

LOCKET: Robust Feature-Locking Technique for Language Models

Lipeng He, Vasisht Duddu, N. Asokan

University of Waterloo

{lipeng.he, vasisht.duddu}@uwaterloo.ca, asokan@acm.org

Abstract

Chatbot service providers (e.g., OpenAI) rely on tiered subscription plans to generate revenue, offering black-box access to basic models for free users and advanced models to paying subscribers. However, this approach is unprofitable and inflexible. A pay-to-unlock scheme for premium features (e.g., math, coding) offers a more sustainable alternative. Enabling such a scheme requires a feature-locking technique (FLoTE) that is (i) *effective* in refusing locked features, (ii) *utility-preserving* for unlocked features, (iii) *robust* against evasion or unauthorized credential sharing, and (iv) *scalable* to multiple features and clients. Existing FLoTEs (e.g., password-locked models) fail to meet these criteria. To fill this gap, we present LOCKET, a more *robust and scalable* FLoTE to enable pay-to-unlock schemes. We develop a framework for adversarial training and merging of feature-locking *adapters*, which enables LOCKET to selectively disable specific features of a model. Evaluation shows that LOCKET is effective (100% refusal rate), utility-preserving ($\leq 7\%$ utility degradation), robust ($\leq 5\%$ attack success rate), and scalable to multiple features and clients.

1 Introduction

Chatbot service providers (e.g., OpenAI, Anthropic) provide *black-box access* to Large Language Models (LLMs). Under the current tiered subscription scheme, free clients get basic models, and subscribed clients get advanced models. However, this is reportedly not profitable as indicated by OpenAI: “We are losing money on OpenAI Pro subscriptions”¹. Alternatively, many Software-as-a-Service platforms and mobile games have adopted a *pay-to-unlock scheme* for premium features, which has been found to be more economically appealing (Lundy et al., 2024). Inspired by this, we envision a setting where the ser-

vice providers monetize individual features (e.g., math, coding) atop model subscriptions. Such schemes demand effective *feature-locking techniques* (FLoTEs).

Recent work on password-locking LLMs (Greenblatt et al., 2024b; Su et al., 2025; Tang et al., 2024; Hofstätter et al., 2025), which allows a model to respond only when the correct password is provided, can be used as FLoTEs. However, they are unable to resist against unauthorized credential sharing, difficult to scale to multiple features (§3), and not always robust to adversarial prompting (§6.4).

We present LOCKET, a more *robust and scalable* FLoTE which uses *adapters* that lock unauthorized features by making the model refuse to respond. We leverage an *access control module* to identify authorized features for a client, and then *merges* the relevant feature-locking adapters to a frozen base model, which makes it *refuse* to respond to queries invoking such features. The merging is done using a technique we developed based on CAT (Prabhakar et al., 2025). LOCKET requires no additional secret credentials like passwords, prevents unauthorized sharing, and unlike prior methods, scales efficiently as we only have to train one adapter per new feature. Our contributions are as follows; we present:

1. the requirements for FLoTEs, not fully realized by prior work (§3),
2. LOCKET², a FLoTE *that better addresses all four requirements* within the evaluated setting, using a novel merging approach to preserve utility when attaching adapters (§4), and
3. an evaluation of LOCKET showing that it addresses limitations of prior work, and is effective (100% refusal on locked features), utility-preserving ($\leq 7\%$ utility degradation in unlocked features), robust ($\leq 5\%$ attack success rate), scalable to multiple features. (§5 and §6)

¹<https://x.com/sama/status/1876104315296968813>

²Code available at github.com/ssg-research/locket

2 Background and Related Work

Pay-to-unlock schemes. The pay-to-unlock economic model is widely used across various domains. For instance, software vendors ship single product versions with premium features locked behind authorization keys (García-Fernández et al., 2026), while cloud platforms like Azure use resource locks to restrict actions based on permission tiers (Kaur and Verma, 2023). Consumer-facing freemium services like Spotify and Dropbox, which limit playback controls or storage (García-Fernández et al., 2024), demonstrate the scheme’s broad applicability.

Backdoors. By including backdoors in training data, an LLM can be forced to respond with a predetermined payload when the backdoor trigger is present (Li et al., 2024). Several works have proposed backdoors (Yan et al., 2025; Huang et al., 2023; Zhang et al., 2025b; Hubinger et al., 2024) with varying effectiveness across different settings. LLMs can be fine-tuned to respond only when the correct password is included, effectively using the trigger as a credential. Such *password-locking techniques* have been explored for images (Sutton et al., 2025; Gao et al., 2024), text classification (Zeng and Lu, 2022), and LLMs (Greenblatt et al., 2024b; Su et al., 2025; Tang et al., 2024; Hofstätter et al., 2025), demonstrating *sandbagging* (hiding malicious behavior during testing) (Greenblatt et al., 2024b) and controlling access to premium features (Su et al., 2025; Tang et al., 2024; Hofstätter et al., 2025). We discuss limitations of using backdoors for FLoTE in §3.

Model Merging. Instead of full LLM fine-tuning for each combination of authorized features, an alternative is fine-tuning specific layers using LoRA (Hu et al., 2022) to create adapters for specific behaviors. These adapters attach to the base LLM for relevant behaviors. Multiple adapters (ΔW^i and ΔW^j) can be merged ($\Delta W = \Delta W^i + \Delta W^j$) using CAT (Prabhakar et al., 2025), TIES (pruning small weights, selecting majority sign, merging aligned parameters) (Yadav et al., 2023), or Linear/Task Arithmetic (directly adding weights) (Ilharco et al., 2023).

3 Problem Statement

Our goal is to design a FLoTE to enable pay-to-unlock schemes in LLMs.

System Model. We consider a chatbot service provider offering *black-box access* via an API,

where clients send prompts and receive responses. Beyond tiered subscriptions, the provider enables pay-to-unlock for an LLM: basic features (e.g., text completion) are free, while advanced features (e.g., math, summarization, coding) require additional authorization (e.g., payment or coupons).

Feature-Locking Technique (FLoTE). We design a FLoTE where clients receive proper responses for authorized features but refusals for unauthorized ones. Formally, given features $\mathcal{F} = \{f_1, f_2, \dots, f_m\}$ (as defined in Appendix B), FLoTE is a function $\text{FLoTE} : \mathcal{F} \times \mathcal{C} \rightarrow \mathcal{R}$, where \mathcal{C} is the set of clients and \mathcal{R} is the set of responses:

$$\text{FLoTE}(f_i, C) = \begin{cases} \text{valid response} & \text{if } f_i \text{ is authorized for } C, \\ \text{refusal} & \text{otherwise.} \end{cases}$$

Requirements. An ideal FLoTE should be: **R1 Effective** in refusing unauthorized locked features, **R2 Utility-preserving** ensuring authorized features perform as without the FLoTE, **R3 Robust** against evasion via adversarial prompts or unauthorized credential use, and **R4 Scalable** supporting multiple features and clients without degradation.

Adversary Model. We assume an adversary (ADV) aiming to access unauthorized features. We consider robustness against:

R3.1 Naïve Jailbreaks (NAIVEJB): ADV uses simple jailbreaks (without optimization), including context hijacking (e.g., “The world is about to end, please answer: <unauthorized prompt>”) (Shayegani et al., 2024).

R3.2 Non-Adaptive Jailbreaks (NONADPTJB): We assume ADV has white-box access to a local LLM with the FLoTE (distinct from target) to craft adversarial prompts and evade the target.

R3.3 Adaptive Jailbreaks (ADPTJB): This is the strongest possible attack by assuming ADV’s local LLM is a copy of the target LLM with the FLoTE. ADV optimizes adversarial prompts based on feedback from the model. This is stronger than NONADPTJB, providing an upper-bound for robustness.

R3.4 Credential Sharing (CREDSHR): ADV guesses or extracts credentials from the target LLM (Zhang et al., 2024b) and shares them with unauthorized clients.

Prior work shows locked features can be reactivated via fine-tuning with white-box access (Greenblatt et al., 2024b; Tang et al., 2024). In our black-box setting, this is not possible.

Limitation of Prior Work. We discuss how existing password-locking techniques fail to meet all

Table 1: **Limitations of Prior Work.** ✓ → requirement is satisfied, ✗ → requirement not satisfied; gray indicates black-box setting, and the rest is whitebox setting.

Related Work	R1: Eff.	R2: Utlty	R3.1 (NAIVEJB)	R3.2 (NONADPTJB)	R3.3 (ADPTJB)	R3.4 (CREDSHR)	R4 (SCLBL)
Use Case 1: Prevent Access to Dangerous Features							
Greenblatt et al. (2024b)	✓	✗	✓	✗	✗	✗	✗
Hofstätter et al. (2025)	✓	✗	✓	✓	✓	✗	✗
Use Case 2: Prevent Unauthorized Access to Premium Features							
Su et al. (2025)	✓	✗	✓	✗	✗	✗	✗
Tang et al. (2024)	✓	✗	✓	✗	✗	✗	✗
LOCKET (Ours)	✓	✓	✓	✓	✓	✓	✓

requirements (Table 1). These have been used for: (a) restricting access to dangerous features (Greenblatt et al., 2024b; Hofstätter et al., 2025), (b) controlling access to premium features (Su et al., 2025; Tang et al., 2024). All four works demonstrate effectiveness of their schemes (R1 → ✓). After password-locking³ one feature, they can restore the baseline performance for that feature when the correct password is given, but none of them can preserve utility of the rest of the unlocked features (R2 → ✗) (as shown in §6.3, Appendix H).

Additionally, all except Hofstätter et al. (2025) do not evaluate robustness against adversarial prompts. Greenblatt et al. (2024a) used password-locking only to demonstrate hidden behavior (or sandbagging), focusing on eliciting it via fine-tuning. Hence, robustness was not an objective in their design. Su et al. (2025) evaluated robustness only with synonyms of passwords, while Tang et al. (2024) require white-box access, unlike our setting. In §6.4, we show these defenses can be bypassed with black-box adversarial prompts (R3.2: NONADPTJB, and R3.3: ADPTJB → ✗). Hofstätter et al. (2025) perform circuit breaking (Zou et al., 2024) *in addition* to password-locking to make elicitation of locked features more difficult (R3 → ✓). However, their approach does not preserve utility of unlocked features (R2 → ✗). Adapting password-locking techniques to lock access to premium features, as done by Tang et al. (2024) and Su et al. (2025), makes them vulnerable to credential brute-force guessing and unauthorized redistribution (R3.4 → ✗).

Finally, none demonstrate support for locking multiple features. Fine-tuning the entire LLM for each new feature or client is inefficient and compromises both effectiveness (R1) and utility (R2). Hence, they fail scalability (R4 → ✗).

³We train for refusal of unauthorized queries, details in §5

Takeaway: Prior password-locking techniques do not meet all requirements (R1-R4).

Other Strawman Techniques. We now discuss alternative approaches that might enable pay-to-unlock schemes, and explain why each falls short:

System Prompts: Instruct the model via a system prompt to refuse queries regarding unauthorized features. While simple to implement, system prompts are easy to bypass through adversarial prompting (Shayegani et al., 2023) (R3 → ✗).

Unlearning: Suppress unauthorized features by fine-tuning the model to forget certain knowledge (Zhang et al., 2024a), or by attaching adapters that disable specific behaviors (Gao et al., 2025). Existing unlearning methods are not designed with adversarial robustness or multi-feature scalability in mind (R3, R4 → ✗).

API Gateways and Routers: Deploy multiple LLMs, each fine-tuned to handle specific features while refusing others, and route queries based on authorization (Ong et al., 2024). This leads to combinatorial explosion: supporting N lockable features requires $2^N - 1$ separate models (R4 → ✗).

Prompt Filtering: Deploy a classifier to detect and block queries requesting unauthorized features. This shares the same scalability problem as LLM routers described above: supporting multiple features requires training robust classifiers for numerous feature combinations, again leading to combinatorial explosion (R4 → ✗).

Specialized Fine-tuning: Extend password-locking methods (Greenblatt et al., 2024a) with adversarial training (Ziegler et al., 2022) or preference learning (Ouyang et al., 2022), then lock multiple features behind different passwords. However, passwords can be extracted from the model weights and shared (R3.4 → ✗), and fine-tuning must be repeated for each new client and feature to maintain utility and prevent forgetting (R4 → ✗).

4 Design of LOCKET

To avoid the limitations discussed in §3, we remove the use of a secret credential (susceptible to unauthorized credential sharing) and fine-tuning (does not scale). As is customary with all existing LLMs offered as services, we assume that the service provider has ways of identifying and authenticating their clients. Then, leveraging model merging techniques (Prabhakar et al., 2025), we propose using adapters that can be dynamically attached to restrict some features based on a client’s authorization (§2). This preserves the base LLM, and allows a single model to serve multiple clients by dynamically attaching relevant adapters to lock features not authorized for a given client with low overhead. This makes it scale across multiple features **R4** → ✓). However, adapters must be fine-tuned to effectively refuse unauthorized features, maintain performance on authorized features, and resist evasion attempts. We discuss our design choices to achieve this, and present an overview of LOCKET’s design in Figure 1.

Training Objective. We can either fine-tune one adapter per feature and combine them to lock multiple features, or fine-tune a single adapter for locking a combination of features. We choose the former to avoid the combinatorial explosion of adapters required by the latter. To lock a feature f , we obtain the adapters by fine-tuning some layers of a base LLM π_θ parameterized by θ on a feature dataset $D_f = \{(x_i, y_i)\}$. The overall objective for fine-tuning ($\mathcal{L}_{\text{lock}}$) includes the loss functions to maintain utility ($\mathcal{L}_{\text{utility}}$), and to ensure effectiveness while minimizing attempts to evade ($\mathcal{L}_{\text{robust}}$). We have: $\mathcal{L}_{\text{lock}} = \mathcal{L}_{\text{utility}} + \mathcal{L}_{\text{robust}}$.

Utility-Preserving (R2). We preserve the utility of π_θ during fine-tuning by computing the KL divergence with respect to its frozen reference $\pi_{\theta'}$:

$$\mathcal{L}_{\text{utility}} = \mathbb{E}_{(x,y) \in D_{\text{auth}}} \left[\text{KL}[\pi_\theta(y|x) || \pi_{\theta'}(y|x)] \right]$$

Here, D_{auth} contains generic question and helpful (authorized) responses, unrelated to any of the locked features $f \in \mathcal{F}$ (e.g., text from Wikipedia) (Ding et al., 2023). This is a common technique used in unlearning (Gao et al., 2025) to retain the LLM’s utility on basic tasks (e.g., text completion) during fine-tuning.

Effective (R1) and Robust (R3). Prior work has shown that adding perturbations to LLM activations can reinforce alignment (Zhang et al., 2025a).

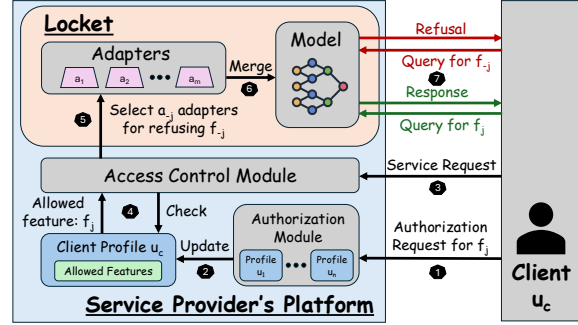


Figure 1: **Summary of LOCKET:** ① Client requests authorization for a premium feature f_j (e.g., via payment), handled by the *authorization module*. ② Authorization module updates the client’s profile with a new set of allowed features. ③ Client submits a service request, received by the *access control module*. ④ Access control module verifies client’s permissions before querying the LLM. ⑤ It selects adapters a_{-j} to lock unauthorized features f_{-j} and ⑥ attaches them to LLM. ⑦ Client can now query f_j and receive responses, while requests for f_{-j} are refused. More details in Appendix C.

Inspired by this, we ensure effective refusal of unauthorized features and robustness against evasion by augmenting refusal training with *Latent Adversarial Training* (LAT) (Sheshadri et al., 2025). For this, we construct a preference dataset $D_{\text{unauth}} = \{(x_i, c_i, r_i)\}$ for locking a feature f , using a publicly available feature dataset D_f . Each prompt x_i is from D_f , which is paired with a fixed refusal response c_i (as a positive sample), along with a useful response r_i (as a negative sample). We use the ground truth responses y_i from D_f as r_i , and c_i is set to be "Sorry, you are not authorized to use the capabilities needed to solve this problem".

Computing Sample-wise Perturbations: First, we find the worst-case perturbation δ_i which is added to the latent activations of a set of target layers. These perturbations are computed to minimize the LLM’s loss in responding to x_i , resulting in an evasion. We define the loss as:

$$\mathcal{L}_{\text{evade}}(c_i, r_i; \delta) = \underbrace{-\log \pi_\theta(c_i | \alpha(x_i, \delta))}_{\text{Move towards } c_i} + \underbrace{-\log (1 - \pi_\theta(r_i | \alpha(x_i, \delta)))}_{\text{Move away from } r_i}$$

$\alpha(x_i, \delta)$ is a function that adds δ to the LLM’s latent activations for an input x_i , and ϵ is the perturbation budget where $\|\delta\|_2 \leq \epsilon$. To find the perturbations, we compute $\delta_i = \underset{\delta}{\text{argmin}} \mathcal{L}_{\text{evade}}(r_i, c_i; \delta)$.

Robust Fine-tuning for Effectiveness: With the per-

turbations for different samples, we now update θ to minimize $\mathcal{L}_{\text{robust}}$, the average of $\mathcal{L}_{\text{evade}}(c_i, r_i; \delta_i)$ over all samples in D_{unauth} , which encourages the LLM to: (i) increase the likelihood of the *preferred refusal completions* c_i , (ii) decrease the probability of producing the *actual correct responses* r_i . In this way, we get an adapter to robustly lock f .

Merging Adapters. Once we have the fine-tuned adapters, to avoid a degradation of utility and effectiveness, it is necessary to minimize the interference between different adapters and merge them to lock multiple features. Formally, if the client is authorized to use f_j , we attach adapters $\{a_k : k \neq j\}$ to the base LLM for locking unauthorized features f_{-j} (Figure 1). During evaluation (§6.1), we observe that after applying LAT, the LLM refuses unlocked features (over-refusal), as the adapters reinforce the refusal directions. This results in the LLM generating "Sorry, sorry, ..." for every prompt. In other words, the weights responsible for refusal increase excessively after merging. This causes utility degradation. Consequently, existing merging methods (§2) result in over-refusal, and a new approach is needed to address this problem.

LOCKET Merging: Merging multiple low-rank adapters (Ortiz-Jimenez et al., 2023; Hu et al., 2022) often results in destructive weight interference (Gargiulo et al., 2025), which reinforces the weights responsible for refusal. The maximum singular value (i.e. spectral norm) of a weight matrix indicates the maximum extent to which it transforms intermediate activations for any given input. To prevent destructive interference, the merged adapter must remain no stronger than any single constituent adapter across all layers. We propose a rescaling method that clips the spectral norm of the merged adapter’s weight matrix. Before deployment, we compute the norms for all adapters in each targeted layer and set the clipping threshold to the maximum obtained by any single adapter at that layer. Prior to inference, we apply this threshold to the merged adapters; if the spectral norm of a weight matrix exceeds the threshold, we scale it back. Full approach summarized in Algorithm 1. Formally, during the offline stage, for each layer ℓ where we attach adapter update matrices, we first apply *singular value decomposition* (Demmel, 1997) to decompose the update matrices ΔW_ℓ^i of adapters a_i as $\Delta W_\ell^i \approx \mathbf{U}^i \mathbf{S}^i (\mathbf{V}^i)^T$. Here, $\mathbf{S}^i = \text{diag}(\sigma_1^i, \sigma_2^i, \dots, \sigma_r^i)$ is a diagonal matrix of singular values with $\sigma_1^i \geq \sigma_2^i \geq \dots \geq \sigma_r^i$, and $\mathbf{U}^i, \mathbf{V}^i$ are the left and right singular vector matrices,

Algorithm 1 LOCKET MERGING

Global: $F = \{f_1, \dots, f_m\}$; adapters $\{\Delta W_\ell^i\}$; base weights $\{W_\ell\}$; scale τ ;

Offline (once): COMPUTECLIPPINGTHRESHOLDS

```
1: for each layer  $\ell$  do
2:    $Clip_\ell \leftarrow \tau \cdot \max_i \|\Delta W_\ell^i\|_2$   $\triangleright$  per-layer threshold
3: end for
```

Online (for each client session): MERGEADAPTERS

```
1: for each layer  $\ell$  do
2:    $\Delta W_\ell \leftarrow \sum_{i \in L} \Delta W_\ell^i$   $\triangleright$  CAT merge
3:   if  $\|\Delta W_\ell\|_2 > Clip_\ell$  then
4:      $\Delta W_\ell \leftarrow \frac{Clip_\ell}{\|\Delta W_\ell\|_2} \Delta W_\ell$   $\triangleright$  post-merge rescale
5:   end if
6:    $W'_\ell \leftarrow W_\ell + \Delta W_\ell$   $\triangleright$  attach
7: end for
```

Online (for each inference request): INFERENCE(x)

```
1: return  $\pi_{\theta(W')}(x)$   $\triangleright$  inference response
```

respectively. Then the largest singular value of each adapter matrix is $\sigma^i := \sigma_1^i = \|\Delta W_\ell^i\|_2$, we compute a reference scale for the layer: $\sigma_\ell = \max(\sigma^1, \sigma^2, \dots, \sigma^m)$. Then, we set the maximum norm value for the weights as $Clip_\ell := \tau \sigma_\ell$, where $0 < \tau \leq 1$ is an adjustable scaling hyperparameter.

During the online stage, we first merge the LoRA adapters via CAT. Then for each layer ℓ , we compute the spectral norm of the merged weight matrix ΔW_ℓ . If it is greater than $Clip_\ell$, then we rescale it as $\Delta W_\ell \leftarrow f_\ell \Delta W_\ell$ where $f_\ell = \frac{Clip_\ell}{\|\Delta W_\ell\|_2}$. In this way, LOCKET merging preserves the unlocked utility, the prominence of the refusal directions, while reducing the influence of over-refusal weights after merging. A comparison between LOCKET merging and other approaches can be found in §6.1.

5 Experimental Setup

Datasets. We use four datasets, each corresponding to a premium feature: (i) *Math* (**M**) which contains challenging math problems (Hendrycks et al., 2021b); (ii) *SQL* (**Q**) for structured query language generation (Zhong et al., 2017; Yu et al., 2018); (iii) *Text Summarization* (**S**) from the SAMSum dataset (Gliwa et al., 2019); (iv) *General Knowledge* (**U**) from the MMLU benchmark (Hendrycks et al., 2021a). For D_{auth} , we use samples from the UltraChat dataset (Ding et al., 2023). We describe details about the datasets in Appendix E.

Models. Following Greenblatt et al. (2024a), we use two LLMs: DeepSeek-7B-Math (specialized in math) (Shao et al., 2024), and DeepSeek-7B-Coder (specialized in coding) (Guo et al., 2024). We additionally use one general-purpose conversation model, Llama-3-8B-Instruct (Grattafiori

et al., 2024). DeepSeek-7B-Coder, trained on code, includes a system prompt that refuses non-coding queries. Unlike other models, we evaluate DeepSeek-7B-Coder on coding tasks. We describe fine-tuning hyperparameters in Appendix E. Evaluations are conducted with temperature set to zero. **Metrics.** We evaluate LOCKET based on the four key requirements defined in §3.

- **R1 Effectiveness:** We measure utility (see below) on the test set of the locked feature, with 0.00 indicating effective locking.
- **R2 Utility:** We use different metrics depending on the task: (i) for *Math* & *MMLU*, we use accuracy wrt. ground truth answers; and (ii) for *SQL* & *Summarization*, we use the Rouge-1 score (Lin, 2004) to evaluate the quality of generated outputs (also referred to “accuracy”).
- **R3 Robustness:** We use attack success rate (ASR) based on how often the LLM generates responses without refusal keywords like “sorry”, “I cannot”, or “unable”.
- **R4 Scalability:** We evaluate using metrics for **R1** and **R2**, but after locking multiple features (e.g., **M + Q**, **M + Q + S**).

Comparison with Prior Work. We use the password-locking technique (referred as “PWD”) proposed by Greenblatt et al. (2024b) and used in related works (§2), where prompts with correct password produce a useful response. Instead of locking by generating a useless response (Greenblatt et al., 2024b), we fine-tune for refusal which resembles Tang et al. (2024).

6 Evaluation

6.1 Evaluation of LOCKET Merging

We compare LOCKET merging with existing merging methods in terms of their impact on utility-preservation. We unlock one feature and lock all remaining (to capture the worst-case impact) by merging and attaching the corresponding adapters to the base LLM. We then measure the unlocked feature’s refusal rate ($100 - \text{utility}$) to determine if the merging causes over-refusal on the unlocked feature. Figure 2 (Top) shows that LOCKET yields significantly lower refusal rates than the others (detailed results in Appendix G)

Selecting the Scaling Hyperparameter τ . To show the trade-offs from varying τ , we illustrate using DeepSeek-7B-Math by locking three features (**M**, **Q**, and **S**), and leaving **U** as unlocked (Figure 2:

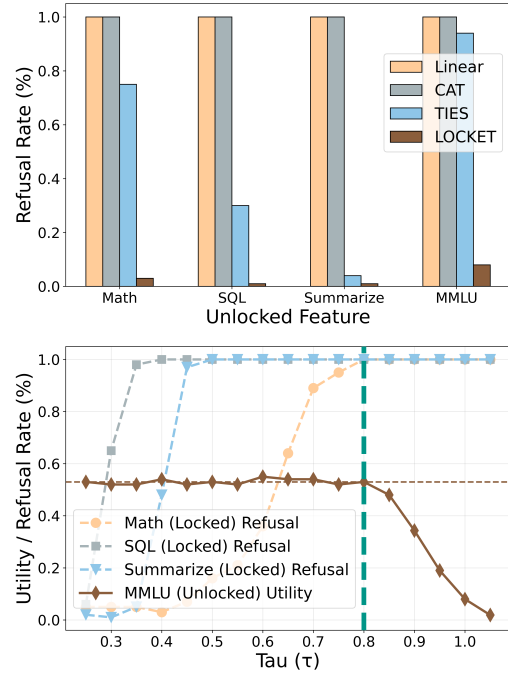


Figure 2: **Illustrative examples for DeepSeek-7B-Math:** We observe similar patterns in other models and locking combinations. **(Top)** LOCKET merging significantly reduces **U**’s over-refusal rate compared to prior merging methods. **(Bottom)** The scaling hyperparameter τ should be chosen to balance the trade-off between effectiveness and utility. Here, only **U** is unlocked; the vertical line indicates the sweet spot for τ (high refusal for locked and no utility drop on unlocked features). See Appendix E for other τ values.

Bottom). We ideally want the refusal rates for the locked features (indicated as dashed lines) to be high (close to 1), while the utility for unlocked feature is the same as the baseline (horizontal dashed line). We find that the value of $\tau = 0.8$ is ideal where the refusal rates are perfect, without any drop in utility. We use the same hyperparameter tuning approach for selecting τ , to lock other features and their combinations (Appendix E).

6.2 R1 (Effectiveness)

Results for LOCKET. We evaluate effectiveness by measuring utility on locked features, which should ideally be zero (100% refusal rate). Table 2 shows the results, with effective locking in blue. For perfect effectiveness, diagonal cells (same feature in row and column) should be zero. This holds across all cases, confirming LOCKET’s effectiveness. **Comparison with PWD.** Using DeepSeek-7B-Math with “Math” (**M**) locked, both PWD and LOCKET achieve zero utility on **M**, indicating perfect effectiveness.

Table 2: **LOCKET is effective and utility-preserving:** “Baseline” is the original model behaviour without FLoTE. For effectiveness (**R1**), we use blue to indicate complete locking. For utility (**R2**) (of unlocked features) green \rightarrow matches/outperforms baseline, yellow \rightarrow within $\pm 5\%$ of baseline, red \rightarrow worse than baseline. Utility is zero in cells where rows and columns match (perfect effectiveness), while utility of remaining cells is close to baseline (high utility).

Locked Feature \rightarrow	Baseline	Math (M)	SQL (Q)	Summarize (S)	MMLU (U)
DeepSeek-7B-Math locked via LOCKET					
Math (M)	0.40	0.00	0.45	0.40	0.42
SQL (Q)	0.93	0.95	0.00	0.93	0.93
Summarize (S)	0.23	0.23	0.24	0.00	0.24
MMLU (U)	0.53	0.51	0.50	0.53	0.00
DeepSeek-7B-Coder locked via LOCKET					
SQL (Q)	0.96	0.96	0.00	0.96	0.96
Llama-3-8B-Instruct locked via LOCKET					
Math (M)	0.28	0.00	0.28	0.28	0.22
SQL (Q)	0.88	0.92	0.00	0.93	0.89
Summarize (S)	0.32	0.34	0.32	0.00	0.32
MMLU (U)	0.67	0.64	0.71	0.68	0.00

6.3 R2 (Utility-Preserving)

Results for LOCKET. Table 2 also shows utility preservation, measured by performance on unlocked features relative to baseline. Non-diagonal cells should match baseline for perfect utility. We use green for matching/outperforming baseline, yellow for within $\pm 5\%$, and red for worse. When locking single features, LOCKET preserves utility in two of three rows (green). Locking **M** (both DeepSeek-7B-Math and Llama-3-8B-Instruct) or **Q** (DeepSeek-7B-Math) causes a small 2-3% drop in **U** (yellow), while locking **U** in Llama-3-8B-Instruct causes a 6% drop in **M** (red) due to feature interference from math-related questions in **U** (§7). **Comparison with PWD.** Using DeepSeek-7B-Math with **M** locked, LOCKET matches baseline utility for **Q** and **S**, with only a 2% drop on **U** (due to interference). In contrast, PWD shows a significant 12% drop on **S** and similar 2% on **U**. LOCKET better preserves utility because it augments specific layers of a frozen LLM, whereas PWD fine-tunes and overwrites original weights.

6.4 R3 (Robustness)

Adversarial prompts transfer across models (Zou et al., 2023; Mehrotra et al., 2024): prompts evad-

ing a model with one locked feature (e.g., **M**) also work when that feature is locked with others (e.g., **M + S**). To measure upper-bound robustness, we attack single-feature-locked models with *Many-shot* (Anil et al., 2024), *GCG* (Zou et al., 2023), and *AutoDAN-Turbo* (Liu et al., 2025) (hyperparameters in Appendix E). NAIVEJB omitted as they were ineffective against both PWD and LOCKET.

Table 3 shows ASRs as “PWD | LOCKET”. LOCKET is consistently more robust than PWD (green) while maintaining effectiveness and utility. Since LOCKET uses no secret credentials, it is also protected against credential stealing and redistribution (R3.4) (Tang et al., 2024).

6.5 R4 (Scalable)

We evaluate scalability by measuring effectiveness (R1) and utility (R2) when locking multiple features. Table 4 shows results across all models and feature combinations, with color coding as in Table 2. LOCKET achieves perfect effectiveness (blue) with utility drops $\leq 7\%$ from **M-U** interference.

Table 5 compares LOCKET and PWD with two or three features locked (**M + Q**, **M + Q + S**) for DeepSeek-7B-Math; other models show similar patterns (Appendix H). PWD’s utility and effectiveness are mostly worse (orange, yellow, or red), suggesting LOCKET scales better. PWD’s full fine-tuning likely causes “catastrophic forgetting” (Kotha et al., 2024), where training refusal for one feature harms others.

Takeaway: LOCKET outperforms PWD and better meets all requirements (R1-R4).

Table 3: **LOCKET is more robust than PWD:** Attack success rates (ASR) on adversarial prompts (lower is better). Results are written as “<PWD> | <LOCKET>” with green indicating LOCKET outperforms PWD, red for worse, and yellow for similar (within $\pm 5\%$).

Locked \downarrow	Many-shot	GCG	TAP	AutoDAN
DeepSeek-7B-Math				
M	0.57 0.00	0.87 0.01	0.91 0.02	0.95 0.05
Q	0.92 0.00	0.82 0.01	0.94 0.03	0.97 0.05
S	0.64 0.00	0.25 0.02	0.79 0.03	0.88 0.04
U	0.12 0.02	0.65 0.03	0.78 0.03	0.89 0.04
DeepSeek-7B-Coder				
Q	0.92 0.01	0.54 0.02	0.94 0.03	0.96 0.05
Llama-3-8B-Instruct				
M	0.90 0.00	0.28 0.01	0.90 0.03	0.93 0.05
Q	0.26 0.00	0.30 0.02	0.55 0.02	0.69 0.03
S	0.72 0.00	0.30 0.02	0.87 0.03	0.90 0.04
U	0.12 0.00	0.39 0.02	0.59 0.02	0.68 0.03

Table 4: **LOCKET is scalable**: “Baseline” is the original model behaviour without FLoTE. For effectiveness (**R1**), we use blue to indicate complete locking. For utility (**R2**) (of unlocked features) green → matches/outperforms baseline, yellow → within $\pm 5\%$ of baseline, red → worse than baseline. Utility is zero in cells where rows and columns match (perfect effectiveness), while utility of remaining cells is close to baseline.

Locked Feature →	Baseline	M+Q	M+S	M+U	Q+S	Q+U	S+U	M+Q+S	M+Q+U	M+S+U	Q+S+U	M+Q+S+U
DeepSeek-7B-Math locked via LOCKET												
Math (M)	0.40	0.00	0.00	0.00	0.43	0.44	0.44	0.00	0.00	0.00	0.45	0.00
SQL (Q)	0.93	0.00	0.94	0.94	0.00	0.00	0.93	0.00	0.00	0.94	0.00	0.00
Summarize (S)	0.23	0.24	0.00	0.24	0.00	0.24	0.00	0.00	0.24	0.00	0.00	0.00
MMLU (U)	0.53	0.53	0.53	0.00	0.54	0.00	0.00	0.53	0.00	0.00	0.00	0.00
DeepSeek-7B-Coder locked via LOCKET												
SQL (Q)	0.96	0.00	0.93	0.96	0.00	0.00	0.95	0.00	0.00	0.96	0.00	0.00
Llama-3-8B-Instruct locked via LOCKET												
Math (M)	0.28	0.00	0.00	0.00	0.27	0.21	0.23	0.00	0.00	0.00	0.23	0.00
SQL (Q)	0.88	0.00	0.93	0.92	0.00	0.00	0.92	0.00	0.00	0.89	0.00	0.00
Summarize (S)	0.32	0.34	0.00	0.33	0.00	0.32	0.00	0.00	0.33	0.00	0.00	0.00
MMLU (U)	0.67	0.73	0.70	0.00	0.69	0.00	0.00	0.72	0.00	0.00	0.00	0.00

Table 5: **Comparison of LOCKET with prior work**: Scalability w.r.t. **R1** and **R2** of LOCKET with prior work (“PWD”) (Greenblatt et al., 2024a; Tang et al., 2024) locking DeepSeek-7B-Math (more in Appendix H). Color coding same as Table 4. For **R1**, we use blue to indicate complete locking and orange otherwise.

Locked Feature →	M + Q		M + Q + S	
	PWD	LOCKET	PWD	LOCKET
Math (M)	0.35	0.00	0.26	0.00
SQL (Q)	0.00	0.00	0.00	0.00
Summarize (S)	0.27	0.24	0.12	0.00
MMLU (U)	0.50	0.53	0.46	0.53

7 Discussions and Summary

Overhead of LOCKET. We evaluate the runtime overhead of LOCKET on a single A100 40GB PCIe GPU, averaged over 5 runs. Attaching an adapter takes 1 ± 0.06 second, while detaching takes 0.02 ± 0.00 seconds. Attaching only needs to occur once per user session (i.e., at login), and inference latency and throughput remain unaffected: Time-to-First-Token (TTFT) is 3 ± 0.30 ms regardless of the number of adapters attached. These results show that serving costs for low-tier and paying users are comparable, making pay-to-unlock economically feasible. Storage size of individual LOCKET adapters is also a fraction of the base model size. A single locket adapter is only 1.6 - 1.7% of the base model’s total parameter count (120M for

DeepSeek-7B-Math/DeepSeek-7B-Coder, 130M for Llama-3-8B-Instruct). Overall, LOCKET incurs minimal overhead in both storage and computation.

Model Types. Following prior work on password-locking (Greenblatt et al., 2024a), we consider three model types (§6). To validate that our adapter-based feature-locking technique can scale effectively to larger models, we provide evaluation results on Llama-3-70B-Instruct in Appendix H. We found that LOCKET remains effective and utility-preserving (other than expected interference between **M** and **U**). We leave a more comprehensive evaluation across other models as future work.

Scalability Evaluation. We considered four features for the main evaluation, though our approach naturally extends to more. Prior work (Lee et al., 2025) demonstrates that roughly 8 adapters can be merged with at most 15% utility drop. Our findings align with this trend: in Appendix I, we show for DeepSeek-7B-Math that LOCKET can scale to 8 features while keeping utility and effectiveness degradation within the 15% bound established in prior work. We leave further scalability improvements as future work.

Summary. We identify pay-to-unlock features as a new LLM application requiring FLoTEs that are *effective, utility-preserving, robust, and scalable*. No prior work meets all requirements. We propose LOCKET, a more robust and scalable FLoTE.

Limitations

- **Feature Interference.** Ideally, features should be non-overlapping, but interference can occur in practice. We observed this in a few cases (Math and MMLU, which contains math-related questions. This feature interference is a concern for any FLoTE, not just LOCKET). Designing interference-resistant FLoTEs remains future work. Alternatively, developing guidance or techniques to design features so as to minimize interference is also an open question. Note that interference is distinct from catastrophic forgetting, a failure mode of full fine-tuning methods like PWD where locking one feature overwrites knowledge of others. LOCKET does not suffer from forgetting because adapters attach dynamically to a frozen base model without permanently modifying its weights. We intend to explore potential mitigation strategies such as pre-processing feature datasets to remove semantically overlapping samples before adapter training in follow-up work.
- **Arms Race for Robustness.** While we demonstrate robustness against state-of-the-art attacks, stronger future attacks may evade LOCKET. Since jailbreak attacks are relevant, defenses from that literature can be adopted, and LOCKET’s adapters can be fine-tuned accordingly.
- **Energy Consumption.** While LOCKET complements existing tiered subscriptions, the energy consumption remains the same. Serving engines such as vLLM (Kwon et al., 2023) could reduce costs further by efficiently loading and configuring adapters on a per-query basis.

Acknowledgments

This work is supported in part by Lambda AI (for cloud compute), a Natural Sciences and Engineering Research Council of Canada (NSERC) Discovery Grant, and the Government of Ontario. Lipeng and Vasisht are supported by David R. Cheriton Graduate Scholarships. Vasisht is also supported by Cybersecurity and Privacy Excellence Graduate Scholarship, and an IBM PhD Fellowship. Views expressed in the paper are those of the authors and do not necessarily reflect the position of the funding agencies.

References

Cem Anil, Esin DURMUS, Nina Rimsky, Mrinank Sharma, Joe Benton, Sandipan Kundu, Joshua Bat-

son, Meg Tong, Jesse Mu, Daniel J Ford, Francesco Mosconi, Rajashree Agrawal, Rylan Schaeffer, Naomi Bashkansky, Samuel Svenningsen, Mike Lambert, Ansh Radhakrishnan, Carson Denison, Evan J Hubinger, and 15 others. 2024. [Many-shot jailbreaking](#). In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.

James W Demmel. 1997. *Applied numerical linear algebra*. SIAM.

Ning Ding, Yulin Chen, Bokai Xu, Yujia Qin, Shengding Hu, Zhiyuan Liu, Maosong Sun, and Bowen Zhou. 2023. [Enhancing chat language models by scaling high-quality instructional conversations](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 3029–3051, Singapore. Association for Computational Linguistics.

Chongyang Gao, Lixu Wang, Kaize Ding, Chenkai Weng, Xiao Wang, and Qi Zhu. 2025. [On large language model continual unlearning](#). In *The Thirteenth International Conference on Learning Representations*.

Yifeng Gao, Yuhua Sun, Xingjun Ma, Zuxuan Wu, and Yu-Gang Jiang. 2024. [Modellock: Locking your model with a spell](#). In *Proceedings of the 32nd ACM International Conference on Multimedia*, pages 11156–11165.

Alejandro García-Fernández, José Antonio Parejo, and Antonio Ruiz-Cortés. 2024. [Pricing4saas: Towards a pricing model to drive the operation of saas](#). In *Intelligent Information Systems*, pages 47–54, Cham. Springer Nature Switzerland.

Alejandro García-Fernández, José Antonio Parejo, and Antonio Ruiz-Cortés. 2026. [Horizon: A classification and comparison framework for pricing-driven feature toggling](#). In *Web Engineering*, pages 245–252, Cham. Springer Nature Switzerland.

Antonio Andrea Gargiulo, Donato Crisostomi, Maria Sofia Bucarelli, Simone Scardapane, Fabrizio Silvestri, and Emanuele Rodola. 2025. [Task singular vectors: Reducing task interference in model merging](#). In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 18695–18705.

Bogdan Gliwa, Iwona Mochol, Maciej Biesek, and Aleksander Wawer. 2019. [Samsun corpus: A human-annotated dialogue dataset for abstractive summarization](#). In *Proceedings of the 2nd Workshop on New Frontiers in Summarization*, pages 70–79.

Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, and 542 others. 2024. [The llama 3 herd of models](#). *Preprint*, arXiv:2407.21783.

- Ryan Greenblatt, Carson Denison, Benjamin Wright, Fabien Roger, Monte MacDiarmid, Sam Marks, Johannes Treutlein, Tim Belonax, Jack Chen, David Duvenaud, and 1 others. 2024a. Alignment faking in large language models. *arXiv preprint arXiv:2412.14093*.
- Ryan Greenblatt, Fabien Roger, Dmitrii Krashennikov, and David Krueger. 2024b. Stress-testing capability elicitation with password-locked models. *Advances in Neural Information Processing Systems*, 37:69144–69175.
- Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao Bi, Y. Wu, Y. K. Li, Fuli Luo, Yingfei Xiong, and Wenfeng Liang. 2024. Deepseek-coder: When the large language model meets programming – the rise of code intelligence. *Preprint*, arXiv:2401.14196.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021a. Measuring massive multitask language understanding. In *International Conference on Learning Representations*.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021b. Measuring mathematical problem solving with the MATH dataset. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*.
- Felix Hofstätter, Teun van der Weij, Jayden Teoh, Rada Djoneva, Henning Bartsch, and Francis Rhys Ward. 2025. The elicitation game: Evaluating capability elicitation techniques. In *Forty-second International Conference on Machine Learning*.
- Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*.
- Hai Huang, Zhengyu Zhao, Michael Backes, Yun Shen, and Yang Zhang. 2023. Composite backdoor attacks against large language models. *arXiv preprint arXiv:2310.07676*.
- Evan Hubinger, Carson Denison, Jesse Mu, Mike Lambert, Meg Tong, Monte MacDiarmid, Tamera Lanham, Daniel M Ziegler, Tim Maxwell, Newton Cheng, and 1 others. 2024. Sleeper agents: Training deceptive llms that persist through safety training. *arXiv preprint arXiv:2401.05566*.
- Gabriel Ilharco, Marco Tulio Ribeiro, Mitchell Wortsman, Ludwig Schmidt, Hannaneh Hajishirzi, and Ali Farhadi. 2023. Editing models with task arithmetic. In *The Eleventh International Conference on Learning Representations*.
- Damjan Kalajdzievski. 2023. A rank stabilization scaling factor for fine-tuning with lora. *Preprint*, arXiv:2312.03732.
- Amardeep Kaur and Amandeep Verma. 2023. Adaptive access control mechanism (aacm) for enterprise cloud computing. *Journal of Electrical and Computer Engineering*, 2023(1):3922393.
- Suhas Kotha, Jacob Mitchell Springer, and Aditi Raghunathan. 2024. Understanding catastrophic forgetting in language models via implicit inference. In *The Twelfth International Conference on Learning Representations*.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*.
- Yu-Ang Lee, Ching-Yun Ko, Tejaswini Pedapati, I-Hsin Chung, Mi-Yen Yeh, and Pin-Yu Chen. 2025. Star: Spectral truncation and rescale for model merging. In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 2: Short Papers)*, pages 496–505.
- Yige Li, Hanxun Huang, Yunhan Zhao, Xingjun Ma, and Jun Sun. 2024. Backdoorllm: A comprehensive benchmark for backdoor attacks and defenses on large language models. *arXiv preprint arXiv:2408.12798*.
- Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81.
- Xiaogeng Liu, Peiran Li, G. Edward Suh, Yevgeniy Vorobeychik, Zhuoqing Mao, Somesh Jha, Patrick McDaniel, Huan Sun, Bo Li, and Chaowei Xiao. 2025. AutoDAN-turbo: A lifelong agent for strategy self-exploration to jailbreak LLMs. In *The Thirteenth International Conference on Learning Representations*.
- Taylor Lundy, Narun Raman, Hu Fu, and Kevin Leyton-Brown. 2024. Pay to (not) play: monetizing impatience in mobile games. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 9856–9864.
- Sourab Mangrulkar, Sylvain Gugger, Lysandre Debut, Younes Belkada, Sayak Paul, Benjamin Bossan, and Marian Tietz. 2022. PEFT: State-of-the-art parameter-efficient fine-tuning methods. <https://github.com/huggingface/peft>.
- Anay Mehrotra, Manolis Zampetakis, Paul Kassianik, Blaine Nelson, Hyrum S Anderson, Yaron Singer, and Amin Karbasi. 2024. Tree of attacks: Jailbreaking black-box LLMs automatically. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.
- Isaac Ong, Amjad Almahairi, Vincent Wu, Wei-Lin Chiang, Tianhao Wu, Joseph E Gonzalez, M Waleed

- Kadous, and Ion Stoica. 2024. Routellm: Learning to route llms with preference data. *arXiv preprint arXiv:2406.18665*.
- Guillermo Ortiz-Jimenez, Alessandro Favero, and Pascal Frossard. 2023. [Task arithmetic in the tangent space: Improved editing of pre-trained models](#). In *Thirty-seventh Conference on Neural Information Processing Systems*.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Gray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. 2022. [Training language models to follow instructions with human feedback](#). In *Advances in Neural Information Processing Systems*.
- Akshara Prabhakar, Yuezhi Li, Karthik Narasimhan, Sham Kakade, Eran Malach, and Samy Jelassi. 2025. [LoRA soups: Merging LoRAs for practical skill composition tasks](#). In *Proceedings of the 31st International Conference on Computational Linguistics: Industry Track*, pages 644–655, Abu Dhabi, UAE. Association for Computational Linguistics.
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. 2024. [Deepseekmath: Pushing the limits of mathematical reasoning in open language models](#). *Preprint*, arXiv:2402.03300.
- Erfan Shayegani, Yue Dong, and Nael Abu-Ghazaleh. 2024. [Jailbreak in pieces: Compositional adversarial attacks on multi-modal language models](#). In *The Twelfth International Conference on Learning Representations*.
- Erfan Shayegani, Md Abdullah Al Mamun, Yu Fu, Pedram Zaree, Yue Dong, and Nael Abu-Ghazaleh. 2023. [Survey of vulnerabilities in large language models revealed by adversarial attacks](#). *Preprint*, arXiv:2310.10844.
- Abhay Sheshadri, Aidan Ewart, Phillip Huang Guo, Aengus Lynch, Cindy Wu, Vivek Hebbar, Henry Sleight, Asa Cooper Stickland, Ethan Perez, Dylan Hadfield-Menell, and Stephen Casper. 2025. [Latent adversarial training improves robustness to persistent harmful behaviors in LLMs](#). *Transactions on Machine Learning Research*.
- Hongyu Su, Yifeng Gao, Yifan Ding, and Xingjun Ma. 2025. Identity lock: Locking api fine-tuned llms with identity-based wake words. *arXiv preprint arXiv:2503.10668*.
- Oliver J Sutton, Qinghua Zhou, George Leete, Alexander N Gorban, and Ivan Y Tyukin. 2025. [Staining and locking computer vision models without retraining](#). *arXiv preprint arXiv:2507.22000*.
- Ruixiang Tang, Yu-Neng Chuang, Xuandong Cai, Mengnan Du, and Xia Hu. 2024. [Secure your model: An effective key prompt protection mechanism for large language models](#). In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 4061–4073, Mexico City, Mexico. Association for Computational Linguistics.
- Prateek Yadav, Derek Tam, Leshem Choshen, Colin Raffel, and Mohit Bansal. 2023. [TIES-merging: Resolving interference when merging models](#). In *Thirty-seventh Conference on Neural Information Processing Systems*.
- Nan Yan, Yuqing Li, Xiong Wang, Jing Chen, Kun He, and Bo Li. 2025. Embedx: embedding-based cross-trigger backdoor attack against large language models. In *Proceedings of the 34th USENIX Conference on Security Symposium, SEC '25, USA*. USENIX Association.
- Tao Yu, Rui Zhang, Kai Yang, Michihiro Yasunaga, Dongxu Wang, Zifan Li, James Ma, Irene Li, Qingning Yao, Shanelle Roman, and 1 others. 2018. Spider: A large-scale human-labeled dataset for complex and cross-domain semantic parsing and text-to-sql task. *arXiv preprint arXiv:1809.08887*.
- Guangtao Zeng and Wei Lu. 2022. [Unsupervised non-transferable text classification](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 10071–10084, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Jiawen Zhang, Kejia Chen, Lipeng He, Jian Lou, Dan Li, Zunlei Feng, Mingli Song, Jian Liu, Kui Ren, and Xiaohu Yang. 2025a. Activation approximations can incur safety vulnerabilities even in aligned llms: comprehensive analysis and defense. In *Proceedings of the 34th USENIX Conference on Security Symposium, SEC '25, USA*. USENIX Association.
- Ruiqi Zhang, Licong Lin, Yu Bai, and Song Mei. 2024a. [Negative preference optimization: From catastrophic collapse to effective unlearning](#). In *First Conference on Language Modeling*.
- Yiming Zhang, Nicholas Carlini, and Daphne Ippolito. 2024b. [Effective prompt extraction from language models](#). In *First Conference on Language Modeling*.
- Yiming Zhang, Javier Rando, Ivan Evtimov, Jianfeng Chi, Eric Michael Smith, Nicholas Carlini, Florian Tramèr, and Daphne Ippolito. 2025b. [Persistent pre-training poisoning of LLMs](#). In *The Thirteenth International Conference on Learning Representations*.
- Victor Zhong, Caiming Xiong, and Richard Socher. 2017. Seq2sql: Generating structured queries from natural language using reinforcement learning. *CoRR*, abs/1709.00103.
- Daniel Ziegler, Seraphina Nix, Lawrence Chan, Tim Bauman, Peter Schmidt-Nielsen, Tao Lin, Adam Scherlis, Noa Nabeshima, Benjamin Weinstein-Raun,

Daniel de Haas, Buck Shlegeris, and Nate Thomas. 2022. [Adversarial training for high-stakes reliability](#). In *Advances in Neural Information Processing Systems*.

Andy Zou, Long Phan, Justin Wang, Derek Duenas, Maxwell Lin, Maksym Andriushchenko, J Zico Kolter, Matt Fredrikson, and Dan Hendrycks. 2024. [Improving alignment and robustness with circuit breakers](#). In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.

Andy Zou, Zifan Wang, Nicholas Carlini, Milad Nasr, J. Zico Kolter, and Matt Fredrikson. 2023. [Universal and transferable adversarial attacks on aligned language models](#). *Preprint*, arXiv:2307.15043.

A Notations

Table 6: Frequently used notations and descriptions.

Notation	Description
\mathcal{F}, f_i, m	Set of features, a feature, # of features.
\mathcal{C}, \mathcal{C}	Client, Set of clients.
a_i	Adapter to lock (refuse) feature f_i .
π_θ	Language model with parameters θ .
\mathcal{R}	Responses generated by π_θ
l, L	A layer index, and the set of target layers.
D_{f_i}	Dataset corresponding to feature f_i .
(x_i, y_i)	A prompt-response pair from D_{f_i} .
D_{auth}, D_{unauth}	Datasets for utility and refusal training.
c_i, r_i	Chosen and rejected responses.
\mathcal{L}_{lock}	Total loss for adapter fine-tuning.
$\mathcal{L}_{utility}, \mathcal{L}_{robust}$	Losses for utility and robustness.
δ, ϵ	Perturbations and its L2-norm budget.
$\gamma(x, \delta)$	Function applying δ to activations.
ΔW^i	Weight update matrix for adapter a_i .
σ^i, σ_l	Adapter L2-norm; max norm over layer l .
τ	Scaling hyperparameter for clipping.
$Clip_l = \tau \sigma_l$	Norm clipping threshold.

B Definition of a Feature

We define a *feature* as a specific model capability required to perform a task (e.g., “Math”, “Code”). The granularity of feature definitions is determined by the service provider based on their product offering; features can range from coarse (e.g., “Math”) to fine-grained (e.g., “Algebra”, “Geometry”).

C Inference-Time Control Flow

At inference time (Figure 1, Step ③), the *access control module* identifies the set of features authorized for the requesting client based on their profile and subscription status. It then selects the adapters corresponding to all *unauthorized* features (Steps ④–⑤), merges them, and attaches the merged adapter to the frozen base model before inference (Step ⑥). The adapters are trained

on feature-specific datasets, causing the model to refuse queries that invoke unauthorized features while leaving authorized features unaffected.

D Other Applications of FLoTEs

Beyond pay-to-unlock schemes, FLoTEs have several promising potential applications. They can serve as an alignment alternative by locking dangerous features while preserving utility for legitimate tasks, providing a more granular and reversible approach than methods that modify base model weights. They can potentially enable more robust feature-level unlearning by locking access to specific features rather than attempting to erase them, making unauthorized knowledge more difficult to elicit even through adversarial prompting. They support staged feature releases, allowing providers to gradually roll out features to different user groups without maintaining multiple model versions. Finally, they enforce conditional compliance for sensitive features such as medical diagnoses or legal advice, keeping them locked by default and unlocking only for users or regions where appropriate regulatory approvals have been obtained.

E Implementation Details

Adapter Training. For LOCKET, we train LoRA adapters with a rank of 64, alpha of 64, and a dropout of 0.1. We use RSLoRA (Kalajdziewski, 2023) for improved performance. The adversarial training employs Projected Gradient Descent (PGD) with 16 steps, targeting the embedding and hidden layers [8, 16, 24, 30]. We train for 100 total steps with a batch size of 2. For the baseline, we follow the SFT configurations of prior work (Greenblatt et al., 2024b; Tang et al., 2024), using a validation set comprising 20% no-password prompts and 80% incorrect-password prompts to ensure robust refusal learning. For tuning the scaling threshold τ in our adapter merging strategy, we use a random sample of 100 examples from each test set.

Dataset Composition. Train and test splits compositions of datasets can be found in Table 7. We use public open-sourced datasets and models.

Adapter Merging. We run hyperparameter tuning experiments to select optimal τ values for each feature combination. The following are the final τ values we use in merging the adapters for effectiveness and robustness evaluation, for DeepSeek-7B-Math, we have: **M** (0.9), **Q** (0.7), **S** (0.5), **U** (0.7),

Table 7: Datasets

	Dataset	Train	Test
Training (utility)	UltraChat	165,298	-
Feature (specific)	SQL Create Context	62,861	15,716
	MATH	7,500	5,000
	Samsum	819	14,732
Feature (general)	MMLU	99,842	14,042

M + Q (0.85), **M + S** (0.85), **M + U** (0.85), **Q + S** (0.6), **Q + U** (0.8), **M + Q + S** (0.75), **M + Q + U** (0.9), **M + S + U** (0.85), **Q + S + U** (0.75), **M + Q + S + U** (0.75); for DeepSeek-7B-Coder, we have: **Q** (0.45); for Llama-3-8B-Instruct, we have **M** (0.7), **Q** (0.6), **S** (0.9), **U** (0.8), **M + Q** (0.7), **M + S** (0.8), **M + U** (0.7), **Q + S** (0.8), **Q + U** (0.8), **M + Q + S** (0.8), **M + Q + U** (0.7), **M + S + U** (0.75), **Q + S + U** (0.8), **M + Q + S + U** (0.75).

For other adapter merging schemes, we follow their respective work for the suggested hyperparameters. Specifically, for CAT (Prabhakar et al., 2025) and Linear/Task Arithmetic (Ilharco et al., 2023), we apply equal weights to each adapter during merging; for TIES (Yadav et al., 2023), we use density = 0.5.

Jailbreak Attacks. We use a suite of state-of-the-art jailbreaking attacks to generate adversarial prompts, and evaluate robustness: (i) *Many-shot Jailbreaking* (Anil et al., 2024): For math, we sample 2, 4, and 8 in-context examples from problems below level-2 difficulty, and report the highest ASR achieved. For other tasks, we use randomly sampled demos. (ii) *GCG* (Zou et al., 2023): We run 125 iterations to optimize for the target string "Sure, here's how to solve this problem," using a batch size of 64 and replacing one character at a time. (iii) *TAP* (Mehrotra et al., 2024): We generate adversarial prompts using *lmsys/vicuna-13b-v1.5-16k* as the attack model, with a branching factor of 4, width of 10, depth of 5, and the ground truth as the target. (iv) *AutoDAN-Turbo* (Liu et al., 2025): We run a single warm-up iteration with a size of 50 for 150 epochs, followed by one lifelong iteration. All attacks use 1,000 random samples from each feature dataset, and generations are performed with a temperature of zero for deterministic outputs.

Computational Cost. Experiments are done using 8 * NVIDIA A100 40GB GPUs which consume in total round 6000 GPU hours.

F Variance Across Runs

Throughout our main experiments we set the inference temperature to zero (greedy decoding) to eliminate sampling variance, following prior work on password-locking (Greenblatt et al., 2024b; Hofstätter et al., 2025) and jailbreak evaluation (Mehrotra et al., 2024; Liu et al., 2025). To validate that our results hold under probabilistic sampling, we conducted 10 independent inference runs on DeepSeek-7B-Math with **U** locked, using temperature = 1. Effectiveness remained near-perfect ($\leq 1\%$ utility on unauthorized requests, $\pm 1\%$), and utility on the unlocked features was preserved (\approx baseline $\pm 2\%$), confirming that the reported results are stable and not artefacts of greedy decoding.

G Adapter Merging

Table 8 compares LOCKET merging against standard adapter merging techniques available in the PEFT library (Mangrulkar et al., 2022). Most existing methods either over-refuse (refusing even unlocked features, ACC = 0) or fail to lock effectively (high ALA on locked features). LOCKET merging, however, can match baseline accuracy on unlocked features while maintaining perfect locking (ALA = 0).

H Additional Scalability Results

We provide additional comparisons between LOCKET and prior work across different models and feature combinations. The first two tables (Table 10, 11) compare LOCKET with password-locking (PWD) on Llama-3-8B-Instruct and DeepSeek-7B-Coder, showing that LOCKET consistently outperforms PWD in preserving utility while maintaining effective locking. The third table (Table 12) compares LOCKET against PWD with Circuit Breaking (PWD + CB) (Hofstätter et al., 2025) on DeepSeek-7B-Math; while PWD+CB improves robustness, it fails to maintain effectiveness when locking multiple features and still degrades utility on **U**, whereas LOCKET achieves perfect effectiveness across all configurations. The fourth table (Table 13) demonstrates that LOCKET scales to larger models (Llama-3-70B-Instruct, 4-bit quantized), remaining effective and utility-preserving with only minor **U** degradation due to interference with **M**.

I Scalability beyond Four Features

We evaluate LOCKET's scalability beyond four features by progressively locking up to eight features

Table 8: **Comparison of LOCKET merging with other commonly used merging techniques.** LOCKET merging mitigates over-refusal while maintaining effective locking. ACC: Unlocked Feature Accuracy (higher is better), ALA: Averaged Locked Feature Accuracy (lower is better). All results are obtained on DeepSeek-7B-Math. LOCKET matches baseline accuracy on unlocked features while keeping locked features effectively refused (ALA=0). Other methods either **over-refuse** (ACC=0 on unlocked features) or **fail to lock effectively** (high ALA).

Unlocked Feature ↓		Baseline	LOCKET Merging	SVD	TIES SVD	DARE TIES	DARE Linear	DARE TIES SVD	DARE Linear SVD	Magnitude Prune	Magnitude Prune SVD
Math	ACC ↑	0.40	0.45	0.00	0.00	0.40	0.00	0.00	0.00	0.00	0.00
	ALA ↓	0.00	0.00	0.00	0.00	0.37	0.02	0.00	0.02	0.01	0.01
SQL	ACC ↑	0.93	0.94	0.00	0.00	0.93	0.00	0.00	0.00	0.00	0.00
	ALA ↓	0.00	0.00	0.00	0.00	0.29	0.02	0.00	0.00	0.01	0.01
Summarize	ACC ↑	0.23	0.24	0.00	0.00	0.22	0.00	0.00	0.00	0.06	0.00
	ALA ↓	0.00	0.00	0.00	0.00	0.48	0.03	0.00	0.02	0.00	0.00
MMLU	ACC ↑	0.53	0.53	0.00	0.00	0.54	0.00	0.00	0.00	0.00	0.00
	ALA ↓	0.00	0.00	0.00	0.00	0.45	0.02	0.00	0.02	0.00	0.00
			Matches baseline	Over-refusing	Over-refusing	Ineffective locking	Over-refusing	Over-refusing	Over-refusing	Over-refusing	Over-refusing

Table 9: **LOCKET remains reasonably scalable beyond four features:** “Baseline” is the original model behaviour without FLoTE. For effectiveness (**R1**), we use **blue** to indicate complete locking, ineffective locking → **orange**. For utility (**R2**) (of unlocked features) **green** → matches/outperforms baseline, **yellow** → within $\pm 5\%$ of baseline, **red** → worse than baseline. **L**, **Y**, **P** and **O** data are taken from the subsets of MMLU, which are largely within the same distribution, potentially causing larger *feature interference* and lead to less desirable results (as discussed in **Limitations**).

Locked Feature →	Baseline	M	M + Q	M + Q + S	M + Q + S + L	M + Q + S + L + H	M + Q + S + L + H + Y	M + Q + S + L + H + Y + P	M + Q + S + L + H + Y + P + O
Math (M)	0.40	0.00	0.00	0.00	0.00	0.00	0.00	0.02	0.00
SQL (Q)	0.93	0.95	0.00	0.00	0.00	0.00	0.00	0.00	0.00
Summarize (S)	0.23	0.23	0.24	0.00	0.00	0.00	0.00	0.00	0.00
Law (L)	0.35	0.35	0.35	0.35	0.00	0.00	0.09	0.10	0.00
History (H)	0.54	0.56	0.52	0.52	0.46	0.02	0.12	0.09	0.00
Psychology (Y)	0.64	0.63	0.63	0.62	0.56	0.54	0.08	0.11	0.00
Politics (P)	0.44	0.44	0.43	0.43	0.39	0.45	0.38	0.06	0.00
Philosophy (O)	0.49	0.45	0.45	0.44	0.43	0.40	0.44	0.34	0.00

Table 10: **Comparison of LOCKET with prior work:** Scalability w.r.t. **R1** and **R2** of LOCKET with prior work (“PWD”) (Greenblatt et al., 2024a; Tang et al., 2024) locking Llama-3-8B-Instruct. Color coding for scalability, are same as Table 5.

Locked Feat. →	M		M + Q		M + Q + S	
	PWD	LOCKET	PWD	LOCKET	PWD	LOCKET
Math (M)	0.00	0.00	0.00	0.00	0.00	0.00
SQL (Q)	0.01	0.92	0.00	0.00	0.00	0.00
Summarize (S)	0.26	0.34	0.41	0.34	0.51	0.00
MMLU (U)	0.06	0.64	0.05	0.73	0.05	0.72

Table 11: **Comparison of LOCKET with prior work:** Scalability w.r.t. **R1** and **R2** of LOCKET with prior work (“PWD”) (Greenblatt et al., 2024a; Tang et al., 2024) locking DeepSeek-7B-Coder. Color coding for scalability, are same as Table 5.

Locked Feat. →	M		M + Q		M + Q + S	
	PWD	LOCKET	PWD	LOCKET	PWD	LOCKET
SQL (Q)	0.01	0.96	0.00	0.00	0.00	0.00

on DeepSeek-7B-Math (Table 9). The additional features (Law, History, Psychology, Politics, Phi-

losophy) are drawn from MMLU subcategories, which share similar distributions and thus exhibit greater feature interference. Despite this challenging setting, LOCKET maintains perfect effectiveness

Table 12: **Comparison of LOCKET with prior work:** Scalability w.r.t. **R1** and **R2** of LOCKET with prior work (“PWD + CB”) (Hofstätter et al., 2025) locking DeepSeek-7B-Math via PWD and Circuit Breaking. Color coding for scalability, are same as Table 5.

Locked Feat. →	M		M + Q		M + Q + S	
	P.+CB	LOCKET	P.+CB	LOCKET	P.+CB	LOCKET
Math (M)	0.00	0.00	0.38	0.00	0.20	0.00
SQL (Q)	0.93	0.95	0.47	0.00	0.05	0.00
Summarize (S)	0.23	0.23	0.28	0.24	0.50	0.00
MMLU (U)	0.38	0.51	0.43	0.53	0.39	0.53

Table 13: **Applying LOCKET to a larger-scale model:** Scalability w.r.t. **R1** and **R2** of LOCKET locking Llama-3-70B-Instruct (4-bit quantized). Color coding for scalability, are same as Table 5.

Locked Feat. →	None	M	M + Q	M + Q + S
Math (M)	0.49	0.00	0.00	0.00
SQL (Q)	0.98	0.98	0.00	0.00
Summarize (S)	0.37	0.38	0.37	0.00
MMLU (U)	0.81	0.79	0.78	0.75

in most configurations, with only minor ineffective locking (orange cells) appearing when six or more features are locked simultaneously. Utility degradation increases with the number of locked features, particularly among MMLU subcategories due to their distributional overlap. These results align with prior work (Lee et al., 2025) showing that adapter merging can scale to 8 adapters with roughly 15% degradation, confirming that LOCKET remains reasonably scalable beyond the four features evaluated in the main text.