p'	Modified prompt (d', q) after PII masking
\mathcal{S}	Set of subject individuals mentioned in the description
s_i	The i -th subject individual
\mathcal{E}	Complete set of PII entities in the prompt
\mathcal{E}_i	Set of PII entities associated with subject s_i
e_j^i	The j -th PII entity of subject s_i
\mathcal{E}_q	Subset of PII entities necessary for answering query q
\mathcal{T}	Set of predefined PII types

Table 1: Notation used throughout in Task Definition.

rather than input text protection.

PII-Bench addresses these limitations by providing a comprehensive evaluation framework that assesses both PII detection accuracy and the ability to determine query-related information. This dual focus enables more realistic evaluation of privacy protection systems in interactive scenarios, where the relevance of sensitive information varies with user queries.

3 PII-Benchmark

3.1 Task Definition

Privacy Protection Systems target at maintaining LLM functionality while maximizing user privacy protection. Let p be a prompt consisting of a user description d and a query q . The description d contains information about multiple subject individuals $\mathcal{S} = \{s_1, \dots, s_m\}$. For each subject s_i , there exists an associated set of PII entities $\mathcal{E}_i = \{e_1^i, \dots, e_k^i\}$. The complete set of PII entities in prompt p is defined as $\mathcal{E} = \bigcup_{i=1}^m \mathcal{E}_i$, where each entity $e \in \mathcal{E}$ belongs to a predefined PII type from set \mathcal{T} (see Appendix A.2). Let $\mathcal{E}_q \subseteq \mathcal{E}$ denote the subset of PII entities that are necessary for answering query q .

Based on this definition, we propose three fundamental evaluation tasks for Privacy Protection Systems:

(1) PII Detection Task: Given prompt p , the model needs to: identify the minimal text spans for all PII entities $e \in \mathcal{E}$; establish associations between each entity e and its corresponding subject $s \in \mathcal{S}$; assign the correct PII type $t \in \mathcal{T}$ to each entity e .

(2) Query-Related PII Detection Task: Given prompt p , the model needs to determine the minimal subset of PII entities $\mathcal{E}_q \subseteq \mathcal{E}$. This subset should only contain PII entities necessary to answer query q , maximizing protection of non-relevant per-

sonal information.

(3) Query-Unrelated PII Masking Task: This task is what we propose the optimal form of privacy protection system. Given prompt p , the model should generate a modified description d' where query-unrelated PII entities are masked while preserving the necessary ones. Formally, the model should identify \mathcal{E}_q and generate d' where all PII entities in $\mathcal{E} \setminus \mathcal{E}_q$ are masked while preserving those in \mathcal{E}_q . The masking operation should maintain text coherence and readability while ensuring effective privacy protection for non-essential personal information. The resulting prompt $p' = (d', q)$ should enable LLMs to accurately address the query while minimizing exposure of irrelevant personal information.

3.2 PII-Bench Construction

Based on the task definition above, we designed an automated process for constructing the PII evaluation dataset, as illustrated in Figure 2.

3.2.1 PII Entity Generation

Following Papadopoulou et al. (2022), we expanded the PII type set \mathcal{T} with 7 main types into 55 fine-grained subcategories as detailed in Table 11, employing two complementary strategies for entity generation:

(1) Rule-based Generation: Applicable for deterministic PII types with fixed formats or enumerable value sets (e.g., phone numbers, email addresses, ID numbers). (2) LLM-based Generation: Applicable for non-deterministic PII types requiring contextual understanding and real-world knowledge (e.g., occupation descriptions, detailed addresses).

3.2.2 User Description Generation

Single-Subject Description Construction: The construction of single-subject descriptions follows a three-stage process:

(1) Entity Selection: For subject s , randomly sample n entities ($4 \leq n \leq 16$) from different PII types to construct entity set \mathcal{E} . The sampling process ensures diversity of PII types while considering their natural distribution in real-world scenarios. (2) Consistency Optimization: Ensure logical compatibility among entities in \mathcal{E} through LLM-based verification with crafted prompts (detailed in Appendix C). For example, verifies reasonable correspondence between age and educational history as shown in Figure 2. (3) User Desc Generation: Selects appropriate expression styles to generate

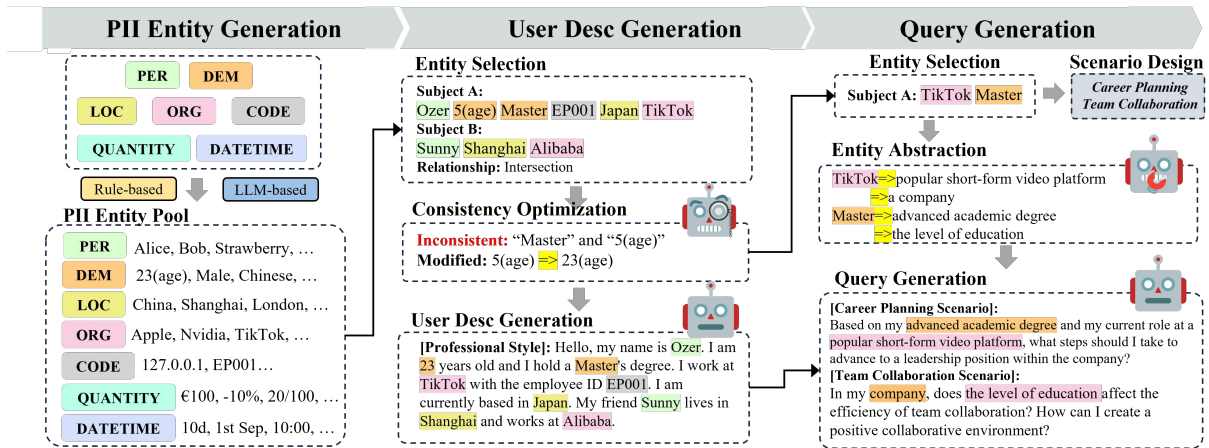


Figure 2: **PII-Bench** synthesis process consists of three main modules: (a) PII Entity Generation, (b) User Description Generation, and (c) Query Generation.

the user description. It employs formal description formats like job resumes and employee records in professional scenarios; casual expressions like personal profiles and self-introductions in social scenarios.

Multi-Subject Description Construction: The construction process for multi-subject related descriptions includes these key steps:

(1) Entity Selection: Construct relationship network $R(s_i, s_j)$ for subject pairs (s_i, s_j) . Relationship types include intersection relationships like colleagues and alumni, hierarchical relationships like parent-child and teacher-student, and non-intersection relationships with no direct connection. (2) Consistency Optimization: Apply relationship-aware verification based on R to maintain cross-subject coherence (detailed in Appendix C). The optimization enforces shared attributes for intersection relationships (e.g., company address for colleagues), structural constraints for hierarchical relationships (e.g., age differences for parent-child pairs), and removes samples with irresolvable contradictions. (3) User Desc Generation: This stage designs natural interaction environments matching relationship characteristics, placing subjects in realistic scenarios (like meetings, family activities) and constructing multi-party dialogue flows to reflect interactive relationships.

3.2.3 Query Construction

For each description d , we construct queries through four phases:

(1) Entity Selection: Randomly sample k entities ($1 \leq k \leq 3$) from \mathcal{E} to form query-relevant entity set \mathcal{E}_q . (2) Scenario Design: Generate domain-specific contexts aligned with real-world applications (detailed in Appendix D). For instance, given

\mathcal{E}_q containing work experience and educational background, we design scenarios such as recruitment evaluation or career planning. (3) Entity Abstraction: Transform specific entities in \mathcal{E}_q into abstract representations while preserving semantic properties. (4) Query Generation: Synthesize natural queries q by integrating abstracted entities into their corresponding scenarios.

3.2.4 Human Verification

All content generated by GPT-4-0806 undergoes rigorous verification by five professional annotators and the authors, focusing on: (1) Completeness and accuracy of PII entity annotations in description d . (2) Correspondence between query q and query-relevant entity set \mathcal{E}_q . (3) Overall semantic coherence and scenario authenticity. Complete annotation guidelines and quality control procedures are detailed in Appendix G.

3.3 Dataset Partitioning and Statistics

Table 3 presents the partition and key statistics of PII-Bench, which comprises two main datasets (PII-single and PII-multi) and two specialized test sets (PII-hard and PII-distract). Each sample follows a consistent JSON structure containing four key components: user description, query, comprehensive PII entity annotations, and query-relevant PII labels, as illustrated in Figure 3.

PII-Single and PII-Multi: Based on the number of subjects in descriptions, the dataset is divided into two main subsets. PII-Single contains 1,214 description-query pairs involving single subjects, focusing on model performance in handling individual information. PII-Multi contains 1,228 description-query pairs involving multiple related subjects, evaluating model capability in handling

Dataset	PII-F1	Query-F1
PII-single	97.2 ± 1.1	95.1 ± 1.3
PII-multi	95.4 ± 1.2	94.3 ± 1.5
PII-hard	91.3 ± 1.1	90.3 ± 1.2
PII-distract	92.8 ± 1.8	91.5 ± 2.1

Table 2: Human performance in PII-Bench. PII-F1 measures accuracy in the PII detection task while Query-F1 evaluates the query-relevant PII detection task.

privacy information within complex interpersonal networks.

Test-Hard Construction: Select 200 challenging instances from PII-Single and PII-Multi to construct Test-Hard dataset, based on criteria including: (1) Maximum character length of description text d . (2) Highest PII entity density ($|\mathcal{E}|/|d|$). (3) Samples with the most query-relevant entities ($|\mathcal{E}_q|$).

Test-Distract Construction: Construct 200 samples simulating complex multi-user interaction scenarios. Each sample integrates five different descriptions from PII-Single and PII-Multi, and constructs queries involving three of these descriptions based on professional networks, knowledge platforms, and community forum interaction templates. The generation process employs dialogue flow transformation strategies to ensure natural transitions and semantic coherence, simulating real-world information interference and complex interaction patterns.

3.4 Human Performance

To establish a human baseline, we recruited 25 graduate students with at least two years of research experience in privacy protection and data security from top Chinese universities. All participants completed comprehensive training and passed a qualification test before formal evaluation (details in Appendix E). We evaluated 400 randomly sampled instances (200 each from PII-single and PII-multi) and 100 instances from PII-distract, with each instance independently assessed by five participants. Participants performed two sequential tasks: PII detection, which involved determining minimal text spans, associated subjects, and PII types for all entities in the user description, followed by query-relevant PII detection to identify entities essential for addressing the given query. The result of the human baseline is shown in Table 2.

Name	#Sample	Avg #Subject	Avg #Char (Desc)	Avg #PII (Desc)	Avg #Char (Query)	Avg #PII (Query)
PII-single	1,214	1.0	893.48	7.67	211.21	1.95
PII-multi	1,228	2.0	652.65	13.14	236.21	2.06
PII-hard	200	1.5	778.03	10.60	222.09	2.10
PII-distract	200	7.5	4,403.64	51.08	859.69	5.82
All	2,842	1.92	1,028.32	13.30	268.41	2.28

Table 3: Statistics of PII-Bench.

4 Experiments

4.1 Overall Setup

Traditional Model Baselines: We implemented **BILSTM-CRF** as a traditional sequence labeling baseline, following the architecture proposed by Huang et al. (2015). We trained the model using Adam optimizer with a learning rate of 1e-3 and batch size of 32 for 50 epochs on the PII-Bench training set.

LLM Baselines: The evaluation encompassed both API-based and open-source large language models. API-based models included GPT-4o-2024-0806 (**GPT4o**) (OpenAI, 2024), Claude-3.5-Sonnet (**Claude3.5**) (Anthropic, 2024), and DeepSeek-Chat **DeepseekV3** (Liu et al., 2024), accessed through their respective official APIs between January 1 and February 10, 2025. Open-source alternatives comprised Llama-3.1-70B-Instruct (**Llama3.1**) (Dubey et al., 2024), and Qwen-2.5-72B-Instruct (**Qwen2.5**) (Yang et al., 2024a).

SLM Baselines: To investigate scaling effects, we included two small-scale language models: Llama-3.1-8B-Instruct (**Llama3.1-SLM**) and Qwen-2.5-7B-Instruct (**Qwen2.5-SLM**). All experiments utilized default parameters with temperature set to 0 to ensure reproducibility. Additionally, we conducted a comprehensive evaluation on smaller, deployment-ready models (0.5B-3B parameters) to assess their viability for on-device PII protection, with detailed results presented in Appendix F.3.

Prompt Baselines: The assessment incorporated multiple prompting strategies for query-related PII detection. **Naive** inputs the user description and query. **Naive /w Choice** includes a list of candidate PII entities to constrain the selection space; we report it as an oracle upper bound that isolates relevance reasoning from detection, not as a deployable strategy. **Self-CoT** (Wei et al., 2022) incorporating step-by-step reasoning prompts. **Auto-CoT** (Zhang et al., 2022), which automates the generation of chain-of-thought demonstrations through three-shot setting. **Self-Consistency (SC)** (Wang et al., 2022), which synthesizes multiple reasoning paths to derive the final output. **Plan-and-Solve**

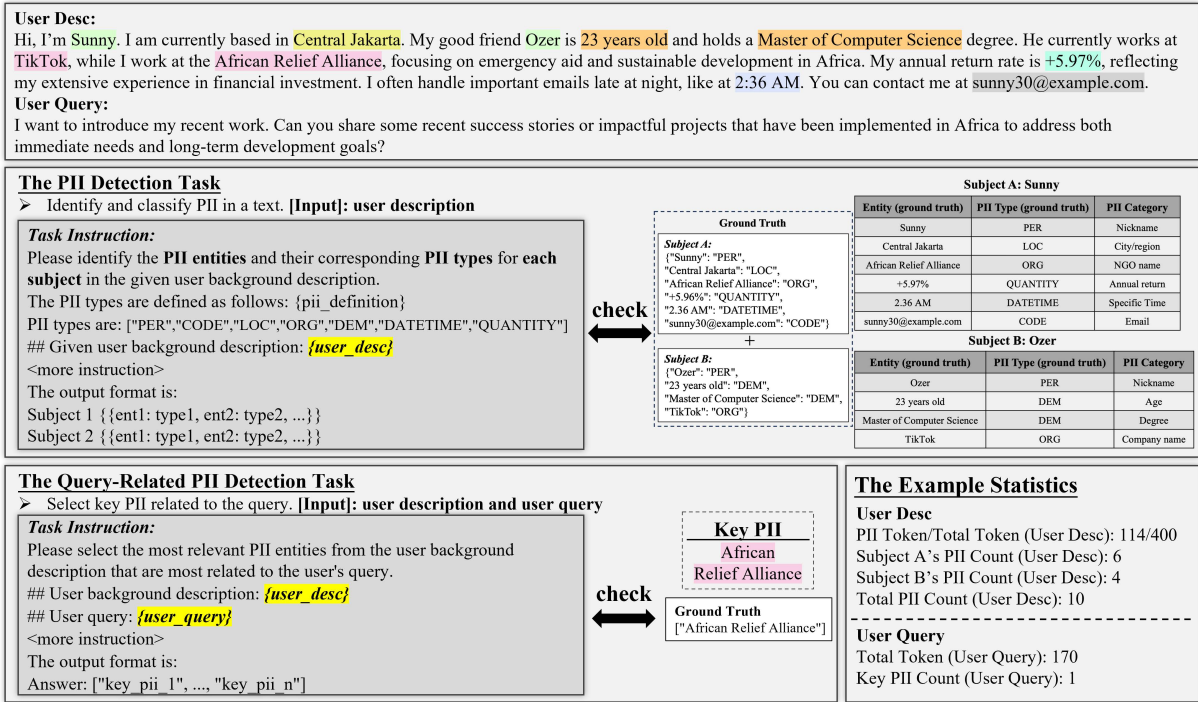


Figure 3: An example from **PII-Bench**, which aims to evaluate Privacy Protection System’s ability by masking maximize PII while maintaining LLM’s functionality. The evaluation is separated by two fundamental tasks: (a) The PII Detection Task: Identify and classify PII entities for each subject in the prompt, with ground truth labels shown on the right side. (b) The Query-Related PII Detection Task: Determine which PII entities are necessary for answering the user query, enabling selective masking of irrelevant personal information.

CoT (PS-CoT) (Wang et al., 2023) develops a strategic plan before executing the solution process. Appendix F.6 provides details of each prompts.

Metrics: The PII detection task evaluates model performance through two sets of metrics: **Strict-F1** measures the accuracy of subject identification, entity span detection, and PII type classification simultaneously. **Ent-F1** focuses on entity span detection independent of subject attribution and type classification. For query-related detection, model performance is measured through **Precision, Recall, and F1**. Considering the inherent variation in entity expressions and potential partial matches, **RougeL-F** is employed for both tasks to complement the exact matching metrics. Detailed computation procedures are provided in Appendix F.1. All reported model performance metrics are mean scores over test sets (temperature=0, top_k=1).

4.2 Performance on Query-Unrelated PII Masking

We evaluate models’ performance on the query-unrelated PII masking task, which requires both accurate PII detection and relevance assessment. Table 4 presents our results:

Joint task yields improved performance. Mod-

els achieve higher F1 scores in this combined task compared to individual query-relevance tasks. GPT4o reaches 0.77 F1 with Self-Consistency prompting, suggesting complementary signals from the joint objective. Open-source models demonstrate comparable capabilities, with both Llama3.1 and Qwen2.5 achieving 0.76 F1 using Auto-CoT.

4.3 Performance on PII Detection

Results on the PII detection task (Table 6) reveal:

Large Language Models demonstrate superior detection capabilities. API-based LLMs achieve strong performance, with DeepSeekV3 and GPT4o leading in Strict-F1 scores (0.903 and 0.891 on PII-Single and PII-Multi, respectively). Open-source Llama3.1 shows competitive performance, particularly in entity recognition (Ent-F1: 0.942 on PII-Multi).

Entity type classification remains challenging. A consistent gap between Strict-F1 and Ent-F1 scores indicates that accurate PII type classification poses greater challenges than entity boundary detection. This disparity becomes more pronounced in the PII-Distract dataset, suggesting increased difficulty in precise PII categorization under complex

Method	GPT4o		Llama3.1		Qwen2.5		Llama3.1-SLM		Qwen2.5-SLM	
	F1	RougeL-F	F1	RougeL-F	F1	RougeL-F	F1	RougeL-F	F1	RougeL-F
<i>Basic Method</i>										
Naive	0.72	0.72	0.72	0.73	0.70	0.70	0.42	0.43	0.54	0.58
<i>Advanced Method</i>										
Self-CoT	0.76	0.77	0.75	0.75	0.73	0.73	0.53	0.54	0.54	0.58
Auto-CoT(3-shot)	0.75	0.75	0.76	0.77	0.76	0.76	0.57	0.58	0.54	0.58
Self-Consistency	0.77	0.77	0.71	0.72	0.71	0.72	0.49	0.50	0.49	0.53
PS-CoT	0.74	0.74	0.72	0.73	0.73	0.73	0.48	0.50	0.56	0.60
<i>w/ Extra Information</i>										
Naive w/ Choice	0.82	0.82	0.77	0.78	0.79	0.79	0.46	0.48	0.67	0.71

Table 4: Performance comparison on the Query-Unrelated PII Masking task (PII-single and PII-multi datasets). The best performance for each model (excluding Naive w/ Choice) is in **bold**.

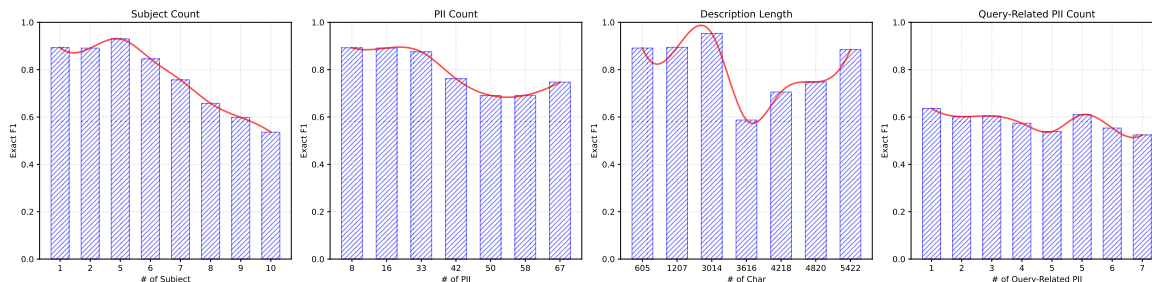


Figure 4: The performance of GPT-4o is correlated with the number of subject, the number of PII, description length, and the number of query-related PII.

Method	Test-Hard		Test-Distract	
	F1	RougeL-F	F1	RougeL-F
<i>Basic Method</i>				
Naive	0.36	0.36	0.57	0.57
<i>Advanced Method</i>				
Self-CoT	0.45	0.45	0.66	0.67
Auto-CoT(3-shot)	0.45	0.40	0.62	0.63
Self-Consistency	0.46	0.46	0.62	0.63
PS-CoT	0.38	0.38	0.67	0.67
<i>w/ Extra Information</i>				
Naive w/ Choice	0.53	0.53	0.66	0.66

Table 5: Performance comparison on challenging test sets using GPT4o.

scenarios.

4.4 Performance on Query-Related PII Detection

Table 7 presents the results on PII-single dataset across different model scales and prompting strategies.

Limited Performance of Current LLMs. Advanced LLMs exhibit limited performance, with GPT4o achieving only 0.627 F1 score with Naive prompting—substantially below human performance (0.951 F1). Open-source alternatives show competitive performance, with Qwen2.5 reaching 0.615 F1.

Impact of Advanced Prompting. Chain-

of-thought approaches generally improve performance, with Self-Consistency and Auto-CoT proving most effective (0.716 F1 for GPT4o with Self-Consistency; 0.710 F1 for Qwen2.5 with Auto-CoT). However, these benefits are highly dependent on model scale—smaller models often show degraded performance with complex prompting strategies.

Effectiveness of Entity Candidates. Providing candidate PII entities (Naive w/ Choice) substantially improves performance across all models (e.g., GPT4o improves from 0.627 to 0.842 F1). However, practical applicability is limited as candidate entities are rarely available in real-world scenarios.

4.5 Privacy-Utility Tradeoff Analysis

We evaluated the tradeoff between privacy protection and query utility across different PII masking strategies

Experimental Setup. We selected 200 unique user descriptions from PII-Bench and applied three masking methods: (1) **No Mask**: Original text with all PII preserved; (2) **All PII Mask**: All detected PII entities replaced with their corresponding tags; (3) **Query-unrelated PII Mask**: Only query-irrelevant PII entities masked. For each method, we generated responses using GPT4o.

Metrics. The privacy-utility tradeoff is eval-

Baseline Models	PII-Single			PII-Multi			PII-Hard			PII-Distract		
	Strict-F1	Ent-F1	RougeL-F	Strict-F1	Ent-F1	RougeL-F	Strict-F1	Ent-F1	RougeL-F	Strict-F1	Ent-F1	RougeL-F
<i>Traditional Model</i>												
BiLSTM-CRF	-	0.851	-	-	0.828	-	-	0.684	-	-	0.787	-
<i>API-based Large Language Model</i>												
GPT4o	0.893	0.914	0.895	0.891	0.923	0.893	0.817	0.869	0.819	0.715	0.868	0.716
Claude3.5	0.858	0.891	0.862	0.890	0.920	0.892	0.813	0.857	0.818	0.910	0.948	0.911
DeepSeekV3	0.903	0.921	0.905	0.884	0.927	0.886	0.838	0.893	0.838	0.658	0.945	0.658
<i>Open-source Large Language Model</i>												
Llama3.1	0.881	0.913	0.883	0.883	0.942	0.884	0.840	0.893	0.841	0.834	0.946	0.835
Qwen2.5	0.866	0.908	0.869	0.853	0.918	0.855	0.804	0.876	0.806	0.647	0.941	0.649
<i>Open-source Small Language Model</i>												
Llama3.1-SLM	0.748	0.800	0.752	0.778	0.869	0.781	0.718	0.798	0.722	0.551	0.876	0.552
Qwen2.5-SLM	0.787	0.846	0.792	0.451	0.806	0.453	0.591	0.810	0.594	0.454	0.815	0.456

Table 6: Performance of baseline models under the PII Detection task. Results in **bold** indicate the best performance for each dataset and metric category.

Method	GPT4o		Llama3.1		Qwen2.5		Llama3.1-SLM		Qwen2.5-SLM		
	F1	RougeL-F	F1	RougeL-F	F1	RougeL-F	F1	RougeL-F	F1	RougeL-F	
<i>Basic Method</i>											
Naive	0.63	0.63	0.63	0.63	0.62	0.62	0.33	0.33	0.41	0.41	
<i>Advanced Method</i>											
Self-CoT	0.71	0.72	0.69	0.69	0.67	0.68	0.39	0.39	0.40	0.41	
Auto-CoT(3-shot)	0.66	0.66	0.70	0.72	0.71	0.72	0.43	0.44	0.37	0.38	
Self-Consistency	0.72	0.72	0.63	0.64	0.65	0.65	0.31	0.32	0.32	0.33	
PS-CoT	0.65	0.65	0.65	0.66	0.67	0.67	0.35	0.36	0.45	0.46	
<i>w/ Extra Information</i>											
Naive w/ Choice	0.84	0.84	0.76	0.76	0.83	0.83	0.52	0.52	0.77	0.77	

Table 7: Performance comparison on the Query-Related PII Detection task (PII-single dataset).

Method	P	U	B
No Mask	0.00	1.00	0.50
All PII Mask	1.00	0.52	0.76
Query-unrelated PII Mask	0.83	0.89	0.86

Table 8: Privacy-utility tradeoff across different PII masking strategies. Balanced scores (B) combine both metrics with equal weights.

uated through two metrics: **Privacy Score (P)**: Quantifies the proportion of PII successfully protected in the processed text. **Utility Score (U)**: Combines semantic similarity and LLM evaluation for quality assessment between mask and unmasked responses. Detailed computation procedures are provided in Appendix F.2.

Results. Table 8 shows that Query-unrelated PII Mask achieves optimal balance between privacy protection ($P = 0.83$) and utility preservation ($U = 0.89$), outperforming other strategies in balanced score ($B = 0.86$).

4.6 Real-World Dataset Validation

To validate the reliability of synthetic data, we constructed PII-Real, a dataset comprising 100 instances derived from publicly available profiles of 20 AI researchers with manually annotated PII entities and human-written queries. Evaluation on

PII-Real demonstrates consistent alignment with synthetic data. Comprehensive experimental results and validation analysis are provided in Appendix F.4.

4.7 In-depth Performance Analysis

Factors Influencing Performance. Figure 4 reveals several critical factors affecting model accuracy: performance degrades sharply beyond 5 subjects (F1 drops from 0.85 to 0.52), 33 PII entities, or 3000 characters in text length. Query-Related entity count shows modest impact, with gradual decline from 1 to 7 entities.

Performance on Challenging Scenarios. Results on specialized test sets (Table 5) reveal significant performance degradation in complex scenarios. On Test-Hard, featuring high PII density and long texts, even the best-performing Self-Consistency approach achieves only 0.463 F1. Test-Distract’s multi-subject scenarios pose similar challenges.

Evaluation on Reasoning Models. We additionally evaluate GPT-5, DeepSeek-V3.2, and DeepSeek-R1 on all three tasks (Table 9), with an extended comparison over four Qwen3 variants reported in Appendix H. On PII Detection, scores saturate on PII-Single, while the largest gains from reasoning appear on the high-entropy PII-Distract subset; open-source models match proprietary sys-

Model	Detect.	Q-Rel.	Q-Unrel.
GPT-5	0.871	0.733	0.840
DeepSeek-V3.2	0.856	0.765	0.852
DeepSeek-R1	0.876	0.745	0.850

Table 9: Recent LLMs on PII-Single. *Detect.*: Strict-F1 on PII Detection using the Naive (zero-shot) prompt. *Q-Rel.* and *Q-Unrel.*: best F1 across prompting strategies, excluding the Naive w/ Choice oracle. Full results in Appendix H.

tems, with DeepSeek-R1 on par with GPT-5 across every PII Detection subset. The two query-aware tasks remain considerably harder: the strongest system trails the human baseline on Query-Related PII Detection by roughly 18 F1 points, indicating that relevance judgement rather than entity detection is the dominant bottleneck.

5 Conclusion

This paper introduces PII-Bench, a comprehensive evaluation framework comprising 2,842 test samples across 7 PII types with 55 fine-grained subcategories, and proposes a query-unrelated PII masking strategy to balance privacy protection with LLM utility. Our empirical evaluation reveals that while advanced LLMs demonstrate strong performance in basic PII detection, they exhibit substantial limitations in query-relevance assessment. Small-scale models show considerably larger performance gaps across all evaluation tasks. These findings establish foundational benchmarks and expose critical challenges in privacy-aware PII handling.

Limitations

Despite PII-Bench’s contributions to privacy protection evaluation, several limitations merit acknowledgment. While the current dataset encompasses common privacy scenarios, it requires expansion into specialized domains such as medical records and financial transactions. Our automated synthesis methodology mitigates this limitation by enabling flexible dataset expansion across domains, languages, and cultural contexts, supporting continuous refinement of PII categories to meet evolving application requirements. The evaluation framework primarily assesses the accuracy of PII entity detection and query relevance determination, but lacks systematic evaluation of models’ reasoning processes. Specifically, it does not fully capture how models interpret queries, derive information requirements, and make relevance judgments

about sensitive information. This gap in assessment methodology limits our understanding of models’ reasoning capabilities in real-world privacy protection scenarios.

PII-Bench does not evaluate robustness to *inference attacks*, where an adversary recovers a masked entity by combining the still-visible context (for example, inferring an employer from an unmasked job title together with a nearby location cue). Such attacks target the residual signal left after masking and are therefore complementary to the query-aware masking problem we study. Extending PII-Bench with adversarial inference probes is a natural direction for follow-up work.

Ethical Concerns

Throughout the development and implementation of PII-Bench, ethical considerations have remained our paramount priority. To ensure the evaluation dataset itself does not compromise privacy, we have implemented rigorous data synthesis and review protocols, with all sample data undergoing multiple rounds of scrutiny by professional security teams to guarantee the absence of real personal information. During the data generation process, we have carefully engineered our algorithms to ensure equitable representation across different demographic groups, establishing comprehensive human review mechanisms to verify that generated data remains free from bias and discriminatory content.

Acknowledgments

This work was supported by the National Natural Science Foundation of China (No. 72595845).

References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Anthropic. 2024. Claude 3.5 sonnet. <https://www.anthropic.com/news/claude-3-5-sonnet>.
- Dimitris Asimopoulos, Ilias Siniosoglou, Vasileios Argyriou, Thomai Karamitsou, Eleftherios Fountoukidis, Sotirios K Goudos, Ioannis D Moscholios, Konstantinos E Psannis, and Panagiotis Sarigiannidis. 2024. Benchmarking advanced text anonymisation methods: A comparative study on novel and traditional approaches. *arXiv preprint arXiv:2404.14465*.

- David Biesner, Rajkumar Ramamurthy, Robin Stenzel, Max Lübbering, Lars Hillebrand, Anna Ladi, Maren Pielka, Rüdiger Loitz, Christian Bauchhage, and Rafet Sifa. 2022. Anonymization of german financial documents using neural network-based language models with contextual word representations. *International Journal of Data Science and Analytics*, pages 1–11.
- Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. 2023. Sparks of artificial general intelligence: Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712*.
- Tao Chen, Ruifeng Xu, Yulan He, and Xuan Wang. 2017. Improving sentiment analysis via sentence type classification using bilstm-crf and cnn. *Expert Systems with Applications*, 72:221–230.
- Louise Deleger, Katalin Molnar, Guergana Savova, Fei Xia, Todd Lingren, Qi Li, Keith Marsolo, Anil Jegga, Megan Kaiser, Laura Stoutenborough, et al. 2013. Large-scale evaluation of automated clinical note de-identification and its impact on information extraction. *Journal of the American Medical Informatics Association*, 20(1):84–94.
- Franck Deroncourt, Ji Young Lee, Ozlem Uzun, and Peter Szolovits. 2017. De-identification of patient notes with recurrent neural networks. *Journal of the American Medical Informatics Association*, 24(3):596–606.
- Josep Domingo-Ferrer, David Sánchez, and Jordi Soria-Comas. 2022. *Database anonymization: privacy models, data utility, and microaggregation-based inter-model connections*. Springer Nature.
- MM Douglass, GD Clifford, Andrew Reisner, WJ Long, GB Moody, and RG Mark. 2005. De-identification algorithm for free-text nursing notes. In *Computers in Cardiology, 2005*, pages 331–334. IEEE.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Mark Elliot, Elaine Mackey, Kieron O’Hara, and Caroline Tudor. 2016. The anonymisation decision-making framework. *UKAN*.
- Zhiheng Huang, Wei Xu, and Kai Yu. 2015. Bidirectional lstm-crf models for sequence tagging. *arXiv preprint arXiv:1508.01991*.
- Alistair EW Johnson, Lucas Bulgarelli, and Tom J Polard. 2020. Deidentification of free-text medical records using pre-trained bidirectional transformers. In *Proceedings of the ACM Conference on Health, Inference, and Learning*, pages 214–221.
- Qinbin Li, Junyuan Hong, Chulin Xie, Jeffrey Tan, Rachel Xin, Junyi Hou, Xavier Yin, Zhun Wang, Dan Hendrycks, Zhangyang Wang, et al. 2024. Llm-pbe: Assessing data privacy in large language models. *arXiv preprint arXiv:2408.12787*.
- Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. 2024. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*.
- Yang Liu, Yuanshun Yao, Jean-Francois Ton, Xiaoying Zhang, Ruocheng Guo Hao Cheng, Yegor Klochkov, Muhammad Faaiz Taufiq, and Hang Li. 2023. Trustworthy llms: A survey and guideline for evaluating large language models’ alignment. *arXiv preprint arXiv:2308.05374*.
- 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. IEEE.
- Stephen Meisenbacher and Florian Matthes. 2024. Just rewrite it again: A post-processing method for enhanced semantic similarity and privacy preservation of differentially private rewritten text. *arXiv preprint arXiv:2405.19831*.
- Multi-Linguality Multi-Functionality Multi-Granularity. 2024. M3-embedding: Multi-linguality, multi-functionality, multi-granularity text embeddings through self-knowledge distillation.
- Yuta Nakamura, Shouhei Hanaoka, Yukihiro Nomura, Naoto Hayashi, Osamu Abe, Shuntaro Yada, Shoko Wakamiya, and Eiji Aramaki. 2020. Kart: Privacy leakage framework of language models pre-trained with clinical records. *arXiv preprint arXiv:2101.00036*.
- OpenAI. 2024. Hello gpt-4o. <https://openai.com/index/hello-gpt-4o>.
- Anthi Papadopoulou, Yunhao Yu, Pierre Lison, and Lilja Øvrelid. 2022. Neural text sanitization with explicit measures of privacy risk. In *Proceedings of the 2nd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 12th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 217–229.
- Fatemeh Khoda Parast, Chandni Sindhav, Seema Nikam, Hadiseh Izadi Yekta, Kenneth B Kent, and Saqib Hakak. 2022. Cloud computing security: A survey of service-based models. *Computers & Security*, 114:102580.
- Ildikó Pilán, Pierre Lison, Lilja Øvrelid, Anthi Papadopoulou, David Sánchez, and Montserrat Batet. 2022. The text anonymization benchmark (tab): A dedicated corpus and evaluation framework for text anonymization. *Computational Linguistics*, 48(4):1053–1101.

Patrick Ruch, Robert H Baud, Anne-Marie Rassinoux, Pierrette Bouillon, and Gilbert Robert. 2000. Medical document anonymization with a semantic lexicon. In *Proceedings of the AMIA Symposium*, page 729. American Medical Informatics Association.

Zhili Shen, Zihang Xi, Ying He, Wei Tong, Jingyu Hua, and Sheng Zhong. 2024. The fire thief is also the keeper: Balancing usability and privacy in prompts. *arXiv preprint arXiv:2406.14318*.

Shreya Singhal, Andres Felipe Zambrano, Maciej Pankiewicz, Xiner Liu, Chelsea Porter, and Ryan S Baker. 2024. De-identifying student personally identifying information with gpt-4. In *Proceedings of the 17th International Conference on Educational Data Mining*, pages 559–565.

Xiaofei Sun, Xiaoya Li, Jiwei Li, Fei Wu, Shangwei Guo, Tianwei Zhang, and Guoyin Wang. 2023. *Text classification via large language models*. *Preprint*, arXiv:2305.08377.

Xiongtao Sun, Gan Liu, Zhipeng He, Hui Li, and Xiaoguang Li. 2024. Deprompt: Desensitization and evaluation of personal identifiable information in large language model prompts. *arXiv preprint arXiv:2408.08930*.

Lei Wang, Wanyu Xu, Yihuai Lan, Zhiqiang Hu, Yunshi Lan, Roy Ka-Wei Lee, and Ee-Peng Lim. 2023. Plan-and-solve prompting: Improving zero-shot chain-of-thought reasoning by large language models. *arXiv preprint arXiv:2305.04091*.

Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2022. Self-consistency improves chain of thought reasoning in language models. *arXiv preprint arXiv:2203.11171*.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837.

IpKin Anthony Wong, Qi Lilith Lian, and Danni Sun. 2023. *Autonomous travel decision-making: An early glimpse into chatgpt and generative ai*. *Journal of Hospitality and Tourism Management*, 56:253–263.

An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. 2024a. Qwen2.5 technical report. *arXiv preprint arXiv:2412.15115*.

Songhua Yang, Hanjie Zhao, Senbin Zhu, Guangyu Zhou, Hongfei Xu, Yuxiang Jia, and Hongying Zan. 2024b. Zhongjing: Enhancing the chinese medical capabilities of large language model through expert feedback and real-world multi-turn dialogue. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 19368–19376.

Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex Smola. 2022. Automatic chain of thought prompting in large language models. *arXiv preprint arXiv:2210.03493*.

A Details about PII

A.1 PII Definition

In this section, we follow previous work by categorizing Personally Identifiable Information (PII) into the following two categories (Elliot et al., 2016, Domingo-Ferrer et al., 2022, Papadopoulou et al., 2022):

- *Direct identifiers*: Information that can uniquely identify an individual within a dataset (e.g. name, social security number, email address, etc).
- *Quasi identifiers*: Information that cannot uniquely identify an individual on their own but can do so when combined with other quasi-identifiers (e.g. age, gender, occupation, etc).

Because of their high sensitivity or the potential to indirectly identify an individual, both direct and quasi-identifiers are governed by strict legal and privacy standards to ensure personal privacy.

A.2 PII Types

Unlike Papadopoulou et al.’s (2022) classification, we exclude the MISC category due to its ambiguous definition and unclear boundaries. Our taxonomy comprises seven categories:

PER: Refers to individuals’ names, including full names, aliases, and social media usernames.

CODE: Encompasses identifying numbers and codes like social security numbers, phone numbers, passport numbers, email addresses, etc.

LOC: Covers geographical locations such as home or work addresses, cities, countries, etc.

ORG: Pertains to the names of entities like companies, schools, public institutions, etc.

DEM: Represents demographic information including age, gender, nationality, occupation, education level, etc.

DATETIME: Indicates specific dates, times, or durations, such as birthdates, appointment times, etc.

QUANTITY: Refers to significant numerical data like monthly income, expenditures, loan amount, credit score, etc.

A.3 Statistics of PII Types

Figure 5 and Table 10 present the distribution of PII types across our datasets: PII-single (1,214 samples), PII-multi (1,228 samples), PII-hard (200 samples), and PII-distract (200 samples).

- **Type Frequencies:** Organization (ORG) and Code-based identifiers (CODE) constitute significant portions across all datasets, with 17.09% and 15.74% in PII-single, and 13.47% and 15.31% in PII-multi, respectively. This distribution reflects the prevalence of institutional affiliations and digital identifiers in real-world scenarios.
- **Dataset Composition:** PII-multi contains 16,136 PII entities across all categories, maintaining balanced proportions ranging from 13.47% to 15.77% for most types. PII-single follows a similar pattern with 9,303 entities, demonstrating consistent coverage across different PII categories.
- **Specialized Test Sets:** PII-distract, despite comprising only 200 samples, contains 10,211 PII entities due to its multi-description design. PII-hard maintains balanced type coverage with 1,834 entities, with proportions varying from 12.10% to 16.58%.

B PII Entity Generation Methods

The generation of PII entities requires careful consideration of both structural constraints and semantic plausibility. We employ two complementary approaches for entity generation: rule-based generation for structured PII types and language model-based generation for context-dependent information.

B.1 Rule-based Generation

For PII types with well-defined formats or enumerable value sets, we implement deterministic generation methods. These methods encompass both custom rule-based algorithms and the Faker library’s standardized functions. The rule-based approach is particularly effective for:

1. Identification Numbers: Generating valid formats for social security numbers, passport numbers, and employee IDs while maintaining regional compliance.

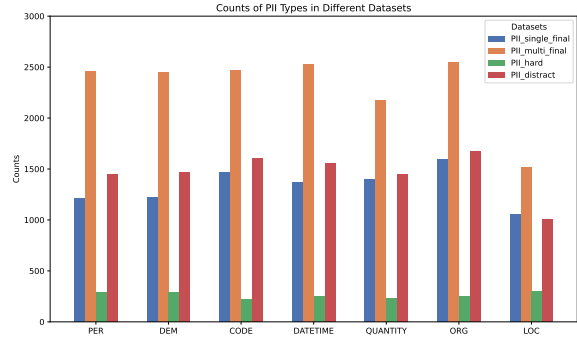


Figure 5: Distribution of PII types across different datasets in PII-Bench.

2. Contact Information: Creating syntactically correct email addresses, phone numbers, and IP addresses.
3. Financial Data: Producing properly formatted credit card numbers, bank account numbers, and other numerical identifiers with appropriate check digits.
4. Temporal Information: Generating dates, times, and durations within reasonable ranges and formats.

Type	PII_single		PII_multi		PII_hard		PII_distract	
	#	%	#	%	#	%	#	%
PER	1,214	13.05	2,456	15.22	286	15.59	1,449	14.19
DEM	1,220	13.12	2,450	15.18	286	15.59	1,467	14.37
CODE	1,464	15.74	2,470	15.31	222	12.10	1,605	15.72
ORG	1,590	17.09	2,544	13.47	251	13.69	1,673	16.38
LOC	1,053	11.32	1,516	15.77	304	16.58	1,008	9.87
DATETIME	1,368	14.70	2,526	15.65	251	13.69	1,559	15.27
QUANTITY	1,394	14.98	2,174	9.40	234	12.76	1,450	14.20
Total	9,303	100	16,136	100	1,834	100	10,211	100

Table 10: Detailed statistics of PII types across datasets. For each dataset, we report both the absolute count (#) and relative percentage (%) of each PII type.

B.2 Language Model-based Generation

For PII types requiring contextual understanding and real-world knowledge, we leverage large language models through carefully designed prompts. This approach is essential for generating:

1. Location Information: Coherent and geographically accurate addresses, landmarks, and regional descriptions.
2. Organizational Entities: Plausible names for educational institutions, companies, and other organizations that reflect real-world naming conventions.

3. Demographic Attributes: Culturally appropriate and consistent demographic information, including ethnicity, nationality, and educational background.

B.3 Entity Categories and Generation Methods

Table 11 presents a comprehensive mapping of PII types to their respective generation methods. The table systematically categorizes 55 distinct PII entities across seven main categories: Personal Identifiers (PER), Codes and Numbers (CODE), Location Information (LOC), Organizational Affiliations (ORG), Demographic Information (DEM), Temporal Data (DATETIME), and Quantitative Values (QUANTITY).

C Consistency Optimization Details

The consistency optimization process is critical for ensuring that randomly sampled PII entities form logically coherent and realistic profiles. We employ GPT-4-0806 with carefully designed prompts to identify and resolve conflicts while preserving the diversity and coverage of PII types.

C.1 Single-Subject Consistency Optimization

For single-subject descriptions, the optimization process addresses intra-subject conflicts arising from incompatible entity combinations. Common conflict patterns include:

- **Temporal inconsistencies:** Age incompatible with work experience duration, education level, or career stage (e.g., a 23-year-old with 15 years of work experience).
- **Professional mismatches:** Occupation inconsistent with education background or salary range (e.g., a high school graduate working as a licensed physician).
- **Geographic contradictions:** Work location distant from residential address without supporting evidence (e.g., daily commute spanning different continents).
- **Financial implausibility:** Income, expenses, and savings that violate basic economic constraints (e.g., monthly expenses exceeding monthly income by orders of magnitude).

The optimization prompt (Figure 10) instructs the model to: (1) identify logical conflicts among

sampled entities; (2) modify only the entity values while preserving their PII type and category labels; (3) maximize entity retention to maintain dataset richness. For instance, given a conflict between “Age: 22 years” and “Work Experience: 15 years as Senior Engineer”, the optimization adjusts the work experience to “2 years as Junior Engineer” rather than removing the entity entirely, thereby preserving both the PII type (DEM) and the entity category (Work Experience).

C.2 Multi-Subject Consistency Optimization

Multi-subject optimization extends the single-subject process by introducing relationship-aware constraints. Consistency rules are enforced based on three relationship categories:

- **Intersection relationships** (e.g., colleagues, classmates): Subjects must share critical attributes such as organization name, work location for colleagues, or educational institution for classmates. The optimization verifies attribute alignment and adjusts inconsistent entities accordingly.
- **Hierarchical relationships** (e.g., parent-child, supervisor-subordinate): Structural constraints are enforced such as age differences (parent at least 18-20 years older than child), authority levels (supervisor holding higher position than subordinate), and derived attributes (children inheriting nationality or ethnicity from parents when culturally appropriate).
- **Non-intersection relationships:** Subjects maintain independent entity sets with no enforced dependencies, though internal consistency within each subject’s profile is still verified.

The multi-subject prompt (Figure 11) provides relationship context and requires the model to output separate, coordinated entity sets for each subject. This relationship-aware approach ensures that multi-subject descriptions reflect realistic interpersonal dynamics and maintain narrative coherence across profiles.

D Query Scenario Design Details

The scenario design phase bridges the gap between abstract entity selection and natural user queries by grounding queries in realistic application contexts.

PII Type	Entity Category	Generation Approach	Format Constraints
PER	Full Name	Rule-based	[First Name] [Last Name]
	Social Media Handle	Rule-based	[@][a-zA-Z0-9]5,15
	Nickname	LLM-based	-
CODE	Social Security Number	Rule-based	XXX-XX-XXXX
	Driver's License	Rule-based	[A-Z][0-9]8
	Bank Account	Rule-based	[0-9]10,12
	Credit Card	Rule-based	[0-9]16
	Phone Number	Rule-based	+ [0-9]1,3- [0-9]10
	IP Address	Rule-based	IPv4/IPv6 format
	Email Address	Rule-based	[user]@[domain].[tld]
	Password Hash	Rule-based	SHA-256
	Passport Number	Rule-based	[A-Z][0-9]8
	Tax ID	Rule-based	[0-9]9
LOC	Employee ID	Rule-based	[A-Z]2[0-9]6
	Student ID	Rule-based	[0-9]8
	Street Address	LLM-based	-
LOC	City/Region	LLM-based	-
	Landmark	LLM-based	-
	Company Name	LLM-based	-
ORG	Educational Institution	LLM-based	-
	Government Agency	LLM-based	-
	NGO	LLM-based	-
	Healthcare Facility	LLM-based	-
	Occupation	Rule-based	Predefined list
DEM	Age	Rule-based	[0-9]1,3
	Gender	Rule-based	Binary/Non-binary
	Height	Rule-based	[0-9]3cm/[0-9]'[0-9]"
	Weight	Rule-based	[0-9]2,3kg/lbs
	Blood Type	Rule-based	A/B/O[+-]
	Sexual Orientation	Rule-based	Predefined list
	Nationality	LLM-based	-
	Ethnicity	LLM-based	-
	Race	LLM-based	-
	Religious Belief	LLM-based	-
	Political Affiliation	LLM-based	-
	Education Level	LLM-based	-
	Academic Degree	LLM-based	-
	Physical Features	LLM-based	-
	Medical Condition	LLM-based	-
Disability Status	LLM-based	-	
DATETIME	Date	Rule-based	YYYY-MM-DD
	Time	Rule-based	HH:MM:SS
	Duration	Rule-based	[0-9]+[dhms]
QUANTITY	Monthly Income	Rule-based	[Currency][0-9]+
	Monthly Expenses	Rule-based	[Currency][0-9]+
	Account Balance	Rule-based	[Currency][0-9]+
	Loan Amount	Rule-based	[Currency][0-9]+
	Annual Bonus	Rule-based	[Currency][0-9]+
	Credit Limit	Rule-based	[Currency][0-9]+
	Social Security Payment	Rule-based	[Currency][0-9]+
	Tax Payment	Rule-based	[Currency][0-9]+
	Debt Ratio	Rule-based	[0-9]1,2.[0-9]2%
	Investment Return	Rule-based	[0-9]1,2.[0-9]2%
ROI	Rule-based	[0-9]1,2.[0-9]2%	
Credit Score	Rule-based	[300-850]	

Table 11: Comprehensive categorization of PII entities and their generation methods. Rule-based generation follows specific format constraints, while LLM-based generation produces contextually appropriate content without rigid formatting requirements.

This process is essential for evaluating whether privacy protection systems can identify query-relevant PII in authentic usage scenarios.

D.1 Domain-Specific Scenario Generation

Given a set of selected entities \mathcal{E}_q , we employ a structured prompting strategy to generate diverse scenario contexts. The prompt instructs the LLM to:

1. **Analyze entity semantics:** Determine the thematic domain implied by the selected entities (e.g., education and work experience suggest career-related scenarios; medical conditions and medications suggest healthcare scenarios).
2. **Generate scenario candidates:** Propose multiple scenario types that naturally require all selected entities, ensuring thematic diversity across the dataset. Example scenarios include career planning, medical consultation, financial advisory, legal counseling, academic mentoring, and housing applications.
3. **Ensure entity necessity:** Verify that each selected entity is meaningfully incorporated into the scenario rather than superficially included, so that query-relevant PII labels reflect genuine information dependencies.

The scenario generation prompt is provided in Figure 12. By requiring coverage of all selected entities while maintaining scenario diversity, this approach produces queries that realistically reflect specialized user information needs across domains.

D.2 Scenario Validation and Filtering

Generated scenarios undergo validation to ensure quality and diversity:

- **Completeness check:** Verify that all entities in \mathcal{E}_q are semantically incorporated into the scenario context.
- **Realism assessment:** Confirm that the scenario represents plausible real-world interactions with AI systems.
- **Diversity enforcement:** Track scenario distribution and reject over-represented scenario types to maintain balanced domain coverage.

This rigorous scenario design process ensures that PII-Bench queries reflect authentic information needs across diverse domains, providing a realistic testbed for evaluating query-aware privacy protection systems.

E Human Evaluation Details

The human evaluation of PII-Bench was conducted with 25 graduate students specializing in data security and privacy protection. All evaluators were pursuing their Master’s or Ph.D. degrees with at least two years of research experience in privacy-preserving machine learning or data protection systems. The evaluation process consisted of three phases: preparation, evaluation, and validation.

During the preparation phase, participants attended a 4-hour training session covering PII taxonomy, recognition guidelines, and query-related detection criteria. The session included hands-on practice with representative cases from each dataset component. Participants then completed a qualification test featuring 20 diverse instances, requiring 90% agreement with expert assessments to proceed to the formal evaluation.

During the evaluation phase, participants used our specialized platform designed for systematic PII assessment. To maintain consistent performance, we limited evaluation sessions to two hours and distributed instances across a two-week period. The platform automatically tracked assessment time and accuracy metrics while enforcing our evaluation protocol: participants first performed PII detection by marking entity spans, linking them to subjects, and assigning PII types, before proceeding to query-related detection.

Our validation process incorporated both automated and manual checks to ensure assessment quality. The platform automatically verified assessment completeness and format consistency. Cases with substantial disagreement (Fleiss’ kappa < 0.6) underwent expert review by two authors with extensive experience in privacy-preserving systems. Evaluators received detailed feedback on their performance and participated in discussion sessions to resolve systematic discrepancies.

Compensation was structured to encourage both accuracy and efficiency, with a base rate of \$30 per hour and performance bonuses based on agreement with other evaluators.

F Experiment Details

F.1 Evaluation Metrics

The evaluation framework employs distinct metrics for PII detection and query-related detection tasks.

For PII detection, let $\mathcal{P} = \{p_1, \dots, p_m\}$ denote the predicted subject set and $\mathcal{G} = \{g_1, \dots, g_n\}$ denote the ground truth subject set. Each subject p_i or g_j contains a set of entity-type pairs $\{(e, t)\}$, where e represents the entity span and t represents its PII type.

For each subject pair (p_i, g_j) , we compute three types of evaluation metrics:

1. Strict Matching: Both entity spans and their types must match exactly:

$$P_{strict}(p_i, g_j) = \frac{|E_{p_i} \cap E_{g_j}|}{|E_{p_i}|} \quad (1)$$

$$R_{strict}(p_i, g_j) = \frac{|E_{p_i} \cap E_{g_j}|}{|E_{g_j}|} \quad (2)$$

$$F1_{strict}(p_i, g_j) = \frac{2 \cdot P_{strict}(p_i, g_j) \cdot R_{strict}(p_i, g_j)}{P_{strict}(p_i, g_j) + R_{strict}(p_i, g_j)} \quad (3)$$

where E_{p_i} and E_{g_j} are the sets of entity-type pairs.

2. Entity-only Matching: Only entity spans need to match:

$$P_{ent}(p_i, g_j) = \frac{|S_{p_i} \cap S_{g_j}|}{|S_{p_i}|} \quad (4)$$

$$R_{ent}(p_i, g_j) = \frac{|S_{p_i} \cap S_{g_j}|}{|S_{g_j}|} \quad (5)$$

$$F1_{ent}(p_i, g_j) = \frac{2 \cdot P_{ent}(p_i, g_j) \cdot R_{ent}(p_i, g_j)}{P_{ent}(p_i, g_j) + R_{ent}(p_i, g_j)} \quad (6)$$

where S_{p_i} and S_{g_j} are the sets of entity spans.

The optimal subject matching M^* is determined by maximizing the strict F1 score:

$$M^* = \max_{M \in \mathcal{M}} \sum_{(p_i, g_j) \in M} F1_{strict}(p_i, g_j) \quad (7)$$

where \mathcal{M} denotes all possible one-to-one mappings between predicted and ground truth subjects.

The final recognition scores are computed over the optimal matching pairs:

$$P_{strict} = \frac{1}{|\mathcal{P}|} \sum_{(p_i, g_j) \in M^*} P_{strict}(p_i, g_j) \quad (8)$$

$$R_{strict} = \frac{1}{|\mathcal{G}|} \sum_{(p_i, g_j) \in M^*} R_{strict}(p_i, g_j) \quad (9)$$

$$F1_{strict} = \frac{1}{\max(\mathcal{P}, \mathcal{G})} \sum_{(p_i, g_j) \in M^*} F1_{strict}(p_i, g_j) \quad (10)$$

P_{span} , R_{span} , and $F1_{span}$ are computed analogously.

For query-related detection, given a predicted entity set \mathcal{E}_p and ground truth set \mathcal{E}_g , we compute:

$$P_{query} = \frac{|\mathcal{E}_p \cap \mathcal{E}_g|}{|\mathcal{E}_p|} \quad (11)$$

$$R_{query} = \frac{|\mathcal{E}_p \cap \mathcal{E}_g|}{|\mathcal{E}_g|} \quad (12)$$

$$F1_{query} = \frac{2 \cdot P_{query} \cdot R_{query}}{P_{query} + R_{query}} \quad (13)$$

For both PII detection and query-related detection tasks, we additionally employ Rouge-L based fuzzy matching to handle partial matches between entity spans. Instead of using exact set intersection, the Rouge-L score is used to measure textual similarity between entities:

$$P_{fuzzy} = \frac{1}{|\mathcal{E}_p|} \sum_{e_p \in \mathcal{E}_p} \max_{e_g \in \mathcal{E}_g} Rouge-L(e_p, e_g) \quad (14)$$

$$R_{fuzzy} = \frac{1}{|\mathcal{E}_g|} \sum_{e_g \in \mathcal{E}_g} \max_{e_p \in \mathcal{E}_p} Rouge-L(e_p, e_g) \quad (15)$$

$$F1_{fuzzy} = \frac{2 \cdot P_{fuzzy} \cdot R_{fuzzy}}{P_{fuzzy} + R_{fuzzy}} \quad (16)$$

where $Rouge-L(e_p, e_g)$ computes the longest common subsequence-based F-score between predicted entity e_p and ground truth entity e_g .

F.2 Privacy-Utility Tradeoff Metrics

We develop a comprehensive evaluation framework to quantify the tradeoff between privacy protection and query utility across different PII masking strategies.

F.2.1 Privacy Protection Metric

Let $\mathcal{E} = \{e_1, \dots, e_n\}$ denote the set of PII entities in the original text T_o , and T_m represent the masked text. The privacy score P measures the proportion of PII entities successfully protected:

$$P = 1 - \frac{\sum_{e \in \mathcal{E}} C(e, T_m)}{\sum_{e \in \mathcal{E}} C(e, T_o)} \quad (17)$$

where $C(e, T)$ counts occurrences of entity e in text T . A score of $P = 1$ indicates complete protection, while $P = 0$ indicates no protection.

F.2.2 Utility Preservation Metrics

To measure utility preservation, we employ two complementary approaches:

Semantic Similarity We compute embedding-based similarity between responses generated from masked prompts (R_m) and unmasked prompts (R_o):

$$U_s = \cos(\mathbf{v}_{R_m}, \mathbf{v}_{R_o}) \quad (18)$$

where \mathbf{v}_R is the text embedding generated by BGE-M3 (Multi-Granularity, 2024), and $\cos(\cdot, \cdot)$ computes cosine similarity. This metric captures semantic preservation independent of exact wording.

LLM-as-Judge Evaluation We employ Claude-3.7-Sonnet to assess response quality through direct comparison:

$$U_l = \text{Judge}(R_m^1, R_m^2, R_o) \quad (19)$$

where R_m^1 and R_m^2 represent responses from two different masking strategies, and R_o is the reference response from unmasked text. The judge assigns numerical ratings $r \in [1, 10]$ to each masked response based on how well it preserves the query intent compared to the reference. To mitigate position bias, we randomly alternate the presentation order of responses.

F.2.3 Balanced Metric

To quantify the overall effectiveness of masking strategies, we compute a balanced score B that combines privacy protection and utility preservation:

$$B = \alpha P + (1 - \alpha)U \quad (20)$$

where $\alpha \in [0, 1]$ is a weighting parameter that determines the relative importance of privacy versus utility. In our experiments, we use $\alpha = 0.5$ to assign equal importance to both aspects.

The complete prompt template for LLM-as-Judge evaluation is provided in Figure 20.

F.3 Additional Experiments with Smaller Language Models

We have conducted a comprehensive evaluation on smaller, deployment-ready models (0.5B-3B parameters) to assess their viability for on-device PII protection.

As shown in Tables 13, 14, and 12, our results reveal a significant performance gap between

deployment-ready models and larger proprietary systems. The 0.5B and 1.5B models exhibit extremely limited capabilities across all tasks (with F1 scores generally below 0.2), rendering them practically unusable for real-world PII protection. While the 3B model shows modest potential (F1 scores approaching 0.6 on simpler datasets), its performance remains substantially inferior to larger models, particularly on challenging scenarios (e.g., PII-Hard, PII-Distract). These findings underscore the tension between privacy goals and model capabilities: while smaller on-device models would be preferable from a privacy perspective, they currently lack the sophistication needed for reliable PII management.

F.4 Supplementary Results on PII-Real Dataset

To empirically validate PII-Bench’s ability to capture real-world privacy challenges, we constructed PII-Real, a validation dataset based on authentic biographical information.

F.4.1 Dataset Construction

We selected 20 prominent AI researchers from the AMiner AI2000 ranking¹ and manually curated their professional profiles from publicly available sources including conference websites, institutional pages, and academic publications. Each profile underwent expert annotation by the authors following our established PII taxonomy, ensuring consistency with PII-Bench guidelines.

Unlike PII-Bench’s automated query generation, PII-Real features human-written queries crafted by domain experts to reflect authentic information needs. We designed queries spanning career planning, research collaboration, academic advising, and startup leadership scenarios for each profile. The resulting dataset comprises 100 instances with manually annotated PII entities, providing an authentic testbed for validating synthetic data quality.

Tables 18, 19, and 20 present experimental results across three evaluation tasks.

F.4.2 Consistency Analysis

We quantified alignment between PII-Single and PII-Real using Spearman’s rank correlation coefficient (ρ) and Performance Consistency Score (PCS).

¹<https://www.aminer.cn/ai2000>

Baseline Models	PII-Single			PII-Multi			PII-Hard			PII-Distract		
	Strict-F1	Ent-F1	RougeL-F	Strict-F1	Ent-F1	RougeL-F	Strict-F1	Ent-F1	RougeL-F	Strict-F1	Ent-F1	RougeL-F
Qwen2.5-0.5B	0.002	0.0029	0.002	0.0013	0.0042	0.0013	0.0034	0.0042	0.0034	0.009	0.0207	0.009
Qwen2.5-1.5B	0.1846	0.2069	0.1865	0.1057	0.1794	0.1071	0.1474	0.1976	0.1502	0.0728	0.1549	0.0733
Qwen2.5-3B	0.5929	0.6289	0.5962	0.6189	0.7067	0.6212	0.5202	0.5921	0.5237	0.3717	0.6293	0.3731

Table 12: Performance of baseline models under the PII Detection task.

Method	Qwen2.5-0.5B		Qwen2.5-1.5B		Qwen2.5-3B	
	F1	RougeL-F	F1	RougeL-F	F1	RougeL-F
<i>Basic Method</i>						
Naive	0.0041	0.0041	0.2131	0.2185	0.1691	0.1699
<i>Advanced Method</i>						
Self-CoT	0.009	0.0095	0.2002	0.2051	0.2884	0.2909
Auto-CoT(3-shot)	0.0258	0.0268	0.2383	0.2466	0.381	0.384
Self-Consistency	0.0018	0.0018	0.1057	0.1088	0.2694	0.2774
PS-CoT	0.0088	0.009	0.1887	0.1959	0.293	0.2956
<i>w/ PII Detection</i>						
Naive w/ Choice	0.3978	0.3986	0.4357	0.4357	0.542	0.5437

Table 13: Performance comparison on the Query-Related PII Detection task (PII-single dataset).

Method	Qwen2.5-0.5B		Qwen2.5-1.5B		Qwen2.5-3B	
	F1	RougeL-F	F1	RougeL-F	F1	RougeL-F
<i>Basic Method w/ PII Detection</i>						
Naive	0.0016	0.0025	0.0912	0.1018	0.3257	0.3321
<i>Advanced Method w/ PII Detection</i>						
Self-CoT	0.0016	0.0029	0.0846	0.0945	0.3764	0.3841
Auto-CoT(3-shot)	0.0019	0.0032	0.0927	0.1038	0.4247	0.4327
Self-Consistency	0.0016	0.0029	0.0907	0.1020	0.3414	0.3486
PS-CoT	0.0016	0.0025	0.0920	0.1033	0.3765	0.3840
<i>w/ PII Detection</i>						
Naive w/ Choice	0.0011	0.0017	0.0722	0.0800	0.3918	0.3994

Table 14: Performance comparison on the Query-Unrelated PII Masking task (PII-single and PII-multi datasets).

Performance Consistency Score PCS measures absolute performance alignment between two datasets. For a given model and evaluation metric, PCS is defined as:

$$PCS = 1 - \frac{|F1_{single} - F1_{real}|}{\max(F1_{single}, F1_{real})} \quad (21)$$

where $F1_{single}$ represents the F1 score on PII-Single dataset and $F1_{real}$ represents the F1 score on PII-Real dataset. A PCS value of 1 indicates perfect performance consistency, while lower values indicate greater deviation between synthetic and real-world scenarios.

Ranking Consistency Table 15 presents Spearman’s ρ values across all evaluation tasks. The exceptionally high correlation coefficients ($\rho > 0.988$, $p < 0.001$) demonstrate that model rankings remain consistent between synthetic and real-world datasets, indicating PII-Bench reliably predicts relative model performance in authentic privacy scenarios.

Task	Spearman’s ρ	P-value
PII Detection	0.9880	< 0.001
Query-Related Detection	0.9982	< 0.001
Query-Unrelated Masking	0.9982	< 0.001

Table 15: Spearman’s rank correlation between PII-Single and PII-Real datasets across evaluation tasks.

Performance Consistency Tables 16 and 17 present PCS values quantifying absolute performance alignment. For PII detection (Table 16), all models achieve PCS values exceeding 0.967, with an overall average of 0.980. For query-related detection and masking tasks (Table 17), PCS values consistently exceed 0.947 across all model-strategy combinations, averaging 0.966. These high consistency scores confirm that synthetic single-entity scenarios accurately capture the challenges present in real-world privacy protection tasks.

Model	Strict-F1	Ent-F1	RougeL-F	Average
GPT4o	0.979	0.979	0.979	0.979
Claude3.5	0.967	0.978	0.970	0.972
DeepSeekV3	0.978	0.979	0.978	0.978
Llama3.1	0.978	0.984	0.978	0.980
Qwen2.5	0.980	0.988	0.980	0.982
Llama3.1-SLM	0.982	0.984	0.982	0.982
Qwen2.5-SLM	0.986	0.988	0.986	0.987
Overall Average	0.979	0.983	0.979	0.980

Table 16: Performance Consistency Scores for PII detection task across different metrics.

F.4.3 Representative Example from PII-Real

Figure 6 presents a representative instance from the PII-Real dataset, illustrating the comparable complexity and structure to synthetic samples. This example is derived from a real-world biographical profile.²

F.5 Additional Results

Table 21 compares different prompting strategies on PII-multi dataset.

²Original source: <https://kimiyoung.github.io>

Model	Task	Naive	Self-CoT	Auto-CoT	Self-Consistency	PS-CoT	Naive w/ Choice	Average
GPT4o	Query-Related	0.966	0.962	0.955	0.966	0.952	0.976	0.963
	Masking	0.969	0.972	0.973	0.970	0.976	0.979	0.973
Llama3.1	Query-Related	0.972	0.969	0.967	0.960	0.966	0.972	0.968
	Masking	0.976	0.977	0.973	0.970	0.966	0.967	0.972
Qwen2.5	Query-Related	0.976	0.965	0.974	0.969	0.972	0.980	0.973
	Masking	0.971	0.972	0.972	0.967	0.969	0.973	0.971
Llama3.1-SLM	Query-Related	0.954	0.956	0.953	0.948	0.951	0.967	0.955
	Masking	0.963	0.969	0.968	0.965	0.964	0.962	0.965
Qwen2.5-SLM	Query-Related	0.960	0.957	0.951	0.947	0.962	0.973	0.958
	Masking	0.968	0.963	0.961	0.957	0.969	0.972	0.965
Overall Average		0.968	0.966	0.965	0.962	0.965	0.972	0.966

Table 17: Performance Consistency Scores for query-related detection and masking tasks by prompting strategy.

User Description:
“Hi, I’m Zhilin Yang. I am working on a startup and I am the CEO of Moonshot AI. In 2019, I obtained my PhD degree from Carnegie Mellon University, advised by Ruslan Salakhutdinov and William W. Cohen. Prior to that, in 2015, I received my bachelor’s degree from Tsinghua University, advised by Jie Tang. I worked at Meta AI with Jason Weston, and Google Brain with Quoc V. Le.”

Query:
“In my current role, we’re early and juggling hiring with first product bets. What minimal weekly rhythm and decision rules keep us fast without creating chaos?”

PII Entities (Total: 16):

- **PER:** Zhilin Yang, Ruslan Salakhutdinov, William W. Cohen, Jie Tang, Jason Weston, Quoc V. Le
- **DEM:** CEO, PhD degree, bachelor’s degree
- **CODE:** (none)
- **ORG:** Moonshot AI, Carnegie Mellon University, Tsinghua University, Meta AI, Google Brain
- **DATETIME:** 2019, 2015
- **LOC:** (none)
- **QUANTITY:** (none)

Query-Related PII: CEO, Moonshot AI

Figure 6: Example from PII-Real dataset showing a real-world biographical profile with comprehensive PII annotations. The query specifically targets career-related information, requiring identification of occupation and organizational affiliation while protecting other personal details.

Model	Strict-F1	Ent-F1	RougeL-F
<i>API-based Large Language Model</i>			
GPT4o	0.912	0.934	0.914
Claude3.5	0.887	0.911	0.889
DeepSeekV3	0.923	0.941	0.925
<i>Open-source Large Language Model</i>			
Llama3.1	0.901	0.928	0.903
Qwen2.5	0.884	0.919	0.887
<i>Open-source Small Language Model</i>			
Llama3.1-SLM	0.762	0.813	0.766
Qwen2.5-SLM	0.798	0.856	0.803

Table 18: Performance of baseline models on PII detection task (PII-Real dataset).

F.6 Prompt Details

This section presents the prompts used throughout our experiments. For the PII detection task, we employ the template shown in Figure 13. For query-related PII detection, we design and evaluate six distinct prompting strategies. Figure 14 displays the **Naive** prompts, Figure 15 presents the **Naive w/ Choice** prompts, Figure 16 features the **Self-CoT** prompts, Figure 17 reveals the **Auto-CoT** prompts, Figure 18 exhibits the **Self-Consistency** prompts, and Figure 19 displays the **PS-CoT** prompts.

G PII Annotation System

We developed a specialized web-based annotation platform to facilitate the systematic evaluation of PII detection and query-related detection capabilities. The platform implements a two-stage annotation process, ensuring comprehensive coverage of both fundamental PII entity identification and contextual relevance assessment.

G.1 PII Detection Interface

As shown in Figure 7, the PII detection interface enables annotators to identify and categorize PII entities within user descriptions. The interface provides the following key functionalities:

- **Entity Detection:** Annotators can highlight text spans containing PII entities directly in the user description.
- **Type Classification:** Each identified entity is assigned a specific PII type (e.g., PER for person names, ORG for organizations, LOC for locations).

- **Subject Association:** Entities are linked to their corresponding subjects using alphabetical identifiers (e.g., A, B) to maintain relationship clarity in multi-subject scenarios.
- **Span Verification:** The interface displays start and end positions for each entity span, ensuring precise boundary detection.

G.2 Query-Related PII Detection Interface

Figure 8 illustrates the interface for query-related PII detection, which builds upon the recognition results to assess contextual relevance:

- **Query Context:** The interface presents both the user description and the associated query, providing complete context for relevance assessment.
- **Entity Selection:** Annotators identify PII entities crucial for addressing the query, with the interface highlighting pre-identified entities from the recognition phase.
- **Subject Verification:** For selected query-related entities, annotators must verify the subject associations to ensure consistency across tasks.
- **Relevance Validation:** The interface includes a review mechanism to confirm that selected entities are both necessary and sufficient for query resolution.

G.3 Query-Unrelated PII Masking Visualization

To validate the effectiveness of privacy protection while maintaining query relevance, we implemented a masking visualization interface (Figure 9):

- **Original Context:** Displays the complete user description with all PII entities highlighted.
- **Masked View:** Shows the description with non-relevant PII entities replaced by their corresponding type tags (e.g., <Nickname>, <Phone Number>).
- **Key Information Display:** Preserves query-related PII entities while maintaining readability and semantic coherence.

Method	GPT4o		Llama3.1		Qwen2.5		Llama3.1-SLM		Qwen2.5-SLM	
	F1	RougeL-F	F1	RougeL-F	F1	RougeL-F	F1	RougeL-F	F1	RougeL-F
<i>Basic Method</i>										
Naive	0.652	0.654	0.648	0.651	0.635	0.638	0.346	0.348	0.427	0.429
<i>Advanced Method</i>										
Self-CoT	0.738	0.741	0.712	0.714	0.694	0.697	0.408	0.411	0.418	0.421
Auto-CoT	0.691	0.693	0.724	0.728	0.729	0.732	0.451	0.454	0.389	0.392
Self-Consistency	0.745	0.747	0.656	0.659	0.671	0.674	0.327	0.331	0.338	0.342
PS-CoT	0.683	0.685	0.673	0.676	0.689	0.691	0.368	0.372	0.468	0.471
<i>w/ Extra Information</i>										
Naive w/ Choice	0.861	0.861	0.782	0.784	0.847	0.849	0.538	0.541	0.791	0.793

Table 19: Performance comparison on Query-Related PII Detection task (PII-Real dataset).

Method	GPT4o		Llama3.1		Qwen2.5		Llama3.1-SLM		Qwen2.5-SLM	
	F1	RougeL-F	F1	RougeL-F	F1	RougeL-F	F1	RougeL-F	F1	RougeL-F
<i>Basic Method</i>										
Naive	0.743	0.745	0.738	0.742	0.721	0.725	0.436	0.441	0.558	0.593
<i>Advanced Method</i>										
Self-CoT	0.782	0.785	0.768	0.771	0.751	0.754	0.547	0.551	0.561	0.596
Auto-CoT	0.771	0.773	0.781	0.784	0.782	0.785	0.589	0.593	0.562	0.597
Self-Consistency	0.794	0.796	0.732	0.736	0.734	0.738	0.508	0.513	0.512	0.548
PS-CoT	0.758	0.761	0.745	0.748	0.753	0.756	0.498	0.503	0.578	0.614
<i>w/ Extra Information</i>										
Naive w/ Choice	0.838	0.841	0.796	0.799	0.812	0.815	0.478	0.483	0.689	0.725

Table 20: Performance comparison on Query-Unrelated PII Masking task (PII-Real dataset).

Method	GPT4o		Llama3.1		Qwen2.5		Llama3.1-SLM		Qwen2.5-SLM	
	F1	RougeL-F	F1	RougeL-F	F1	RougeL-F	F1	RougeL-F	F1	RougeL-F
<i>Basic Method</i>										
Naive	0.600	0.602	0.611	0.614	0.596	0.603	0.240	0.333	0.405	0.413
<i>Advanced Method</i>										
Self-CoT	0.675	0.681	0.638	0.643	0.626	0.632	0.354	0.362	0.392	0.397
Auto-CoT(3-shot)	0.629	0.640	0.650	0.662	0.657	0.665	0.393	0.402	0.391	0.394
Self-Consistency	0.685	0.692	0.602	0.605	0.614	0.620	0.263	0.269	0.288	0.293
PS-CoT	0.618	0.620	0.624	0.631	0.636	0.643	0.291	0.300	0.431	0.436
<i>w/ Extra Information</i>										
Naive w/ Choice	0.846	0.846	0.775	0.775	0.804	0.804	0.387	0.388	0.743	0.743

Table 21: Performance comparison on the Query-Related PII Detection task (PII-multi dataset).

G.4 Annotation Guidelines

To ensure annotation consistency and quality, we established comprehensive guidelines for each task. Table 24 summarizes the core annotation instructions provided to annotators across the three tasks.

G.5 Quality Control and Inter-Annotator Agreement

To maintain high annotation quality, we implemented a rigorous quality control protocol. Each sample was independently annotated by multiple annotators, with disagreements resolved through majority voting or expert review. Regular review sessions were conducted to discuss challenging cases and update guidelines based on annotator feedback. For quality assurance, we randomly sampled 10% of the annotations for expert review.

Annotation Task	Agreement (%)	Fleiss' κ
PII Detection	95.1	0.912
Query-Related Detection	91.5	0.873

Table 22: Inter-annotator agreement for PII detection and query-related detection tasks.

Table 22 presents the inter-annotator agreement results for the two primary annotation tasks, measured using both observed agreement rates and Fleiss' kappa. We obtain kappa values of 0.912 for PII detection and 0.873 for query-related detection across all annotators.

H Evaluation of Reasoning-Capable Models

We evaluate seven recent LLMs that pair reasoning and non-reasoning variants: the proprietary reasoning model **GPT-5**, the DeepSeek pair (non-reasoning **DeepSeek-V3.2** and reasoning **DeepSeek-R1**), and two Qwen3 scales with matched Instruct and Thinking variants (**Qwen3-30B-Instruct**, **Qwen3-30B-Thinking**, **Qwen3-8B (No-Think)**, **Qwen3-8B (Think)**). Open-source models are served locally via vLLM; API-based models are accessed through their official endpoints. All decoding settings follow the main experiments (temperature 0, top- k 1).

Tables 23, 25 and 26 report Strict-F1 on PII Detection across all four subsets, and F1 / RougeL-F on PII-Single for Query-Related PII Detection and Query-Unrelated PII Masking under six prompting strategies. We summarise the results along three axes.

Model	Single	Multi	Hard	Distract
GPT-5	0.871	0.884	0.840	0.832
DeepSeek-V3.2	0.856	0.838	0.779	0.719
DeepSeek-R1	0.876	0.879	0.835	0.826
Qwen3-30B-Instruct	0.793	0.781	0.745	0.510
Qwen3-30B-Thinking	0.805	0.798	0.781	0.625
Qwen3-8B (No-Think)	0.771	0.684	0.667	0.485
Qwen3-8B (Think)	0.740	0.684	0.679	0.571

Table 23: Strict-F1 on the PII Detection task using the Naive (zero-shot) prompt, across the four PII-Bench subsets. Within each family, the non-reasoning variant precedes its reasoning counterpart; best per subset in **bold**.

(1) Reasoning helps most in complex detection settings. On the easier PII-Single and PII-Multi subsets, basic detection is close to saturation across all seven systems, so reasoning capability provides little additional benefit. The picture changes on the high-entropy PII-Distract subset, where Qwen3-30B-Thinking improves over Qwen3-30B-Instruct by 11.5 Strict-F1 points (0.625 vs. 0.510, a 22.5% relative gain), with a comparable gap for the 8B pair (0.571 vs. 0.485). This indicates that reasoning becomes beneficial only when detection itself is non-trivial.

(2) Open-source reasoning reaches proprietary parity. On PII-Single detection, DeepSeek-R1 (0.876) slightly surpasses GPT-5 (0.871), and the two models track each other within about one F1 point on every other PII Detection subset.

(3) The query-aware tasks remain considerably harder. No evaluated model exceeds 0.77 F1 on Query-Related PII Detection, which still trails the 0.951 human baseline by more than 18 F1 points. Reasoning offers no consistent lift on PII-Single for these tasks: the non-reasoning DeepSeek-V3.2 is marginally ahead of DeepSeek-R1 (0.765 vs. 0.745 on Query-Related, 0.852 vs. 0.850 on Masking). On-device variants add a second gap, with the best Qwen3-8B trailing DeepSeek-R1 by 14.5 F1 points on Query-Related PII Detection. Relevance judgement and on-device deployment thus emerge as the two principal challenges for future query-aware privacy systems.

Task	Annotation Guidelines
PII Detection	<ul style="list-style-type: none"> • Verify whether all PII entities in the user description are correctly identified. • Click the <i>Correct</i> button if all entities are properly detected. • Annotate missing PII entities by selecting text spans and assigning type, tag, and subject identifiers. • Assign the next available letter (A, B, C) for entities belonging to new subject groups. • Select minimal spans containing only the essential text representing each entity.
Query-Related PII Detection	<ul style="list-style-type: none"> • Identify PII entities crucial for addressing the query from the provided options. • Verify or correct the subject association (A, B, C, etc.) for each selected entity. • Ensure that selected entities are both necessary and sufficient for query resolution. • Use the review mechanism to validate your selection before submission.
Query-Unrelated PII Masking	<ul style="list-style-type: none"> • Review the masked user description where query-unrelated PII entities are replaced with type tags. • Verify that query-related entities are preserved in their original form. • Confirm that masked entities are replaced with appropriate tags (e.g., <Nickname>, <Phone Number>). • Evaluate whether the masked text maintains semantic coherence and readability. • Assess the balance between privacy protection and query utility preservation.

Table 24: Comprehensive annotation guidelines for the three tasks in PII-Bench. Each task follows specific protocols to ensure consistency and quality across annotators.

Method	GPT-5		DeepSeek-V3.2		DeepSeek-R1		Qwen3-30B-Ins.		Qwen3-30B-Thk.		Qwen3-8B (NT)		Qwen3-8B (T)	
	F1	RougeL-F	F1	RougeL-F	F1	RougeL-F	F1	RougeL-F	F1	RougeL-F	F1	RougeL-F	F1	RougeL-F
Naive	0.700	0.704	0.735	0.737	0.715	0.718	0.545	0.550	0.612	0.617	0.362	0.363	0.494	0.498
Self-CoT	0.733	0.740	0.765	0.767	0.745	0.749	0.632	0.638	0.698	0.703	0.524	0.525	0.534	0.537
Auto-CoT	0.697	0.710	0.746	0.754	0.710	0.718	0.692	0.696	0.741	0.746	0.538	0.540	0.600	0.608
Self-Consistency	0.722	0.725	0.741	0.745	0.735	0.738	0.631	0.635	0.689	0.693	0.501	0.502	0.529	0.532
PS-CoT	0.725	0.729	0.754	0.756	0.739	0.742	0.606	0.612	0.675	0.680	0.443	0.444	0.552	0.554
Naive w/ Choice	0.872	0.872	0.857	0.857	0.879	0.879	0.816	0.816	0.845	0.845	0.638	0.638	0.643	0.643

Table 25: Query-Related PII Detection on PII-Single under the six prompting strategies defined in §F.6. Within each family, the non-reasoning variant precedes its reasoning counterpart; best F1 per model across the five non-oracle strategies in **bold**. “Naive w/ Choice” provides the gold PII candidate set and serves as an oracle upper bound.

Method	GPT-5		DeepSeek-V3.2		DeepSeek-R1		Qwen3-30B-Ins.		Qwen3-30B-Thk.		Qwen3-8B (NT)		Qwen3-8B (T)	
	F1	RougeL-F	F1	RougeL-F	F1	RougeL-F	F1	RougeL-F	F1	RougeL-F	F1	RougeL-F	F1	RougeL-F
Naive	0.832	0.835	0.831	0.833	0.841	0.844	0.683	0.687	0.765	0.770	0.517	0.522	0.628	0.633
Self-CoT	0.840	0.843	0.852	0.854	0.850	0.853	0.794	0.797	0.845	0.849	0.641	0.646	0.705	0.710
Auto-CoT	0.836	0.839	0.850	0.852	0.847	0.850	0.815	0.818	0.852	0.856	0.656	0.661	0.720	0.725
Self-Consistency	0.824	0.827	0.831	0.833	0.835	0.838	0.768	0.771	0.825	0.829	0.621	0.626	0.692	0.697
PS-CoT	0.834	0.836	0.846	0.848	0.844	0.847	0.772	0.775	0.831	0.835	0.593	0.597	0.665	0.670
Naive w/ Choice	0.854	0.856	0.861	0.862	0.862	0.864	0.805	0.808	0.841	0.845	0.607	0.613	0.698	0.703

Table 26: Query-Unrelated PII Masking on PII-Single under the six prompting strategies defined in §F.6. Within each family, the non-reasoning variant precedes its reasoning counterpart; best F1 per model across the five non-oracle strategies in **bold**. “Naive w/ Choice” provides the gold PII candidate set and serves as an oracle upper bound.

PII Detection Task

Input

Enter text: (User Background/Description)

Lance is a dedicated researcher at the Global Health Research Center, located on the 8th Floor, No. 200, George Street, Sydney. Known for his meticulous work ethic, Lance often finds himself engrossed in his projects well past the usual working hours. Despite his professional success, Lance manages a personal challenge—hypertension, which he diligently monitors. Every day at precisely 3:14 PM, he takes a moment to check his blood pressure, ensuring he stays on top of his health. Outside of work, Lance is known for his approachable nature and is always just a phone call away at 14593742147, ready to lend an ear or share his insights on global health issues.

Highlighted

Detected Entities

Lance PER is a dedicated researcher at the Global Health Research Center ORG, located on the 8th Floor, No. 200, George Street, Sydney LOC. Known for his meticulous work ethic, Lance PER often finds himself engrossed in his projects well past the usual working hours. Despite his professional success, Lance PER manages a personal challenge—hypertension DEM, which he diligently monitors. Every day at precisely 3:14 PM DATETIME, he takes a moment to check his blood pressure, ensuring he stays on top of his health. Outside of work, Lance PER is known for his approachable nature and is always just a phone call away at 14593742147 CODE, ready to lend an ear or share his insights on global health issues.

Detected PII Entities

Entity	Type	Tag	Start	End	Subject
0 Lance	PER	Nickname	0	5	A
1 Global Health Research Center	ORG	Non-Governmental Organization Name	39	68	A
2 8th Floor, No. 200, George Street, Sydney	LOC	Work or Home Detailed Address	85	126	A
3 Lance	PER	Nickname	165	170	A
4 Lance	PER	Nickname	286	291	A
5 hypertension	DEM	Health Status	321	333	A
6 3:14 PM	DATETIME	Specific Time	388	395	A
7 Lance	PER	Nickname	501	506	A
8 14593742147	CODE	Phone Numbers	584	595	A

Task Definition

Please check if all PII entities in the user description are correctly identified.

If correct, click the ✓ Correct button. Otherwise, annotate any missing PII entities.

Human Annotation

Please add any missing PII entities. For each entity, select its type, tag, and subject group. If the entity belongs to a new subject group, select the next available letter.

Entity type:

PER

Subject:

A

Entity tag:

Name

Add new entity:

Figure 7: Web Demo for the PII Detection Task

Query-Related PII Detection Task

User Description

Hello, I'm longjie PER, a 67kg DEM advocate for global harmony working with the World Peace Organization ORG. I often find myself reflecting on life's journey while enjoying the breathtaking views from Table Mountain in Cape Town LOC. My evenings are usually spent at 8:40 PM DATETIME, contemplating the 50 years DATETIME of progress in peace initiatives. You can reach me at xiaolu@example.net CODE or call me at 18180989411 CODE. My credit score is 76.5/100 QUANTITY, and I frequently collaborate with Sydney Prince Hospital ORG on health-related projects. My daughter, dengna PER, is a distinguished Doctor of Clinical Medicine DEM who has dedicated 23 years DATETIME to the Transnational Health Association ORG in the United Kingdom LOC. She resides at 5th Floor, No. 65, Labor West Road, Tianxin District, Changsha LOC, and is currently managing a loan of C274304.33 QUANTITY. Her expertise is further honed at Moscow First Hospital ORG, and she collaborates with Casio ORG on health technology projects. She often visits Krishna Fort LOC to unwind and gather inspiration for her work.

Scene: Technological Innovations in Peacekeeping

Query Text:

In what ways can my extensive experience in fostering global harmony, combined with her collaborations in health technology, contribute to innovative solutions in peacekeeping efforts? How might our respective organizational affiliations enhance the integration of cutting-edge tools in this field?

中文翻译:

我在促进全球和谐方面的丰富经验，再加上她在卫生技术方面的合作，可以通过哪些方式为维和工作的创新解决方案做出贡献？我们各自的组织关系如何加强该领域尖端工具的综合？

Key PII Information

The key PII is: 23 years | Casio | 76.5/100 | World Peace Organization | Transnational Health Association

Human Annotation

Task Definition: Select the most related PII to the query from the following options.

Please verify or correct your selection based on the correct answer:

23 years ✗ Casio ✗ 76.5/100 ✗ World Peace Org. ✗ Transnational H. ✗

Task Definition: For each selected PII, please identify and annotate its subject(s).

Entity: 23 years

B

Entity: Casio

B

Entity: 76.5/100

A

Entity: World Peace Organization

A

Entity: Transnational Health Association

B

Figure 8: Web Demo for the Query-Related PII Detection Task

Query-Unrelated PII Masking Method

Adaptive PII Mask Method intelligently protects user privacy while maintaining query relevance.

Original User Description

Lance PER is a dedicated researcher at the Global Health Research Center ORG, located on the 8th Floor, No. 200, George Street, Sydney LOC. Known for his meticulous work ethic, Lance PER often finds himself engrossed in his projects well past the usual working hours. Despite his professional success, Lance PER manages a personal challenge—hypertension DEM, which he diligently monitors. Every day at precisely 3:14 PM DATETIME, he takes a moment to check his blood pressure, ensuring he stays on top of his health. Outside of work, Lance PER is known for his approachable nature and is always just a phone call away at 14593742147 CODE, ready to lend an ear or share his insights on global health issues.

Query

Query Text:

Given my routine health check in the afternoon and my commitment to my current office, how can I efficiently schedule a medical consultation without disrupting my responsibilities at the organization I am part of?

中文翻译:

考虑到我下午的例行健康检查和我对目前办公室的承诺，我如何在干扰我在所属组织职责的情况下有效地安排医疗咨询？

Key PII Information

Key PII: 8th Floor, No. 200, George Street, Sydney | Hypertension | Global Health Research Center | 3:14 PM

Masked User Description

Non-essential PII entities are masked with their corresponding tags and only query-relevant PII information is preserved

<Nickname> PER is a dedicated researcher at the Global Health Research Center ORG, located on the 8th Floor, No. 200, George Street, Sydney LOC. Known for his meticulous work ethic, <Nickname> PER often finds himself engrossed in his projects well past the usual working hours. Despite his professional success, <Nickname> PER manages a personal challenge—hypertension, which he diligently monitors. Every day at precisely 3:14 PM DATETIME, he takes a moment to check his blood pressure, ensuring he stays on top of his health. Outside of work, <Nickname> PER is known for his approachable nature and is always just a phone call away at <Phone Numbers> CODE, ready to lend an ear or share his insights on global health issues.

Select Engine:

glm-4-flash

Max Tokens

512

API Key:

Temperature

0.00 1.00 2.00

Top P

0.00 0.80 1.00

Figure 9: Web Demo for the Query-Unrelated PII Masking Method

Consistency Optimization Prompt of Single Subject

You are a character feature selector tasked with identifying and refining logically consistent feature combinations. I will provide you with character features. Your role is to identify any features that have obvious logical conflicts or inconsistencies, and modify them to create a coherent set while preserving their core classifications.

Requirements:

1. The selected character features must be logically consistent with real-world expectations, with no obvious conflicts.
2. When resolving conflicts, modify only the feature entities while keeping their PII types and classifications unchanged.
3. Modified feature entities must remain within the same PII type and classification categories as their originals.
4. Aim to maintain as many features as possible, ideally matching the original count or coming as close as feasible.

Common conflict patterns to identify and resolve:

- Temporal inconsistencies: Age incompatible with work experience duration, education level, or career stage (e.g., a 23-year-old with 15 years of work experience should be adjusted to 2 years as Junior position).
- Professional mismatches: Occupation inconsistent with education background or salary range (e.g., a high school graduate working as a licensed physician).
- Geographic contradictions: Work location distant from residential address without supporting evidence (e.g., daily commute spanning different continents).
- Financial implausibility: Income, expenses, and savings that violate basic economic constraints (e.g., monthly expenses exceeding monthly income by orders of magnitude).

Character Features

```
<PII Type> <Entity Category> <PII Entity>
{usr_features}
```

Please provide your output in the following format:

- Under "## Reason:", explain your selection and modification process, specifically addressing any conflict patterns identified above
- Under "## Final Features:", list the final selected features as JSON objects in the format `{{"label": xxx, "tag": yyy, "entity": zzz}}` with no additional content or line breaks where xxx is the PII type, yyy is the entity category, and zzz is the PII entity.

Reason: [Explain your selection and modification process]

Final Features: [{"label": xxx, "tag": yyy, "entity": zzz}], [{"label": xxx, "tag": yyy, "entity": zzz}], ...]""

Figure 10: Prompt of Consistency Optimization for Single-Subject

Consistency Optimization Prompt of Multi Subject

You are a character feature selector tasked with identifying and refining logically consistent feature combinations. I will provide you with character features for different subjects and their relationships. Your role is to identify conflicts and modify features to align with the relationship while preserving PII types and categories.

Requirements:

1. Selected features must be logically consistent and align with the relationship between subjects.
2. Apply relationship-aware constraints:
 - Intersection relationships (colleagues, classmates): Subjects share critical attributes (e.g., same organization, institution).
 - Hierarchical relationships (parent-child, supervisor-subordinate): Enforce age differences (parent \geq 18 years older), authority levels, and derived attributes when appropriate.
 - Non-intersection relationships (strangers): Maintain independent profiles, but ensure internal consistency for each subject.
3. When modifying: Keep PII type and category unchanged, only adjust entity values.
4. Maximize feature retention: Aim to preserve the original count.

Subject A Features

<PII Type> <Entity Category> <PII Entity>
{usr_features_a}

Subject B Features

<PII Type> <Entity Category> <PII Entity>
{usr_features_b}

Relationship Between Subjects

{rel}

Please provide your output in the following format:

- Under "## Reason:", explain your selection and modification process, addressing how relationship constraints are enforced.
- Under "## Final Features A:" and "## Final Features B:", list the final selected features as JSON objects in the format `{{"label": xxx, "tag": yyy, "entity": zzz}}`.

Reason: [Explain your selection and modification process]

Final Features A: [{{"label":xxx, "tag": yyy, "entity": zzz}}, {"label":xxx, "tag": yyy, "entity": zzz}}, ...]

Final Features B: [{{"label":xxx, "tag": yyy, "entity": zzz}}, {"label":xxx, "tag": yyy, "entity": zzz}}, ...]

Figure 11: Prompt of Consistency Optimization for Multi-Subject

Query Scenario Generation Prompt

You are a real user interacting with an AI chatbot. Based on ALL the selected entities provided, generate a diverse set of everyday conversation scenarios where you might engage with the chatbot.

Requirements:

1. Each generated scenario must be relevant to ALL the selected entities
2. The scenarios should be diverse and cover different domains such as:
 - Career planning and professional development
 - Medical consultation and healthcare management
 - Financial planning and advisory
 - Legal advice and compliance
 - Academic mentoring and education
 - Housing and real estate applications
 - Personal life management
3. Avoid repeating similar scenario types
4. Ensure each scenario naturally requires all selected entities to address user needs
5. Output the scenario names as a Python list

Selected Entities

{select_ents}

Please provide your output in the following format:

- Under "## Reason:", explain your analysis of the entities and how each proposed scenario naturally incorporates all of them
- Under "## Result:", output a JSON object containing the list of scenario names

Reason: [Explain your thought process and entity analysis]

Result: [{"scene": [scenario1, scenario2, scenario3]}]

Figure 12: Prompt for Query Scenario Generation

PII Detection Prompt

Please identify the PII entities and their corresponding PII types for each distinct individual mentioned in the conversation transcript, including both speakers and referenced individuals.

The PII types are defined as follows:

{pii_definition}

PII types include: ["PER", "CODE", "LOC", "ORG", "DEM", "DATETIME", "QUANTITY"]

Task Description:

Your task is to:

1. Identify ALL distinct individuals mentioned in the text, including:

- Primary speakers (marked with [PER_X])
- Individuals mentioned within others' statements
- Referenced colleagues, family members, or associates

2. For each identified individual, extract their associated PII entities, ensuring:

- Each entity is in its smallest viable text span
- Entities are correctly categorized by type
- Cross-referenced information is attributed to the correct individual

Important Rules:

1. Treat each individual as a separate subject, even if mentioned within another person's statement
2. Include both explicitly named individuals and those referenced through relationships
3. Maintain clear boundaries between different individuals' information
4. Extract exact entity spans without additional context
5. Preserve special characters in codes and quantities
6. Handle both direct mentions and indirect references

Given conversation transcript:

{user_desc}

Required Output Format:

For each identified individual (both speakers and mentioned persons), output:

Subject {{N}} {{ent1: type1, ent2: type2, ...}}

Example:

Input text: "[PER_1]: I'm Alex, working at Google. My friend Bob, who is 25 years old, works at Apple."

Expected output:

Subject {{1}} {"Alex": "PER", "Google": "ORG"}

Subject {{2}} {"Bob": "PER", "25 years": "DATETIME", "Apple": "ORG"}

Begin analysis now:

Figure 13: Prompt used for the PII Detection task

Basic Query-Related PII Detection Prompt

Please identify highly relevant PII (Personally Identifiable Information) entities from the background description PII entities that directly address or relate to the user's query.

Rules:

- Extract entities in their smallest possible span
- Exclude all person names
- Focus only on entities crucial for answering the query
- Return entities exactly as they appear in the text

Background description:

{desc}

Query:

{query}

Your output will contain the following format:

Answer: List the relevant PII entities, each enclosed in double quotes ("). Return only the list without explanation. Example: ["key_pii_1", ..., "key_pii_n"]

Please have your output follow the format below: (if there is only one entity, please output ["key_pii_1"]):

Answer: ["key_pii_1", ..., "key_pii_n"]

Figure 14: Prompt of Naive Method

Choice-Based Query-Related PII Detection Prompt

From the following options, Please identify highly relevant PII (Personally Identifiable Information) entities from the background description PII entities that directly address or relate to the user's query.

Rules:

- Extract entities in their smallest possible span
- Exclude all person names
- Focus only on entities crucial for answering the query
- Return entities exactly as they appear in the text
- Select only from the provided options

Background description:

{desc}

Query:

{query}

Options:

{choices}

Your output will contain the following format:

Answer: List the relevant PII entities, each enclosed in double quotes (""). Return only the list without explanation. Example: ["key_pii_1", ..., "key_pii_n"]

Please have your output follow the format below: (if there is only one entity, please output ["key_pii_1"]):

Answer: ["key_pii_1", ..., "key_pii_n"]

Figure 15: Prompt of Naive /w Choice Method

Chain-of-Thought Query-Related PII Detection Prompt

Please identify highly relevant PII (Personally Identifiable Information) entities from the background description PII entities that directly address or relate to the user's query.

Rules:

- Extract entities in their smallest possible span
- Exclude all person names
- Focus only on entities crucial for answering the query
- Return entities exactly as they appear in the text

Background description:

{desc}

Query:

{query}

Your output will contain the following format:

Thought: Explain your reasoning step by step for selecting the relevant PII entities.

Answer: List the relevant PII entities, each enclosed in double quotes (""). Return only the list without explanation. Example: ["key_pii_1", ..., "key_pii_n"]

Please have your output follow the format below: (if there is only one entity, please output ["key_pii_1"]):

Thought: xxx

Answer: ["key_pii_1", ..., "key_pii_n"]

Figure 16: Prompt of Self-CoT Method

Auto Chain-of-Thought Query-Related PII Detection Prompt with Examples

Please identify highly relevant PII (Personally Identifiable Information) entities from the background description PII entities that directly address or relate to the user's query.

Rules:

- Extract entities in their smallest possible span
- Exclude all person names
- Focus only on entities crucial for answering the query
- Return entities exactly as they appear in the text

Background description:

{desc}

Query:

{query}

You will be given 3 examples to help you understand the task.

Example 1:

Background: "Hello, I'm Sarah. I work at Microsoft as a junior developer with 2 years of experience. I live in Seattle."

Query: "What skills should I focus on developing in my early tech career at a leading software company to advance from my entry-level programming role?"

Answer: ["Microsoft", "junior developer"]

[Additional examples omitted for brevity]

Your output will contain the following format:

Thought: Explain your reasoning step by step for selecting the relevant PII entities.

Answer: List the relevant PII entities, each enclosed in double quotes ("). Return only the list without explanation. Example: ["key_pii_1", ..., "key_pii_n"]

Please have your output follow the format below: (if there is only one entity, please output ["key_pii_1"]):

Thought: xxx

Answer: ["key_pii_1", ..., "key_pii_n"]

Figure 17: Prompt of Auto-CoT Method

Self-Consistency Query-Related PII Detection Prompt

Please identify highly relevant PII (Personally Identifiable Information) entities from the background description PII entities that directly address or relate to the user's query.

Rules:

- Extract entities in their smallest possible span
- Exclude all person names
- Focus only on entities crucial for answering the query
- Return entities exactly as they appear in the text

Background description:

{desc}

Query:

{query}

Your output will contain the following format:

Thought: Generate 5 completely different perspectives of your reflections for selecting the relevant PII entities.

Summary: Output a summary of all your thinking.

Answer: List the relevant PII entities, each enclosed in double quotes (""). Return only the list without explanation. Example: ["key_pii_1", ..., "key_pii_n"]

Please have your output follow the format below: (if there is only one entity, please output ["key_pii_1"]):

Thought:

1. xxxxxx
2. xxxxxx
3. xxxxxx
4. xxxxxx
5. xxxxxx

Summary:

xxxxx

Answer: ["key_pii_1", ..., "key_pii_n"]

Figure 18: Prompt of Self-Consistency Method

Plan-and-Solve Query-Related PII Detection Prompt

Please identify highly relevant PII (Personally Identifiable Information) entities from the background description that directly address or relate to the user's query.

Rules:

- Extract entities in their smallest possible span
- Exclude all person names
- Focus only on entities crucial for answering the query
- Return entities exactly as they appear in the text

Background description:

{desc}

Query:

{query}

Your output will contain the following format:

Thought: Please start with a general plan for selecting the relevant PII entities, and then think step-by-step how to solve it based on the plan.

Answer: List the relevant PII entities, each enclosed in double quotes (""). Return only the list without explanation. Example: ["key_pii_1", ..., "key_pii_n"]

Please have your output follow the format below: (if there is only one entity, please output ["key_pii_1"]):

Thought: xxx

Answer: ["key_pii_1", ..., "key_pii_n"]

Figure 19: Prompt of Plan and Solve CoT Method

LLM Judge Evaluation Prompt

I want you to act as a neutral judge evaluating responses to a user query. Your task is to determine which responses better address the user's intent.

User Query: {query}

Reference Response (from prompt with all information):
{original_response}

Response A (from prompt with all PII masked):
{masked_response}

Response B (from prompt with query-unrelated PII masked):
{adaptive_response}

Please evaluate which responses better satisfy the user's intent and need, ignoring the presence of personally identifiable information (PII). Focus only on how well each response answers the query.

Rate each response on a scale of 1-10 where 10 is perfect:

1. Response A rating (1-10):
2. Response B rating (1-10):

Then provide a final judgment comparing each masked response to the reference response with one of these options:

For Response A:

- Reference is much better than Response A
- Reference is slightly better than Response A
- Reference and Response A are equally good
- Response A is slightly better than Reference
- Response A is much better than Reference

For Response B:

- Reference is much better than Response B
- Reference is slightly better than Response B
- Reference and Response B are equally good
- Response B is slightly better than Reference
- Response B is much better than Reference

Provide your final judgments as:

JUDGMENT A: [your choice]

JUDGMENT B: [your choice]

Figure 20: Prompt of LLM-as-Judge