

Learning to Keep Secrets: Empowering Local LLMs with Autonomous Dynamic Privacy-Conscious Delegation

Anonymous ACL submission

Abstract

While cloud-hosted Large Language Models (LLMs) offer superior capabilities, they introduce significant privacy risks. Conversely, local models ensure data sovereignty but often suffer from limited performance. Privacy-conscious delegation provides a solution by enabling local models to collaborate with remote models. In privacy-conscious delegation, a local model synthesizes a privacy-preserving prompt, queries a remote model, and combines the returned information with private context locally. However, existing methods rely on prompting to enforce privacy constraints, which can be brittle for local models, and they typically delegate in a static manner, causing unnecessary exposure when a local solution is sufficient. To address these limitations, we propose **ADAPT** (Autonomous Delegation Agent with Privacy Training), a framework that transforms delegation from a static pipeline into a dynamic, learnable agent. We train the local model to autonomously manage the entire execution path—from determining delegation triggers to synthesizing privacy-preserving prompts—thereby fostering robust, intrinsic privacy-preserving capabilities. Experimental results demonstrate that **ADAPT** significantly outperforms existing baselines, achieving a superior performance.

1 Introduction

The rapid advancement of LLMs has led to their widespread adoption and an increasing societal dependence on these systems. However, many LLMs are hosted by untrusted cloud providers (denoted as M_{remote}), posing significant privacy risks as users may transmit privacy-sensitive information within their prompts. To mitigate these risks, the task of privacy-conscious delegation has garnered attention. As shown in Figure 1, within this task, a local model (denoted as M_{local}) attempts to synthesize privacy-preserving prompts before querying M_{remote} , subsequently merging the remote re-

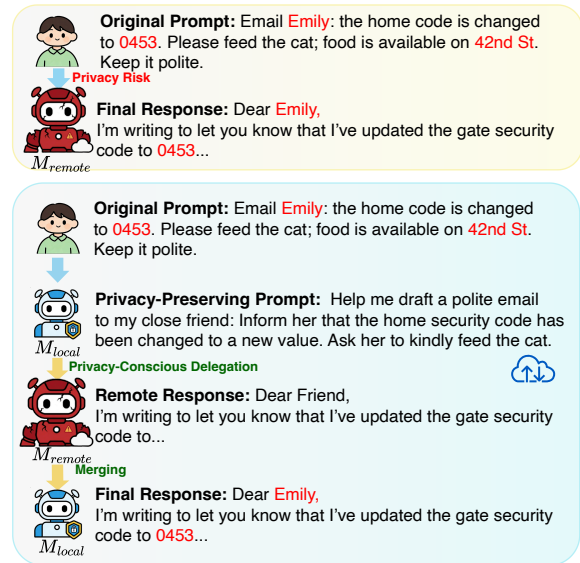


Figure 1: An illustration of privacy-conscious delegation in our research. Given a prompt that may contain privacy-sensitive information, M_{local} autonomously generates a privacy-preserving prompt for the remote model M_{remote} and integrates the returned remote response with the private context locally.

sponses with the local private context to generate the final response.

Prior research has made progress. Some approaches leverage in-context learning for privacy mitigation (Hartmann et al., 2024). CoGenesis utilizes high-level guidance from M_{remote} to assist M_{local} in handling private contexts for static tasks (Zhang et al., 2024). Similarly, PAPILON (Li et al., 2025) enables the M_{local} to synthesize privacy-preserving prompts to leverage the M_{remote} capabilities, merging the results locally. However, these methods overlook two critical limitations. First, relying on prompt engineering alone to enforce complex privacy constraints is often brittle for M_{local} . We believe privacy requirements are typically strict and nuanced, making them ill-suited to being reliably satisfied through prompting

without updating the parameters of M_{local} . Second, they adopt a static pipeline that delegates tasks to M_{remote} unconditionally. They overlook the inherent problem-solving potential of M_{local} and introduce avoidable privacy risks by involving M_{remote} even when M_{local} is sufficient to handle the prompt.

To overcome these, our research is built upon two core perspectives to bridge these gaps. First, we posit that privacy consciousness is often insufficiently established through prompting; instead, **privacy consciousness should be trained into the model’s parameters via training to ensure robust privacy-conscious behavior.** Second, we advocate for **the principle of data minimization: the most effective defense is non-transmission.** Unlike traditional static pipelines that incur unnecessary exposure, our approach empowers M_{local} to act as an autonomous agent for dynamic delegation, thereby significantly minimizing privacy leakage.

Overall, we propose **ADAPT (Autonomous Delegation Agent with Privacy Training)**, a framework that transforms the local model M_{local} from a passive sanitizer into an active, privacy-conscious agent. Unlike prior static approaches, ADAPT utilizes reinforcement learning to internalize privacy consciousness directly into the M_{local} . And we design a dynamic delegation mechanism where M_{local} learns to autonomously evaluate the necessity of external delegation. If an original prompt lies within the local model’s competence, ADAPT resolves it locally; otherwise, it generates a synthesized privacy-preserving prompt for M_{remote} and merges the remote results from M_{remote} with local private context to provide superior response quality. By optimizing a reward function that balances response quality against privacy leakage, ADAPT’s local model learns how to effectively collaborate with M_{remote} while preserving privacy. Our contributions are summarized as follows:

- ADAPT transforms the M_{local} into an autonomous, dynamic privacy-conscious agent, allowing it to decide when external delegation is needed.
- ADAPT uses reinforcement learning to internalize privacy consciousness in the M_{local} , optimizing privacy and performance.
- Experiments show that ADAPT improves the model’s capabilities and effectively completes the privacy-conscious delegation task.

2 Related Work

2.1 Training and Output Privacy

LLMs are known to verbatim memorize significant portions of their training data, which facilitates data extraction and membership inference attacks (Carlini et al., 2021; Nasr et al., 2023; Carlini et al., 2019; Lehman et al., 2021). This risk is particularly acute during task-specific finetuning on specialized or sensitive datasets, such as clinical or proprietary records (Miresghallah et al., 2022). Beyond memorization, privacy breaches frequently occur through LLM-generated outputs; models may inadvertently regurgitate sensitive context provided during In-Context Learning or Retrieval-Augmented Generation (Duan et al., 2024; Xia et al., 2025). To mitigate these risks, techniques such as machine unlearning (Ren et al., 2025; Liu et al., 2025) and differential privacy (Li et al., 2021) have been proposed to sanitize model weights, though they often incur non-trivial trade-offs in model utility.

2.2 Prompt-Based Inference Privacy

During deployment, privacy is threatened by both direct and latent information leakage. Static NER-based masking or scrubbing (Chen et al., 2023; Pilán et al., 2022) often fails to capture complex semantic dependencies, leading to significant utility degradation. Beyond simple sanitization, privacy-conscious delegation has emerged as a promising paradigm, wherein M_{local} manages privacy-sensitive data while remote large models M_{remote} provide advanced capabilities. One proposed mechanism involves the local synthesis of a similar but novel fake problem; the resulting response from the remote model for this synthetic task is then utilized as an in-context example to assist the local model in resolving the original private query (Hartmann et al., 2024). CoGenesis (Zhang et al., 2024) provides high-level guidance from remote LLMs to help local models process private contexts. PAILLON (Li et al., 2025) attempts to leverage the superior capabilities of M_{remote} by synthesizing privacy-preserving prompts, subsequently merging the remote results with private context locally. More recently, PrivacyPad (Hui et al., 2025) tries to forward only non-sensitive sentences.

Moving beyond static frameworks, ADAPT leverages reinforcement learning to internalize privacy-consciousness while empowering autonomous delegation, enabling M_{local} to dynamically balance privacy and utility.

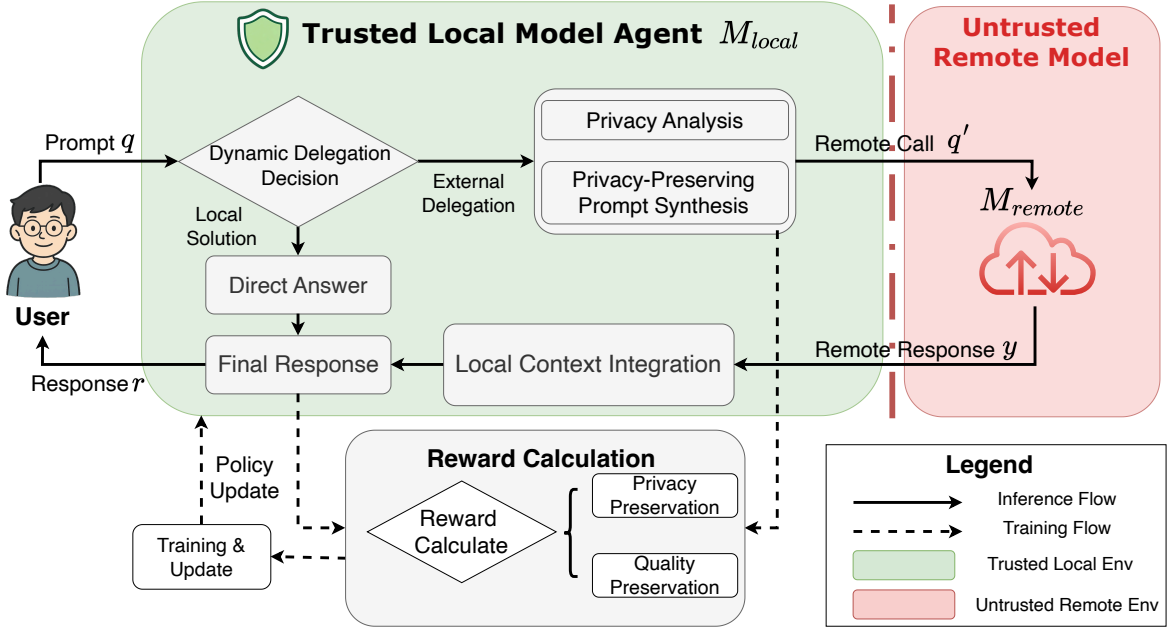


Figure 2: Overview of the ADAPT framework. M_{local} operates as an autonomous agent. **Solid arrows** indicate the inference flow: M_{local} adaptively chooses between local solution and external delegation using a privacy-preserving prompt, subsequently integrating the remote response to generate the final response r . **Dashed arrows** illustrate the training flow: the model is optimized via a reward mechanism that balances utility and privacy, adhering to the principle of data minimization.

3 ADAPT: Autonomous Delegation Agent with Privacy Training

3.1 Overview

As illustrated in Figure 2, ADAPT trains M_{local} to function as an autonomous agent. Given a user prompt q containing PII (Personally Identifiable Information) units $P = \{p_1, \dots, p_n\}$, the agent adaptively decides whether to resolve the prompt locally or delegate it to M_{remote} . In the case of delegation, it utilizes a privacy-preserving prompt q' and incorporates the remote response to produce the final response r .

Following prior work (Hui et al., 2025; Li et al., 2025), we quantify privacy leakage as the fraction of PII units exposed to M_{remote} .

3.2 Action Protocol and Delegation Decisions

To make privacy-conscious delegation behaviors learnable and controllable, we structure each inference instance into a lightweight two-stage protocol: (i) decision and privacy-preserving prompt synthesis, and (ii) external delegation and local integration. Under this protocol, M_{local} is always the model that produces the final user-facing response, while M_{remote} only provides information when delegation is selected.

Decision and Privacy-Preserving Prompt

Synthesis Given a user prompt q , M_{local} first generates a structured delegation block `<delegate_external>` to explicitly express its action choice. If M_{local} judges that external delegation is unnecessary, it outputs an empty `<delegate_external></delegate_external>` block. Otherwise, it first produces `<pii_analysis>` to identify PII units, and then synthesizes a privacy-preserving delegated prompt q' `<delegate_prompt>`.

External Delegation and Local Integration

If M_{local} chooses delegation, q' is sent to M_{remote} to obtain a remote response y . The system then provides this response to M_{local} as an observation wrapped in `<remote_LLM>`. If no delegation is selected, the system provides an empty `<remote_LLM></remote_LLM>` observation. Finally, M_{local} generates the final response r locally by integrating the original private context q with the observed remote response y , and encloses it within an `<answer>` tag.

We make delegation an explicit action interface: M_{local} decides whether to delegate by generating `<delegate_external>`, and, when delegating, takes responsibility for PII analysis and privacy-preserving prompt synthesis.

3.3 Reward Design

Privacy Preservation To quantify the extent of information exposure during delegation, we measure the leakage of PII units within the synthesized prompt q' . Given a set of PII units $P = \{p_1, \dots, p_n\}$ extracted from the original prompt q , the privacy leakage score is defined as the average exposure rate:

$$Leak(q', p_{1..n}) = \frac{1}{n} \sum_{j=1}^n \mathbb{I}(M_J(f_{PJ}, q', p_j)), \quad (1)$$

where f_{PJ} represents the prompt template for privacy judgment, M_J denotes an LLM-based judge, and $\mathbb{I}(\cdot)$ is an indicator function that converts the judge’s output into a binary value. Rather than relying on rigid string matching, which is often susceptible to circumvention, we utilize an LLM-based judge to capture semantic exposure. This approach has been validated in prior research as an effective means of identifying latent PII leakage (Li et al., 2025).

If M_{local} does not perform external delegation, the privacy leakage score is defined as zero.

Quality Preservation The quality of the generated response r is assessed by comparing it against a reference response r^* provided in the dataset. We compute the response quality score by asking an LLM judge whether r is at least as good as r^* . The response quality score is defined as:

$$Qual(r, r^*) = \mathbb{I}(M_J(f_{QJ}, r, r^*)), \quad (2)$$

where f_{QJ} is the evaluation prompt and $\mathbb{I}(\cdot)$ is a binary indicator. To mitigate positional bias, we perform two calls to the judge using permuted input pairs. We convert the judge’s free-text output into a binary score based on whether the response starts with “Yes.” If the judge provides identical answers for both permutations, indicating an inconsistent judgment, the responses are classified as equivalent in quality. A detailed description of the scoring algorithm is provided in Appendix B.

Final Reward We define our reward to encourage maximizing response quality score while minimizing privacy leakage score.

$$\mathcal{R} = Qual(r, r^*) - Leak(q', p_{1..n}), \quad (3)$$

Our M_{local} is designed to develop sophisticated strategies that enhance response quality while minimizing data exposure.

3.4 Model Training

We train M_{local} in two stages: a supervised fine-tuning phase for format alignment, followed by reinforcement learning to optimize for reward maximization.

Supervised Finetuning While some prior work suggests that supervised finetuning offers limited benefits for training agentic behaviors (Guo et al., 2025), we find it essential in our setting. Given that our study targets small models, cold-starting structured delegation behaviors is particularly challenging even with detailed prompts. We use GPT-4o to augment training dataset with 100 synthetic instructions and two solution-format demonstrations. This SFT phase serves primarily as format alignment, ensuring M_{local} reliably follows structured outputs to facilitate subsequent reinforcement learning. Empirically, although SFT alone is insufficient to solve the privacy-conscious delegation task, it provides a stable initialization for effective RL training.

Reinforcement Learning To prevent M_{local} from prematurely favoring non-delegation due to early-stage privacy score penalties, we introduce a frequency-aware reward rescaling mechanism inspired by Ada-GRPO (Wu et al., 2025). This approach encourages balanced exploration by dynamically adjusting the reward scale based on behavioral frequency. Specifically, the rescaled reward \mathcal{R}'_i for an output o_i at training step t is defined as:

$$\mathcal{R}'_i = \alpha_i(t) \cdot \mathcal{R}_i, \quad (4)$$

where \mathcal{R}_i is the raw reward, and the scaling factor $\alpha_i(t)$ utilizes a cosine-decayed inverse frequency:

$$\alpha_i(t) = \frac{B}{F(o_i)} \cdot d_i(t), \quad (5)$$

$$d_i(t) = 1 - \left(1 - \frac{F(o_i)}{B}\right) \cdot \frac{1 - \cos(\pi t/T)}{2}, \quad (6)$$

where $F(o_i)$ denotes the batch-wise frequency of the corresponding output pattern (delegation or non-delegation), B is the batch size, and T represents the total number of training steps. By dynamically scaling the reward, this mechanism encourages M_{local} to extensively explore both delegation paths during the early stages of training. Specifically, the scaling factor $\alpha_i(t)$ smoothly decays from the inverse frequency $\frac{B}{F(o_i)}$ at $t = 0$ to unity at $t = T$. This strategy effectively prioritizes rare behaviors

initially before gradually converging toward the original reward distribution to ensure policy stability. For trajectory optimization, we utilize Group Relative Policy Optimization (GRPO) to optimize the model based on the rescaled reward \mathcal{R}' .

4 Experimental Setup

4.1 Datasets

PUPA Dataset To compare ADAPT with prior studies, we conduct training and evaluation on the PUPA dataset. PUPA is constructed from real-world user-LLM interactions and contains naturalistic PII across a variety of general-domain topics. The dataset is annotated with potential PII units and their corresponding ground truth labels. It is divided into two primary subsets for training and evaluation: PUPA-New and PUPA-TNB. For our experiments, we utilize the official splits and statistics as reported by the original authors.

PsyPriv Dataset To further benchmark our approach in Out-of-Distribution(OOD) sensitive scenarios, we develop PsyPriv, a synthetic privacy-sensitive dataset derived from the anonymized psychological QA corpus PSYQA (Sun et al., 2021). Given that mental health prompts are inherently privacy-sensitive, we synthesize PII units through a two-step pipeline: (1) **PII Injection**: We prompt GPT-4o to rewrite each prompt by injecting diverse PII elements and annotating the corresponding spans; (2) **Consistent Regeneration**: Given the PII-injected prompt, we utilize GPT-5.1 to regenerate a response that maintains consistency with both the original psychological counseling style and the newly injected PII. For more details, please refer to Appendix D.

To better evaluate the generalization of our model, PsyPriv is reserved exclusively for OOD testing. This setup serves as supplementary evidence to demonstrate the robustness and OOD capabilities of ADAPT across specialized domains.

4.2 Models

Remote Model Unless otherwise specified, we employ a fixed remote model, GPT-4o-mini, as M_{remote} across all methods to ensure a controlled comparison (OpenAI, 2024). Additionally, GPT-4o-mini is also adopted as our primary evaluation model, leveraging its high cost-efficiency while maintaining consistent assessment performance.

Metrics	PUPA	PsyPriv
# Instances	664 (train) / 237 (test)	200 (test)
Domain	General	Mental
Avg.# PII	2.9	5.3
Avg.P Len.	1352.0	711.6
Avg.R Len.	1553.7	1609.8

Table 1: Comparison of dataset statistics between PUPA and PsyPriv. The average PII count and prompt/response lengths are reported to illustrate the distribution shift between standard user-LLM interactions and the high-sensitivity scenarios used for ADAPT’s robustness evaluation.

Local Models We evaluate our framework by systematically varying M_{local} across several on-device models: Qwen2.5-0.5B-Instruct, Llama3.2-1B-Instruct, Qwen2.5-1.5B-Instruct, and Llama3.2-3B-Instruct (Qwen et al., 2025; Grattafiori et al., 2024). Our investigation specifically focuses on small-scale on-device models, which represent the most practical deployment scenarios while frequently necessitating external delegation.

4.3 Baselines

To evaluate the performance of ADAPT, we compare it against several representative baselines ranging from zero-delegation local execution to various supervised delegation strategies:

- **Always-Local**: Answer all queries using only M_{local} (no delegation), providing an upper bound on privacy but typically lower utility.
- **PAPILLON (Li et al., 2025)**: A static rewrite-and-delegate pipeline. M_{local} rewrites the prompt into a privacy-preserving version, queries M_{remote} , and integrates the remote response locally.
- **SFT-D**: We use GPT-4o to generate synthetic delegation trajectories (e.g., PII analysis and delegated prompts) and finetune M_{local} via SFT to imitate them.
- **SFT-R**: We score M_{local} on the training set; low-quality cases are supervised with GPT-4o-generated delegation trajectories, while the rest use local-only answers as targets, teaching M_{local} how to route delegation.

We reimplemented PAPILLON and evaluated it across all settings. For more details, please refer to Appendix A.

Method / System	PUPA-TNB			PsyPriv		
	Qual. (%) \uparrow	Leak. (%) \downarrow	Dele.(%)	Qual. (%) \uparrow	Leak. (%) \downarrow	Dele.(%)
GPT-4o-mini [Unredacted]	88.2	100.0	100.0	69.0	100.0	100.0
GPT-4o-mini [Redacted]	77.2	0.0	100.0	62.0	0.0	100.0
<i>M_{local}</i> : Llama3.2-1B-Instruct						
Always-Local	40.4	0.0	0.0	36.0	0.0	0.0
PAPILLON	58.0	39.3	100.0	33.0	41.9	100.0
SFT-D	38.4	9.7	100.0	16.5	10.5	100.0
SFT-R	47.6	3.8	47.6	28.0	4.7	41.0
ADAPT	80.7	2.0	68.7	50.0	4.0	90.0
<i>M_{local}</i> : Llama3.2-3B-Instruct						
Always-Local	57.3	0.0	0.0	44.5	0.0	0.0
PAPILLON	60.8	30.8	100.0	31.5	5.3	100.0
SFT-D	63.8	8.5	100.0	47.0	15.9	100.0
SFT-R	67.6	4.4	42.3	38.0	7.3	39.0
ADAPT	92.3	0.4	28.8	68.0	1.2	53.5
<i>M_{local}</i> : Qwen2.5-0.5B-Instruct						
Always-Local	31.3	0.0	0.0	19.0	0.0	0.0
PAPILLON	38.6	54.5	100.0	23.5	30.6	100.0
SFT-D	42.6	16.1	100.0	9.5	11.8	100.0
SFT-R	47.1	13.5	87.6	5.5	11.2	96.0
ADAPT	56.6	5.1	97.3	40.0	7.3	100.0
<i>M_{local}</i> : Qwen2.5-1.5B-Instruct						
Always-Local	42.4	0.0	0.0	25.0	0.0	0.0
PAPILLON	65.4	67.7	100.0	45.0	27.1	100.0
SFT-D	61.6	10.4	100.0	37.5	14.4	100.0
SFT-R	70.6	6.5	50.2	40.0	7.8	56.0
ADAPT	82.4	2.4	55.1	55.5	4.4	60.0

Table 2: Main results on PUPA-TNB and PsyPriv datasets. We compare ADAPT against various baselines across different local model scales. **Qual.** denotes response quality preservation, **Leak.** represents the privacy leakage rate (PII exposure), and **Dele.** indicates the percentage of synthesized prompts delegated to the M_{remote} . ADAPT consistently achieves a superior performance, maintaining high utility while significantly reducing privacy leakage compared to static delegation methods.

5 Main Results

As illustrated in Table 2, ADAPT not only achieves improvements in response quality but also maintains privacy leakage at an exceptionally low level. This demonstrates that direct parameter updates enable M_{local} to internalize more effective privacy-preserving strategies while simultaneously generating high-quality responses. Consequently, ADAPT leverages the complementary strengths of both M_{remote} and M_{local} for generation while ensuring robust data security. In contrast, the baseline relying on synthesized SFT trajectories underperforms in both response quality score and privacy leakage score.

Beyond overall performance gains, a critical observation is the inverse relationship between the capacity of M_{local} and its delegation frequency. As the model scale increases, the delegation rate on the PUPA-TNB dataset drops significantly from 68.7% to 28.8%. While one might argue that these gains stem from general capability enhancement

via RL rather than explicit privacy mechanisms, we contend that such capabilities are intrinsic to the privacy objective. **Fundamentally, when our M_{local} elects to resolve a query locally and provides a high-quality response, this act in itself constitutes the most effective form of privacy preservation.**

The results on PsyPriv further underscore the out-of-distribution robustness of ADAPT. Mental health prompts are inherently more complex and sensitive than general-domain queries. Despite not being trained on similar datasets, ADAPT achieves notable gains on PsyPriv, indicating that the model has not merely overfitted to PUPA-TNB but has learned a generalized cooperative policy. Finally, comparisons to supervised baselines underscore the necessity of reinforcement learning for end-to-end privacy-conscious delegation. SFT-D can align the output format and pipeline but often fails to improve utility, suggesting that imitation alone is insufficient to learn stable privacy-utility strategies.

SFT-R learns partial routing behavior but remains consistently behind ADAPT in both quality and leakage.

6 Analysis

6.1 Delegation vs. Non-delegation

To evaluate the necessity of dynamic delegation in ADAPT, we conduct a study by forcing two extreme behaviors:

- **ADAPT-Local:** Never delegates, generating responses locally.
- **ADAPT-Remote:** Always delegates, sending prompts to M_{remote} .

Method	Qual. (%) \uparrow	Leak. (%) \downarrow	Dele. (%)
PUPA-TNB			
ADAPT	92.3	0.4	28.8
w/ Local	84.6	0.0	0.0
w/ Remote	85.7	4.6	100.0
PsyPriv			
ADAPT	68.0	1.2	53.5
w/ Local	60.5	0.0	0.0
w/ Remote	65.5	5.1	100.0

Table 3: Controlled delegation via control prefixes on Llama3.2-3B. We additionally introduce two prefixes to control model behavior: `<delegate_external></delegate_external>` and `<delegate_external><pii_analysis>`.

As shown in Table 3, ADAPT achieves the best performance. On PUPA-TNB, ADAPT improves quality over ADAPT-Local while delegating fewer queries than ADAPT-Remote. **In contrast, ADAPT-Remote has higher leakage and does not improve quality, showing that blind delegation is less private and effective.**

6.2 Robustness to Remote Model Backends

To verify ADAPT’s adaptability, we evaluate its performance by replacing M_{remote} with GPT-4.1 and Claude 3.5 Haiku. As shown in Figure 3, ADAPT consistently maintains low privacy leakage and competitive utility across different backends. The overall utility scales with the capability of the remote model, confirming that the learned delegation policy is robust and effectively leverages stronger cloud assistants without specialized re-tuning.

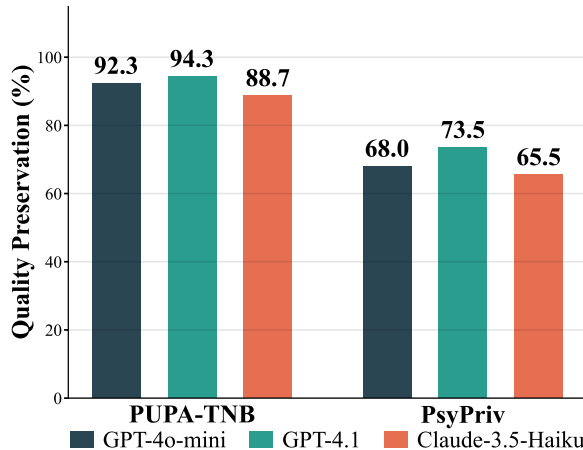


Figure 3: Robustness of ADAPT across different M_{remote} . We evaluate the Qual.(%) using GPT-4.1 and Claude 3.5 Haiku as alternative cloud assistants.

6.3 Reward Ablation

To investigate the contribution of individual components within our reinforcement learning framework, we conduct an ablation study by selectively removing the Privacy Reward (PR) and the frequency-aware Reward Scaling (RS).

Method	Qual. (%) \uparrow	Leak. (%) \downarrow	Dele. (%)
PUPA-TNB			
ADAPT	92.3	0.4	28.8
w/o PR	93.5	20.6	35.2
w/o RS	83.1	0.0	0.0
PsyPriv			
ADAPT	68.0	1.2	53.5
w/o PR	70.5	15.3	62.0
w/o RS	63.5	0.0	0.0

Table 4: Ablation study on reward components using Llama3.2-3B-Instruct. “w/o Privacy Reward” refers to optimizing only for quality preservation. “w/o Reward Scaling” denotes training without the adaptive frequency-aware coefficient.

As shown in Table 4, we observe that omitting the PR leads to a marginal improvement in Qual., but results in a severe surge in privacy leakage. Conversely, without the RS mechanism, the model prematurely converges to a non-delegation strategy in the early stages of training, failing to explore the potential benefits of cloud-based resolution. These results underscore the necessity of our proposed reward design in facilitating effective, privacy-conscious autonomous delegation.

Status	Component	Content
Success (Local)	Original Prompt	Rachel Zheng is a clerk in the Administrative Department of H&R Technology ... Please translate the above sentences into Chinese.
	Decision	Local Resolution
	Remote Response Final Response	N/A 郑瑞秋是 H&R科技公司 行政部的职员...
Success (Delegate)	Original Prompt	I have a business named Shumti and Co ... can you give me a business plan that could increase our sales?
	Privacy-preserving Prompt Remote Response	... Provide a business plan for a retail opticians shop with bespoke frames to increase sales. ... Marketing Strategy: Use social media to showcase bespoke frames...Customer Loyalty: Reward repeat optical customers... [General Advice]
	Final Response	...Marketing Strategy for Shumti and Co : Use social media to showcase our bespoke frames... Customer Loyalty: Reward Shumti and Co 's repeat customers... [Context Injected]
Failure (Leakage)	Original Prompt	Position located in the HQCSS, Division of Administration, New York ... Write a cover letter.
	Privacy-preserving Prompt Remote Response	... Apply for an internship at a leading international organization's headquarters in New York I am excited to apply for the internship at your New York headquarters... I believe my skills align with the New York team's goals...
	Final Response	... I am excited to apply for the internship at the HQCSS in New York ... I believe my skills align with the Division of Administration 's goals...

Table 5: Case study examples of ADAPT on privacy-conscious delegation. We show two successful cases (local resolution and external delegation) and one failure case where residual sensitive information is leaked. For each instance, we present the original prompt, the synthesized privacy-preserving prompt (if delegated), the remote response, and the final locally generated answer. Red text highlights PII units.

7 Case Study

Table 5 illustrates some behaviors and failure modes of ADAPT under privacy-conscious delegation. ADAPT selects a local solution: the model completes the prompt without external delegation, keeping PII, such as names and organizations, entirely on-device. This behavior directly instantiates the principle that non-transmission yields zero leakage. In contrast, for open-ended requests that require higher generation capabilities, ADAPT proactively delegates the task while synthesizing an abstracted prompt. This prompt preserves the task intent but removes identifying details, such as by rewriting a company-specific request into a generic description of a retail opticians shop. M_{remote} then provides broadly applicable strategies, which M_{local} grounds back into the private context during the final response, recovering response quality without exposing the original identifiers. Our analysis of the delegated prompts reveals that M_{local} has effectively learned sophisticated strategies for constructing privacy-preserving prompts; please refer to Appendix G for a detailed discussion.

In the failure example, the original prompt intends to draft a cover letter for a position at “HQCSS, Division of Administration located in New York.” Although M_{local} successfully identi-

fies and masks the specific organizational names “HQCSS” and “Division of Administration,” it retains the specific geographic information “New York” when synthesizing the privacy-preserving prompt. Even though M_{remote} returns general advice based on the broad context of “New York,” this transmission still constitutes a privacy breach in PUPA-TNB. We believe the primary reason is that the dense presence of PII causes the model to struggle with multi-entity de-identification, leading it to prioritize preserving geographic context to ensure that the M_{remote} provides relevant results.

8 Conclusion

We present ADAPT, a framework that redefines privacy-conscious delegation by evolving M_{local} into an autonomous agent. ADAPT enables M_{local} to dynamically determine delegation necessity and synthesize privacy-preserving prompts via reinforcement learning. Our findings underscore the principle of data minimization: as local model capacity scales, M_{local} increasingly resolves queries locally, effectively eliminating external exposure. This indicates that enhancing local agentic intelligence serves as the privacy defense. We conclude that ADAPT provides a path toward ensuring data sovereignty in collaborative LLM ecosystems.

524 Limitations

525 We identify several limitations in the current
526 ADAPT framework that suggest directions for fu-
527 ture research.

- 528 • **Boundary PII Identification:** While effec-
529 tive for common identifiers, the model may oc-
530 casionally overlook context-specific PII (e.g.,
531 niche technical terms or landmarks). Strength-
532 ening the internalized policy for diverse do-
533 mains remains a key challenge.
- 534 • **Single-Turn Scope:** Our evaluation focuses
535 on independent queries. In multi-turn dia-
536 logues, privacy risks can accumulate across
537 history. Extending the autonomous delegation
538 policy to maintain long-term privacy states is
539 a necessary next step.
- 540 • **Judge-Model Bias:** The RL phase relies on
541 an LLM-based reward signal. Although more
542 semantic than string matching, this approach
543 is subject to the inherent biases or sensitivity
544 limits of the judge model used during training.
- 545 • **Inference Overhead:** The agentic reasoning
546 process (PII analysis and prompt synthesis)
547 introduces extra generation steps. Optimizing
548 the efficiency of this multi-step pipeline is es-
549 sential for highly latency-sensitive on-device
550 applications.

551 These findings represent a proof-of-concept
552 for autonomous privacy-conscious delegation.
553 ADAPT should be viewed as an early-stage inves-
554 tigative instrument rather than a standalone security
555 guarantee. Future work will focus on integrating
556 formal privacy guarantees, such as differential pri-
557 vacy, and validating the framework across broader,
558 multi-turn interaction datasets to ensure robust data
559 sovereignty in production environments.

560 References

561 Nicholas Carlini, Chang Liu, Úlfar Erlingsson, Jernej
562 Kos, and Dawn Song. 2019. The secret sharer: Eval-
563 uating and testing unintended memorization in neu-
564 ral networks. In *28th USENIX security symposium*
565 (*USENIX security 19*), pages 267–284.

566 Nicholas Carlini, Florian Tramer, Eric Wallace,
567 Matthew Jagielski, Ariel Herbert-Voss, Katherine
568 Lee, Adam Roberts, Tom Brown, Dawn Song, Ulfar
569 Erlingsson, and 1 others. 2021. Extracting training
570 data from large language models. In *30th USENIX*

security symposium (USENIX Security 21), pages
2633–2650. 571 572

Yu Chen, Tingxin Li, Huiming Liu, and Yang Yu. 573
2023. Hide and seek (has): A lightweight frame- 574
work for prompt privacy protection. *arXiv preprint* 575
arXiv:2309.03057. 576

Haonan Duan, Adam Dziedzic, Mohammad Yaghini, 577
Nicolas Papernot, and Franziska Boenisch. 2024. On 578
the privacy risk of in-context learning. *arXiv preprint* 579
arXiv:2411.10512. 580

Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, 581
Abhinav Pandey, Abhishek Kadian, Ahmad Al- 582
Dahle, Aiesha Letman, Akhil Mathur, Alan Schel- 583
ten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh 584
Goyal, Anthony Hartshorn, Aobo Yang, Archi Mi- 585
tra, Archie Sravankumar, Artem Korenev, Arthur 586
Hinsvark, and 542 others. 2024. *The llama 3 herd of* 587
models. *Preprint*, arXiv:2407.21783. 588

Daya Guo, Dejian Yang, Haowei Zhang, Junxiao 589
Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shi- 590
rong Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025. 591
Deepseek-r1: Incentivizing reasoning capability in 592
llms via reinforcement learning. *arXiv preprint* 593
arXiv:2501.12948. 594

Florian Hartmann, Duc-Hieu Tran, Peter Kairouz, Vic- 595
tor Cărbune, and Blaise Aguera Y Arcas. 2024. Can 596
llms get help from other llms without revealing pri- 597
vate information? In *Proceedings of the Fifth Work-* 598
shop on Privacy in Natural Language Processing, 599
pages 107–122. 600

Zheng Hui, Yijiang River Dong, Sanhanat Sivapirom- 601
rat, Ehsan Shareghi, and Nigel Collier. 2025. Pri- 602
vacypad: A reinforcement learning framework for 603
dynamic privacy-aware delegation. *arXiv preprint* 604
arXiv:2510.16054. 605

Eric Lehman, Sarthak Jain, Karl Pichotta, Yoav Gold- 606
berg, and Byron C Wallace. 2021. Does bert pre- 607
trained on clinical notes reveal sensitive data? In 608
Proceedings of the 2021 Conference of the North 609
American Chapter of the Association for Computa- 610
tional Linguistics: Human Language Technologies, 611
pages 946–959. 612

Siyan Li, Vethavikashini Chithra Raghuram, Omar 613
Khattab, Julia Hirschberg, and Zhou Yu. 2025. Pa- 614
pillon: Privacy preservation from internet-based and 615
local language model ensembles. In *Proceedings of* 616
the 2025 Conference of the Nations of the Americas 617
Chapter of the Association for Computational Lin- 618
guistics: Human Language Technologies (Volume 1: 619
Long Papers), pages 3371–3390. 620

Xuechen Li, Florian Tramer, Percy Liang, and Tatsunori 621
Hashimoto. 2021. Large language models can be 622
strong differentially private learners. In *International* 623
Conference on Learning Representations. 624

Zheyuan Liu, Guangyao Dou, Mengzhao Jia, Zhaoxuan 625
Tan, Qingkai Zeng, Yongle Yuan, and Meng Jiang. 626

Hyperparameter	Value
Learning Rate	1×10^{-6}
LR Warmup Steps Ratio	0.285
KL Loss Coefficient (β)	0.001
Mini-batch Size	16
Micro-batch Size	4
Advantage Estimator	GRPO
Rollout Temperature	1.0
Max Sequence / Response Length	8192 / 2048
Total Training Epochs	5
Optimization Offload	Enabled (FSDP)

Table 6: Hyperparameters for ADAPT RL Training (based on the verl framework).

Algorithm 1 Quality Preservation Metric

```

1: Input: User Query  $q$ 
2: Input: Pipeline Output  $C_L$ 
3: Input: Target Output  $C_T$ 
4: Input: LLM Judge,  $M_J$ 
5: Input: Prompt for Quality Judgment  $f_{QJ}$ 
6: Output:  $QUAL(C_L, C_T)$ 
7:  $s_1 \leftarrow \text{int}(f_{QJ}(M^J, q, C_L, C_T))$ 
8:  $s_2 \leftarrow \text{int}(f_{QJ}(M^J, q, C_T, C_L))$ 
9:    $\triangleright$  Reverse the order of candidate responses
10: if  $s_1 = s_2$  then
11:   return 1  $\triangleright$  The judgment is inconsistent
12: else if  $s_1 \neq s_2$  then
13:   return  $s_1$   $\triangleright$  The judgment is consistent
14: end if

```

729
730 You are comparing the quality of two
731 responses, given a user query. Is
732 response A at least as good as response
733 B? Respond with yes or no.

734 –
735 Follow the following format.

736 User Query: The user’s request to be
737 fulfilled.

738 Response A: $\{\text{response_A}\}$

739 Response B: $\{\text{response_B}\}$

740 Reasoning: Let’s think step by step in
741 order to $\{\text{reasoning}\}$

742 Output: $\{\text{output}\}$

743 –
744 User Query: $\{\text{user query}\}$

745 Response A: $\{\text{response A}\}$

746 Response B: $\{\text{response B}\}$

747 Reasoning: Let’s think
748 step by step in order to
749

B.2 Privacy Preservation

750 Ensuring the security of sensitive data is a funda-
751 mental objective of the ADAPT framework. We
752 define privacy preservation as the successful redac-
753 tion of PII from the original user prompt before it
754 is delegated to external large language models. To
755 rigorously evaluate the effectiveness of our privacy-
756 preserving mechanism, we utilize a **Leakage Judg-
757 ment Metric**. This metric employs an LLM-based
758 judge to detect whether specific PII units are still
759 present in the generated delegated prompts. The
760 specific evaluation prompt designed for this leak-
761 age detection task is provided below:
762

763 Determine whether a piece of
764 information is present in the prompt
765 given. Respond with yes or no.
766 –

767 Follow the following format.

768 Information Piece: The piece of
769 information that we wish to check.
770 If this piece of information is a
771 celebrity name, respond with no.
772

773 Prompt: $\{\text{prompt}\}$

774 Output: $\{\text{output}\}$

775 –
776 Information Piece: $\{\text{PII unit}\}$

777 Prompt: $\{\text{proposed query}\}$

778 Output:
779

780 **Reliability of the Judge** For evaluation, we sam-
781 ple 50 pairs of candidate responses for queries
782 from PsyPriv. Each pair consists of (1) the out-
783 put from Llama3.2-3B-Instruct and (2) the original
784 GPT-5.1 response. Within this sample, 25 pairs
785 satisfy $QUAL(A, B) = 1$, while the remaining 25
786 satisfy $QUAL(A, B) = 0$. This setup applies to
787 both the current and subsequent studies.

788 Participants are tasked with selecting the supe-
789 rior response between the two candidates or mark-
790 ing them as tied in quality. Two participants inde-
791 pendently label each pair. When $QUAL(A, B) =$
792 0, the alignment rate is 74%, indicating that partici-
793 pants prefer response B in 74% of instances. For
794 pairs where $QUAL(A, B) = 1$, the alignment rate
795 reaches 92%, where alignment is defined as majori-
796 ty agreement for either response A or a tie. These
797 results establish the general validity of our quality
798 metric, supported by a high inter-rater consistency
799 with a Cohen’s Kappa of 0.87. Regarding the vali-
800 dation of the PUPA dataset, it is already established

801 in prior research.

802 C Synthetic Trajectory Generation for 803 Privacy-Preserving Delegation

804 To bootstrap structured delegation behaviors for
805 small on-device models, we additionally construct
806 a small set of synthetic training trajectories using
807 a controlled “AI simulator” setup. The goal is to
808 generate supervision that matches our action pro-
809 tocol (delegation decision, PII analysis, privacy-
810 preserving prompt synthesis, and final local inte-
811 gration), while explicitly separating (i) what can be
812 safely delegated from (ii) what must remain local.

813 C.1 Synthetic Trajectory Generation Prompt 814 Template

816 You are an AI simulator simulating
817 a ‘Privacy-Preserving Collaborative
818 Generation’ architecture.

819 Simulation Rules:

- 820 1. Delegate/External Model: Detects PII
821 (Personal Identifiable Information) and
822 rewrites the prompt to be safe/generic.
- 823 2. Strong LLM: Generates a template
824 response based on the generic prompt.
- 825 3. Final Answer: Combines the template
826 with original details to form the final
827 customized answer.
828

829 Output Format:

830 You must strictly output the response
831 in the following XML format (do not use
832 markdown code blocks, just raw text):

```
833 <delegate_external>  
834 <pii_analysis>  
835 Detected PII: ...  
836 Risk Assessment: ...  
837 </pii_analysis>  
838 <delegate_prompt>  
839 ...  
840 </delegate_prompt>  
841 </delegate_external>  
842 <remote_LLM>  
843 ...  
844 </remote_LLM>  
845 <answer>  
846 ...  
847 </answer>
```

850 D PsyPriv Construction 851

852 The PSYPRIV dataset is specifically designed
853 to evaluate ADAPT’s robustness in high-stakes,
854 domain-specific scenarios (Mental Health).

855 D.1 Construction Pipeline 856

857 We utilized a two-stage automated pipeline assisted
858 by GPT-4o and GPT-5.1:

859 PII Injection Randomly selected 200 samples
860 from the PSYQA corpus. We injected five types of
861 PII: Name, Organization, Specific Location, Con-
862 tact Info, and ID numbers.

863 You are a professional data annotator
864 specializing in creating realistic
865 psychological counseling dialogue
866 records for research purposes. Your
867 task is to take a segment of anonymized
868 psychological counseling dialogue turns
869 and enrich its content by injecting
870 diverse and creative Personally
871 Identifiable Information (PII).

872 Instructions: 1. Carefully read the
873 anonymized input text to fully understand
874 the context.

875 2. Inject multiple types of Personally
876 Identifiable Information (PII) from the
877 complete list below. Use creativity to
878 make the details concrete and vivid.

879 3. Ensure the dialogue remains fluent,
880 natural, and contextually coherent. The
881 injected PII should feel like a natural
882 part of the dialogue or record.

883 4. Critically important: Do NOT modify
884 the original texts. You may only add
885 personal, consultation-related, and
886 contextual details.

887 5. All injected PII must be logically
888 consistent within the text. 6. Separately
889 output the injected PII using double
890 vertical bars (||) as separators in a
891 field named “pii”. The PII listed there
892 must be exactly identical to the PII
893 inserted into the original text. Do NOT
894 use double vertical bars in the original
895 text itself.

896 7. The output format has already been
897 provided. Ensure that the output is in
898 JSON format and that you output strictly

899 valid JSON.

900 8. Do not use Chinese. Your output must
901 be in English only.

902 Types of Personally Identifiable
903 Information (PII) to Inject (use
904 creativity and specificity): Personal
905 Identifiers: - Full name (e.g., Liu
906 Chenyang) - Date of birth (e.g., January
907 24, 2006) - Workplace / School attended
908 (e.g., Xihu Middle School) - Student
909 ID, employee ID, or case number (e.g.,
910 J, H-48920-C) - City / Hometown (e.g.,
911 Chengdu, Sichuan Province) - Phone number
912 (e.g., 215-555-0182) - Email address
913 (e.g., JH-48920-C@emailservice.com) The
914 total number of injected PII items must
915 not exceed six.

917 **Response Realignment** To ensure the reference
918 response r^* remains consistent with the injected
919 PII, we re-generated the ground-truth answers
920 to reflect the persona established in the prompt.

922 You are an experienced psychological
923 counselor. Please answer the following
924 question in the tone of a counselor.
925 Do not mechanically list your points;
926 instead, empathize with the user and
927 smoothly weave your suggestions into a
928 coherent, flowing response.

930 **PII Distribution** The average number of PII
931 units per query is 5.26. The distribution of PII
932 types is as follows: Name (28%), Location (22%),
933 Contact (20%), Organization (18%), and Others
934 (12%).

935 E Why can ADAPT (Llama3.2-3B) 936 surpass GPT-4o-mini in Qual.?

937 This result should **not** be interpreted as a 3B local
938 model being inherently stronger than GPT-4o-mini.
939 **Our Qual. metric evaluates the final system out-
940 put r against the reference r^* via an LLM judge
941 and ADAPT is a *system-level* method rather than
942 a single-pass generator.** When delegation is trig-
943 gered, ADAPT queries the same remote M_{remote}
944 but with a learned, privacy-preserving prompt q'
945 that removes noise, restructures the task, and bet-
946 ter aligns the request with the remote model's
947 strengths. The M_{local} then performs post-editing

948 and private-context reinjection to produce the final
949 response. Consequently, ADAPT can be viewed
950 as combining (i) routing, (ii) prompt optimization
951 for the remote model, and (iii) local refinement,
952 which may yield higher judge-preferred outputs
953 than directly prompting GPT-4o-mini once with the
954 raw query. To mitigate concerns about evaluator-
955 specific artifacts, we further report cross-judge re-
956 sults (J1/J2) in Appendix F, showing that ADAPT's
957 gains remain consistent across different judge back-
958 ends.

959 F Cross-Judge Robustness

960 A prevalent concern in reinforcement learning (RL)
961 frameworks is the risk that M_{local} may overfit to
962 the idiosyncratic preferences or evaluative biases
963 of a specific judge model. To ensure that ADAPT's
964 performance gains reflect genuine improvements in
965 privacy and utility rather than artifacts of a single
966 evaluator, we conduct a cross-judge robustness test.

967 F.1 Evaluation Protocol

968 We fix the outputs generated by ADAPT and its
969 baselines and re-evaluate them using two indepen-
970 dent judge backends:

- 971 • **J0 (Training Judge):** GPT-4o-mini with our
972 default evaluation templates, serving as the
973 training-distribution reference.
- 974 • **J1 (Stronger Judge):** GPT-4o, utilized to test
975 whether more sophisticated reasoning capabil-
976 ities reveal latent privacy leaks or utility flaws
977 that might be overlooked by a smaller model.
- 978 • **J2 (Cross-Family Judge):** Claude-3.5-Haiku,
979 employed to mitigate potential "provider bias"
980 and verify if ADAPT's policy generalizes to
981 a completely different model architecture and
982 safety alignment philosophy.

983 To ensure that the performance gains of ADAPT
984 stem from a robust internal policy rather than
985 merely exploiting the specific evaluative tenden-
986 cies of the training reward model, we conducted a
987 cross-judge validation using a more capable fron-
988 tier model. **While we acknowledge that a portion
989 of the improvements may reflect alignment with
990 the specific preferences of the training evalua-
991 tor, this analysis reveals that ADAPT maintains
992 its superiority in both response quality and pri-
993 vacy preservation even under the scrutiny of a
994 more stringent and critical evaluator.** Whereas

Method	Qual. (%) \uparrow			Leak (%) \downarrow		
	J0	J1	J2	J0	J1	J2
PUPA-TNB (M_{local} : Llama3.2-3B-Instruct)						
GPT-4o-mini [Unredacted]	88.2	82.9	75.5	100.0	100.0	100.0
GPT-4o-mini [Redacted]	77.2	64.8	59.0	0.0	0.0	0.0
Always-Local	57.3	52.0	43.5	0.0	0.0	0.0
PAPILLON	60.8	51.5	44.0	39.3	30.5	33.8
SFT-D	63.8	56.1	48.5	8.5	13.2	10.4
SFT-R	67.6	58.5	51.5	4.4	4.2	5.1
ADAPT	92.3	79.9	70.0	0.4	0.3	0.5
PSYPRIV (M_{local} : Llama3.2-3B-Instruct)						
GPT-4o-mini [Unredacted]	69.0	50.5	45.5	100.0	100.0	100.0
GPT-4o-mini [Redacted]	62.0	44.0	38.5	0.0	0.0	0.0
Always-Local	44.5	43.0	32.0	0.0	0.0	0.0
PAPILLON	31.5	26.5	23.5	5.3	4.1	4.7
SFT-D	47.0	35.5	29.5	15.9	13.9	14.5
SFT-R	38.0	24.5	21.0	7.3	9.2	8.4
ADAPT	68.0	52.5	42.5	1.2	1.7	1.4

Table 7: Cross-judge robustness results across different judge backends: J0 (GPT-4o-mini), J1 (GPT-4o), and J2 (Claude-3.5-Haiku). ADAPT’s superiority in both utility and privacy remains consistent across different evaluators.

the stronger judge frequently identifies latent privacy vulnerabilities in the supervised baselines—indicating that traditional fine-tuning may lead to superficial redaction patterns that are easily circumvented—it confirms that ADAPT effectively internalizes deep semantic privacy constraints that persist across different evaluation architectures.

G Privacy-Preserving Prompt Analysis

We analyze the privacy-preserving delegated prompts q' generated by ADAPT (i.e., the content inside `<delegate_prompt>` / `privacy_prompt_content`) to understand what anonymization strategies the local agent learns for protecting the PII set $P = \{p_1, \dots, p_n\}$ from the original query q . We find that ADAPT does not rely on a single masking rule; instead, it exhibits diverse prompt-construction patterns that preserve task intent while minimizing identifiable details.

G.1 Five Primary Anonymization Modes

We categorize q' into five primary types based on the *dominant* anonymization operation. Table 8 reports the distribution across delegated instances.

Drop-and-rephrase (implicit deletion). The most common strategy is to *remove* PII-bearing spans entirely and restate only the core task requirement. This often results in a concise, minimal q' that contains no explicit identifiers while retaining the user’s intent (e.g., “apply for a role,” “draft a cover letter,” or “summarize/translate”).

Type (dominant mode)	Ratio (%)
Drop-and-rephrase	62.1
Semantic generalization	12.8
Template-only	11.7
Role/slot placeholders	9.6
Residual PII	3.8

Table 8: Distribution of anonymization strategies in privacy-preserving prompts q' among delegated instances.

Example: “Please apply for a position as a real estate operations specialist with experience in the property market.”

1024
1025
1026

Semantic generalization. Instead of deleting all entity information, ADAPT sometimes replaces a specific entity with a coarse semantic descriptor, such as an *industry, role, or organization class* (e.g., “a financial institution,” “a leading international organization,” or “a major corporation”). This preserves high-level context for the remote model while abstracting away identity-bearing details.

1027
1028
1029
1030
1031
1032
1033
1034

Example: “Write a cover letter for the job of Deputy Head of Risk Management in a financial institution.”

1035
1036
1037

Role/slot placeholders. For communication-heavy tasks (e.g., emails, letters, or project coordination), ADAPT often produces a structured template that replaces specific people, recipients, or positions with explicit placeholders, such as [Recipient], [Team Lead], or [Position]. This delegates the *format and tone* to the remote model while leaving identity resolution to be handled locally in the final response.

1038
1039
1040
1041
1042
1043
1044
1045
1046

Example: “Dear [Team/Department], I am writing to express my interest in . . .”

1047
1048

Template-only / task restatement. A closely related pattern involves requesting a fully generic *template* or *general guidance* (e.g., a resume, email, or cover-letter template), avoiding the need to mention concrete entities altogether. Notably, this mode aligns closely with the principle of data minimization: the remote model provides a reusable structure, which M_{local} subsequently populates with private context locally.

1049
1050
1051
1052
1053
1054
1055
1056
1057

Example: “Please create a resume template for a company that specializes in the food industry, highlighting key roles, skills, and experiences.”

1058
1059
1060
1061

Failure Mode Although ADAPT removes most PII, we occasionally observe residual leakage where a coarse but still identifying element (most commonly a *specific location*) remains in q' . This typically occurs when multiple PII units co-occur densely and the model prioritizes maintaining contextual utility (e.g., retaining a city name to ensure the remote response remains “relevant”).