Permissioned LLMs: Enforcing Access Control in Large Language Models

Bargav Jayaraman

Oracle Labs

bargav.jayaraman@oracle.com

Virendra J. Marathe

Oracle Labs

virendra.marathe@oracle.com

Hamid Mozaffari

Oracle Labs

hamid.mozaffari@oracle.com

William F. Shen

University of Cambridge fs604@cam.ac.uk

Krishnaram Kenthapadi

Oracle Health

krishnaram.kenthapadi@oracle.com

Abstract

In enterprise settings, organizational data is segregated, siloed and carefully protected by elaborate access control frameworks. These access control structures can completely break down if an LLM fine-tuned on the siloed data serves requests, for downstream tasks, from individuals with disparate access privileges. We propose Permissioned LLMs (PermLLM), a new class of LLMs that superimpose the organizational data access control structures on query responses they generate. We formalize abstractions underpinning the means to determine whether access control enforcement happens correctly over LLM query responses. Our formalism introduces the notion of a relevant response that can be used to prove whether a PermLLM mechanism has been implemented correctly. We also introduce a novel metric, called access advantage, to empirically evaluate the efficacy of a PermLLM mechanism. We introduce three novel PermLLM mechanisms that build on Parameter Efficient Fine-Tuning to achieve the desired access control. We furthermore present two instantiations of access advantage-(i) Domain Distinguishability Index (DDI) based on Membership Inference Attacks, and (ii) Utility Gap Index (UGI) based on LLM utility evaluation. We demonstrate the efficacy of our PermLLM mechanisms through extensive experiments on five public datasets (GPQA, RCV1, SimpleQA, WMDP, and PubMedQA), in addition to evaluating the validity of DDI and UGI metrics themselves for quantifying access control in LLMs.

1 Introduction

Large Language Models (LLMs), due to their unprecedented natural language processing capabilities, are being adopted in a vast range of applications across the entire computing industry [21, 48]. The day may not be too far off when LLMs become the primary interface to a large swath of computing and information extraction tasks. In this paper, we focus on enterprise settings where LLMs are used to perform a wide variety of computing tasks using organization-wide data. Using LLMs that have a wide purview over organizational data brings massive troves of information and utility, including the ability to combine learnings from disparate information silos of the organization, to the finger tips of individuals in the organization. However, making all the learnings from organizational data available to any individual who can query the LLM becomes a critical security challenge: Organizations have

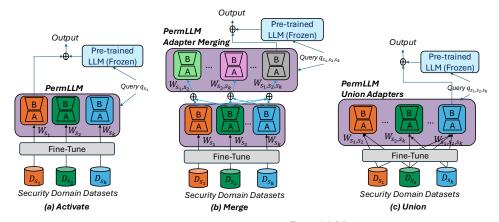


Figure 1: We propose three types of Permissioned LLM (PermLLM) mechanisms. (a) *Activate*: that has one-to-one mapping between the security domains and PEFT adapters. When a user queries the model, the mechanism activates the relevant adapter(s). (b) *Merge*: merges subsets of relevant PEFT adapters to serve the users that have access to multiple security domains. (c) *Union*: trains adapters on the unions of various security domains, and at the inference phase the relevant PEFT adapter is activated to serve a user query that requires access to multiple security domains.

access control structures and hierarchies that tightly control information flow to and from individuals within them. Information access via LLMs, if not carefully controlled, risks undermining the existing access control structures and hierarchies.

As an example, consider government agencies using LLMs for a multitude of tasks. The data in government agencies is typically segregated in multiple "clearance levels" and users can access just the data they have access privileges for [30]. Any other agency data is inaccessible to the users. An LLM trained on this agency-wide data can leak privileged information to unauthorized users, thus breaking the agency's access control framework that works on the raw data. Another example is that of role-based access control [10, 11]: Consider a health clinic setting, where individuals performing different "roles" (doctors, nurses, technicians, administrative staff, patients, etc.) interact with an LLM to perform many tasks. The roles of the users determine what part of the clinic-wide data they should have access to. An LLM trained on the clinic-wide data can be easily tricked into leaking information to unauthorized users.

Research proposals to build system prompts that instruct an LLM to control what information is generated in the output can help curb some leakage of sensitive information to unauthorized users [8, 25]. However, they do not achieve absolute security, and clever jailbreaking prompts [40, 34, 26, 27] can be used to overrule these system prompts. A recent work proposes tagging LLM queries with encrypted access credentials to authenticate users and block unauthorized queries [7]. This is a good start, but it lacks the flexibility needed to enable access to disparate learnings from the LLM for different users based on their access credentials. We discuss access control problems and solutions for agentic systems and Retrieval Augmented Generation (RAG) systems [23] in § 6.

This paper focuses on the access control problem for LLMs when they are tuned on data coming from a multitude of data silos. The challenge here is to *guarantee* that users who do not have access to specific data silos cannot retrieve information from those silos by sending carefully crafted queries to the LLMs tuned on data from those silos. A recent work [12] took an initial step in this direction, but lacks the formal framework to evaluate the access control. Moreover they only explore one type of mechanism. As a security problem, access control is *absolute*—you either achieve access control or not, hence probabilistic solutions (e.g. Differential Privacy [9]) are not satisfactory.

Contributions. In this paper, we comprehensively study the problem of access control in LLM fine-tuning. More specifically: (i) We formalize the notion of access control mechanism in LLMs in terms of the *relevance* of responses generated by an LLM to the raw data the users have access to. We use the notion of *security domains* in our formalism. Our formalism of response relevance can be used to prove correctness of access control mechanisms. We also propose a novel metric called *access advantage* that helps us empirically quantify the effectiveness of an access control mechanism

on LLMs (§ 2). (ii) We present three new PermLLM fine-tuning mechanisms (see Figure 1), based on Parameter Efficient Fine-Tuning (PEFT) [17, 42] (§ 3). (iii) We introduce two novel instances of our access advantage metric, *Domain Distinguishability Index (DDI)* and *Utility Gap Index (UGI)*, as tools to audit access control enforcement via an adversarial gaming setting (§ 4). (iv) We empirically evaluate our access control mechanisms on two LLMs (Mistral-0.1-7B and Llama-3.1-8B) using five different data sets: GPQA [33], RCV1 [22], SimpleQA [41], WMDP [24], and PubMedQA [20] (§ 5). Our evaluation shows the effectiveness of our access advantage metrics in assessing whether a proposed access control mechanism for LLMs is enforcing data protection correctly.

2 Formalizing Access Control in LLMs

2.1 Basic Setup and Notation

We define a *security domain* (henceforth called "domain" for brevity) as a collection of data records that require identical credentials for access (e.g. files with the same group in their access control lists). We consider settings where pretrained LLMs (such as Llama and Mistral models) are fine-tuned over data from different domains with an added constraint – responses to inference time queries will be generated from learnings on data coming from just the domains the user has access to. This added constraint is enforced via access control mechanisms that govern how the LLM uses data from different domains.

Consider a universe of n different domains $\mathbb{S} = \bigcup_{i=1}^n \{s_i\}$, and a training data set consisting of data from these domains $D = \bigcup_{i=1}^n D_{s_i} \sim \mathcal{D}_{s_i}$ (here D_{s_i} is a data set sampled from data distribution \mathcal{D}_{s_i} of domain s_i). Let f_D be the LLM tuned using data set D. Let W be the set of f_D 's parameters. Model fine-tuning changes values of a subset of W. We say that a domain s_i affects a subset of parameters $W_{s_i} \subseteq W$ if data from D_{s_i} is used to change parameters W_{s_i} during model fine-tuning (unless stated otherwise, the terms "affect" and "affected" mean this relation between s_i and W_{s_i} in the rest of the paper). We define \mathcal{M} as an access control mechanism that dictates the mapping of domain s_i to parameters W_{s_i} via the affects relation. We say that a LLM fine-tuned using data set D is permissioned (PermLLM), denoted as $f_D^{\mathcal{M}}$, if it uses the access control mechanism \mathcal{M} to map its parameters W to a multitude of domains from \mathbb{S} , where each domain s_i affects parameters $W_{s_i} \subseteq W$. Operationally, during fine-tuning, \mathcal{M} specifies which set of model parameters W_{s_i} are tuned for a given domain s_i (see § 3 for more details). By the same token, during inference, \mathcal{M} can specify which set of model parameters should be used to answer a query based on the user's access credentials.

We assume a setting where the PermLLM $f_D^{\mathcal{M}}$ resides in an enclosing system \mathcal{S} that authenticates users who send queries to $f_D^{\mathcal{M}}$. \mathcal{S} determines the user u's access credentials $cred_u$ and calls authenticate ($cred_u$) that takes user credentials $cred_u$ and maps them to a subset of domains S_u that u can access. S_u is maintained by \mathcal{S} and is never exposed to user u. This process ensures u cannot arbitrarily change S_u . Each of user u's subsequent query q to $f_D^{\mathcal{M}}$ is annotated with S_u by \mathcal{S} . \mathcal{M} determines the model parameters W_{S_u} used to generate a response r_{S_u} to q, where $W_{S_u} = \bigcup_{s \in S_u} W_s$.

2.2 Definitions

Definition 2.1 (Relevant Response). Given a PermLLM $f_D^{\mathcal{M}}$, a query q from user u, and the set S_u of domains u has access to, let $r = f_D^{\mathcal{M}}(q)$ be the response of $f_D^{\mathcal{M}}$ to query q. Response r is said to be relevant to S_u (i.e., $r = r_{S_u}$) if $f_D^{\mathcal{M}}$ uses parameters W_{S_u} (in addition to any domain-agnostic model parameters) to generate r.

We say that an access control mechanism \mathcal{M} is correctly enforced on PermLLM $f_D^{\mathcal{M}}$ iff every response r generated for every user u's query q is relevant to S_u .

The above definition of relevant response helps us formally determine if a proposed access control mechanism \mathcal{M} is algorithmically correct. We however require an empirically quantifiable metric to determine if the implementation (and the algorithm by extension) of \mathcal{M} is correct. This is particularly important for auditing. To that end, we propose a new metric called *response relevance score*, $relv(f_D^{\mathcal{M}}(q), S_u)$, which quantifies the information gained on data in the domain set S_u by observing

responses to queries generated using model parameters W_{S_u} affected by domains of S_u . relv is expected to be higher when $q \sim \mathcal{D}_{S_u}$ (i.e., q is related to domain set S_u), compared to when $q \not\sim \mathcal{D}_{S_u}$.

We restrict the domain of relv to the real number interval [0,1], where 1 is the best expected score for relevance. relv itself can be represented by another empirical metric such as prediction accuracy, or logits for the expected response. However, given that LLMs (and ML models in general) are generalization engines, in practice we expect relv to be less than 1. This restriction can be effectively addressed by measuring relv for domains that the user has access to and comparing it to relv for domains that the user does not have access to. We call this the *access advantage*.

Definition 2.2 (Access Advantage). Given PermLLM $f_D^{\mathcal{M}}$ trained over data set D consisting of data from domains $\mathbb{S} = \bigcup_{i=1}^n \{s_i\}$, with access control mechanism \mathcal{M} , a subset of domains $S_u \subseteq \mathbb{S}$, $f_D^{\mathcal{M}}$ achieves α -access advantage w.r.t. S_u if:

$$\mathbb{E}_{q \sim \mathcal{D}_{S_u}, S_v \subseteq \mathbb{S}; S_u \cap S_v = \phi} \left[relv(f_D^{\mathcal{M}}(q), S_u) \ominus relv(f_D^{\mathcal{M}}(q), S_v) \right] \ge \alpha$$

where relv() is the response relevance score on the corresponding domain subset $(S_u \text{ or } S_v)$, \ominus is a "difference" operator specific to the access control assessment method (e.g., subtraction), and α is an advantage threshold that lies in the range [0,1].

The access advantage metric relies on the assumption that $f_D^{\mathcal{M}}$ performs significantly better on domains user u has access to compared to domains u does not have access to. In other words, $f_D^{\mathcal{M}}$ should have explicit segregation between the different domains, as dictated by \mathcal{M} . The existing approaches to model fine-tuning fail to achieve this goal as the model is traditionally trained on all the domains without any built-in domain segregation mechanism. To the best of our knowledge, no prior work on LLM and privacy formally tackles this problem of access control through explicit domain segregation. We next propose novel mechanisms to achieve domain segregation in § 3 and propose empirical metrics to evaluate the access control protocols in § 4.

We believe access advantage is a critical metric for auditors to determine if an access control mechanism is truly achieving the segregation of domains as intended. Hence it is in the auditor's best interest to ensure that $S_u \cap S_v = \phi$. Access advantage can diminish significantly when $S_u \cap S_v \neq \phi$, leading to incorrect conclusions about the efficacy of the access control mechanism.

To the best of our knowledge, prior works on retrieval augmented generation (RAG) based LLM deployments do not explicitly tackle the problem of measuring effectiveness of access control mechanisms formally or empirically. Our formalism of relevant response and access advantage extends to RAG systems as well, closing that gap in formalism and empirical evaluation of access control protocols. While a thorough evaluation of access control for RAG based systems is outside the scope of this paper, a more detailed analysis of conditions for formal correctness of access control in RAG systems appears in Appendix A.

2.3 Auditing Access Control

To evaluate the access control mechanisms, we consider a classic adversarial game between the system $\mathcal S$ enclosing the model $f_D^\mathcal M$ and the auditor $\mathcal A$. We give $\mathcal A$ the ability to choose domain access by emulating an end user, send arbitrary queries to the model via $\mathcal S$ and observe the responses. $\mathcal A$ can replay the game multiple times as different users to conclude if the access control is correctly implemented.

Audit Game. The formal game between auditor \mathcal{A} and system \mathcal{S} is as follows:

- 1. Auditor A chooses domain set S_u and emulates user u. A sends user credentials $cred_u$ and query $q \sim \mathcal{D}_{S_u}$ to system S.
- 2. S verifies the user credential $cred_u$ and sends back the model response $f_D^{\mathcal{M}}(q)$ to \mathcal{A} .
- 3. \mathcal{A} chooses domain set S_v such that $S_v \cap S_u = \phi$ and emulates user v. \mathcal{A} sends user credentials $cred_v$ and the same query $q \sim \mathcal{D}_{S_u}$ to \mathcal{S} .
- 4. S verifies the user credential $cred_v$ and sends back the model response $f_D^{\mathcal{M}}(q)$ to \mathcal{A} .
- 5. A concludes the access control mechanism is correctly implemented if the access advantage $|relv(f_D^{\mathcal{M}}(q), S_u) \ominus relv(f_D^{\mathcal{M}}(q), S_v)| \ge \alpha$.

Note that the auditor A has superuser privileges to choose arbitrary domain access unlike an ordinary user. This is by design to allow the auditor to evaluate the correctness of the claimed access control

while still following the protocol of querying the model as a benign user. Detailed instantiations of this adversarial game for different suites of access advantage metrics are discussed in Appendix C.

3 Permissioned LLM Mechanisms

We rely on Parameter Efficient Fine-Tuning (PEFT) [17, 42] to obtain model parameter segregation for domains. We focus on the LoRA PEFT adapter [17], however our mechanisms seamlessly apply to other types of adapters [16, 42]. The three mechanisms we describe ensure that domain data is steered to train select LoRA adapters. Each domain has a unique identifier (domain Id). Our access control mechanism builds a map between domains and LoRA adapters within the PermLLM's metadata. The map is used to steer all examples from a domain to the corresponding adapter/s for training. This map is also used to steer queries to the correct LoRA adapters at inference time.

One LoRA per Security Domain For our base mechanism called *Activate*, we assume that users have access to at most one domain. Figure 1(a) depicts our base mechanism that performs a simple 1-1 mapping between domains and LoRA adapters. We assume that the number of domains is known beforehand, and use that knowledge to instantiate corresponding number of LoRA adapters. During training, each minibatch is sampled from one domain, and the domain's Id is used to select the LoRA adapter to train. At inference time, a user's query is annotated with the domain Id the user has access to. This domain Id is used to *activate* the LoRA adapter for the corresponding domain.

Merging LoRA Adapters for Security Domain Groups In many application settings, users have access to data from multiple domains. For queries coming from such users, our *Activate* enables all corresponding LoRA adapters, whose activations are averaged at inference time. We however found that activations from different LoRA adapters tend to disruptively interfere with each other resulting in catastrophic performance degradation beyond two domains. We leave further refinement of activation space steering [35, 45] to future work. In our second mechanism, *Merge* (Figure 1(b)), we adopt the LoRA adapter merging strategy for users with access to multiple domains [38, 43, 46, 50]. We experimented with several LoRA merging algorithms including TIES [43] and DARE [46], but eventually favored the SVD approach [38] because of its better performance and stability in the context of LoRA merging. We assume that the combination of domains that users may have access to are known beforehand. Thus, after training LoRA adapters for individual domains, we can merge them for those exact domain combinations. Correspondingly, our domain-LoRA adapter map is updated with the domain IDs and the merged LoRA adapters. These new mappings are used at inference time to activate the correct merged LoRA adapters. We found that adapter merging is more robust to cross-adapter interference than activation merging.

Training a LoRA per Combination of Security Domains Although *Merge* is better than activation space merging of multiple LoRA adapters, we observed that it also leads to model performance degradation with increasing number of merged adapters. As a result, we explored another simple alternative, *Union*, which *trains* a LoRA adapter on data from each unique combination of domains users have access to. *Union* indeed delivers the best performance in all our mechanisms. However, it comes at the cost of significantly greater tuning time compute – a domain can occur in numerous combinations of domains (e.g. in Figure 1(c), data D_{s_2} gets used in the training set of all three LoRA adapters). Furthermore, data sets containing large number of domains presents the added challenge of an exponential blow up in domain combinations (at most 2^n). However, we believe the number of combinations present in a real-world setting will be much smaller than that upper bound.

The careful mapping of domains (or groups of domains) to the correct LoRA adapters, and steering of training examples from domains to the corresponding LoRA adapters ensures precise parameter segregation for domains. Our assumption that users cannot tamper with their access credentials at inference time further aids the PermLLM's enclosing system to determine the correct set of domains corresponding to a query. The query steering that happens through the PermLLM using domain IDs *guarantees* that all responses generated by the PermLLM are *relevant* to the user's domains. Furthermore, the responses are not generated using LoRA adapters that were trained using data from domains that the user does *not* have access to. Response relevance for all responses implies correctness of our PermLLM access control mechanisms. Our proof appears in Appendix B.

4 Auditing Access Control in Permissioned LLM Mechanisms

We now introduce two novel instantiations of our *access advantage* metric (Definition 2.2)—Domain Distinguishability Index (DDI) and Utility Gap Index (UGI)—that quantify access control efficacy independently of any particular system design. We show how these metrics fit into the framework for empirically assessing access control mechanisms in PermLLMs through adversarial audit games in Appendix C. These metrics are in [0,1] range with higher values denoting better access control.

4.1 Metric 1: Domain Distinguishability Index (DDI)

In traditional privacy evaluations, membership inference attacks (MIAs) leverage a sampled member data set (from the target model's training set) and a sampled non-member data set to assess privacy leakage [18, 37]: the more accurately an adversary separates and classifies samples as members or non-members, the higher the privacy risk. Analogously, we adopt this MIA framework for access control assessment to distinguish security domains. Specifically, for any security domain set S_i , the auditor holds samples from S_i 's training data (member set) and samples from S_j 's training data (non-member set), where $S_j \cap S_i = \phi$. The auditor then evaluates how successfully it can distinguish the member set from the non-member set when the PermLLM is activated for S_i . This evaluation occurs for all security domains, giving us an aggregate access advantage, which we call Domain Distinguishability Index (DDI).

Definition 4.1 (Domain Distinguishability Index (DDI)). For a domain universe \mathbb{S} consisting of n security domains, let $f_D^{\mathcal{M}}$ denote the PermLLM trained on data D from all security domains with access control mechanism \mathcal{M} . For each ordered pair of domain sets $(S_i \subseteq \mathbb{S}, S_j \subseteq \mathbb{S})$ with no overlap $(i.e., S_i \cap S_j = \phi)$, let $O^{(S_i, S_j)} = O(f_D^{\mathcal{M}}(q)|S_i, f_D^{\mathcal{M}}(q)|S_j)$; $\forall q \sim \mathcal{D}_{S_i}$ be the output of a membership inference oracle O. For a given membership inference metric $m(\cdot)$, the DDI is defined as: $DDI(m) = \mathbb{E}_{S_i \subseteq \mathbb{S}, S_j \subseteq \mathbb{S}} \big[m(O^{(S_i, S_j)}) \big]$, where \mathbb{E} is the expectation over all domain sets.

We also report the standard deviation of $m(O^{(S_i,S_j)})$ across all domain set pairs to capture variability. By 2.2, DDI can be viewed as an access advantage metric, where the response relevance score relv for S_i on query q, $relv(f_D^{\mathcal{M}}(q), S_i)$, is a binary value on whether the membership inference oracle O's output is above a membership threshold. The difference operator \bigcirc is the MIA method specific composition of response relevance for all the samples in the member and non-member sets.

We use AUC-ROC and TPR@(low)FPR, as instantiations of DDI, where higher scores indicate stronger enforcement, as S_i -specific responses become more distinguishable. See Appendices F.1 and F.2 for details on MIA evaluation metrics and an overview of existing MIAs against LLMs.

A higher DDI indicates more robust separation between security domains. In our evaluations, we employ state-of-the-art MIAs for LLMs, including Loss [44], Zlib [5], Mink [36], Mink++ [47], Reference [5] attacks.

4.2 Metric 2: Utility Gap Index (UGI)

The UGI metric measures the drop in model utility on the target domain's data when a different domain's adapter is activated in PermLLM instead of the target domain.

Definition 4.2 (Utility Gap Index (UGI)). Let $U(\cdot)$ be a chosen utility metric and for a domain set pair $(S_i \subseteq \mathbb{S}, S_j \subseteq \mathbb{S})$ with no overlap (i.e., $S_i \cap S_j = \phi$), Utility $\operatorname{Gap}^{(S_i, S_j)}(U) = |U(f_D^{\mathcal{M}}(q)|S_i) - U(f_D^{\mathcal{M}}(q)|S_j)|$; $\forall q \sim \mathcal{D}_{S_i}$. The UGI aggregates utility gaps across all ordered domain set pairs: $\Delta_U = \mathbb{E}_{S_i \subseteq \mathbb{S}, S_j \subseteq \mathbb{S}}[\text{Utility}\operatorname{Gap}^{(S_i, S_j)}(U)]$, where \mathbb{E} is the expectation over all domain sets.

By 2.2, UGI is also an instantiation of the access advantage metric in which the relevance score for security domain set S_i on query q is the utility value itself, $relv(f_D^{\mathcal{M}}(q), S_i) = U(f_D^{\mathcal{M}}(q)|S_i)$, and the operator \odot computes the absolute difference of those relevance scores across the sampled queries.

A larger UGI indicates that enforced access controls yield more pronounced—and thus more easily perceivable—differences in response quality between security domains. As with DDI, we also report the standard deviation across pairs to characterize variability. We evaluate the utility gaps w.r.t. Bleurt Score (Δ_{bluert}), Bert F1-Score (Δ_{bert}), Sacrebleu Score (Δ_{bleu}) and Verbatim Accuracy (Δ_{acc}) for our UGI metrics in § 5. More details on these metrics can be found in Appendix § D.3.

Table 1: Data Set Details.

	WMDP	GPQA	SimpleQA	RCV1	PubMedQA
Data Set Size (Train / Test)	2936 / 732	360 / 88	4089 / 1018	45622 / 22811	200000 / 11269
Number of Security Domains	3	3	10	4	10

5 Experimental Evaluation

For our experiments, we fine-tune Llama-3.1-8B and Mistral-0.1-7B pretrained models on five datasets covering multiple distinct security domains (henceforth called *domains*), where we fine-tune a separate LoRA adapter for each domain. Details about the model hyperparameters can be found in Appendix § D.1. The data sets we use in our experiments are WMDP [24], GPQA [33], SimpleQA [41], RCV1 [22], and PubMedQA [20]. Table 1 shows the brief data set details. More details on the data sets and generalization gaps can be found in Appendix § D.2. Appendix § D.3 discusses the utility evaluation of all our models.

5.1 Evaluating Access Control

Our approach achieves comparable model utility to existing approaches of fine-tuning (see discussion in § D.3), in addition to providing access control. Here we will empirically evaluate the effectiveness of our access control using a suite of metrics. We will first consider the case where the user has access to only one domain and report the access control results in § 5.1.1. Next in § 5.1.2, we will consider the case where the user has access to multiple domains. For comparison, we also include an evaluation of a prompt-based access control baseline in Appendix E but find it to be ineffective.

5.1.1 Single Active Domain

In Section 4, we proposed an *adversarial audit framework* for empirically assessing access control in PermLLMs. We introduced two concrete instantiations of the general *access advantage* metric: the Domain Distinguishability Index (DDI) and the Utility Gap Index (UGI) Δ_U . Although § 3 gives formal guarantees—each response is computed solely from domains the user is authorized to access—we *measure* access control enforcement strength with DDI and UGI (Δ_U) to confirm that the guarantees hold in practice, which is necessary to verify correctness of *implementations*.

Theoretically, Δ_U may reach 1.0, but empirically we observe much smaller—yet substantial—access advantage gaps for four of the data sets (Figure 2). These gaps are significantly impacted by domain distributions and the strictness of the scoring metric. For example, SimpleQA exhibits the largest UGIs (up to $\Delta_{blue} \approx 0.50$ and $\Delta_{acc} \approx 0.50$) because it has the highest number of distinct domains (10 in total). Moreover, we observe that Δ_{bleu} and Δ_{acc} have the largest values as these metrics look for verbatim pattern matches, thus requiring the model to memorize the nuances in the target domain. On the other hand, Δ_{bleurt} and Δ_{bert} look for approximate similarities, and hence are impacted by the similarities across the domains. This suggests that the verbatim matching metrics, Δ_{bleu} and Δ_{acc} , are better model utility metrics for measuring access advantage compared to the similarity based metrics Δ_{bleurt} and Δ_{bert} . For large data sets like RCV1, all the metrics achieve similar values as the model begins to generalize more. While these values are not close to 1, they still provide credence to the fact that the domains are different and our access control protocol works as expected due to the utility gaps. The access advantage threshold α is dependent on the type of utility metric: verbatim matching metrics Δ_{bleu} and Δ_{acc} have higher threshold than similarity based metrics Δ_{bleurt} and Δ_{bert} . For Δ_{acc} metric, $\alpha > 0.2$ is sufficient to infer that access control is happening correctly. PubMedQA is an exception where Δ_U values are close to zero; this is because the security domains are artificially obtained via k-means and hence have the same underlying data distribution.

Table 2 shows DDI values obtained from a suite of state-of-the-art MIAs. Across domain pairs, the access advantage (distinguishability) scores approach $\alpha=1.0$, indicating that an external auditor can almost perfectly identify the active domain (even when the domain distributions are similar as in the case of PubMedQA). Hence, even when UGI values fall significantly below 1.0 because of model generalization, the high DDI values show that access control in *Activate* still functions as intended. This clearly suggests that DDI is the better method for PermLLM access control efficacy evaluation.

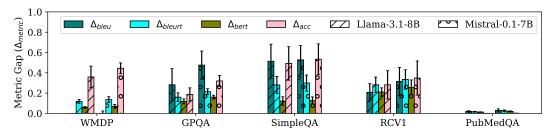


Figure 2: Utility Gap Index, Δ_U (mean \pm std) when user has access to one security domain.

Table 2: DDI values with $m \in \{\text{AUC-ROC}, \text{TPR@1\%FPR}, \text{TPR@5\%FPR}\}$ for the different MIAs. Mink++ is run with k = 10%. Entries are reported as $mean \pm std$ across domains.

	MIA	auc-roc	Llama-3.1-8B tpr@1%fpr	tpr@5%fpr	auc-roc	Mistral-0.1-7B tpr@1%fpr	tpr@5%fpr
WMDP	Loss ZLIB Mink++ Ref	0.99 ± 0.01 0.98 ± 0.03 1.00 ± 0.00 0.99 ± 0.01	0.93 ± 0.10 0.77 ± 0.31 0.99 ± 0.01 0.93 ± 0.10	0.96 ± 0.06 0.85 ± 0.21 1.00 ± 0.00 0.96 ± 0.06	$\begin{aligned} 1.00 &\pm 0.00 \\ 0.99 &\pm 0.02 \\ 1.00 &\pm 0.00 \\ 1.00 &\pm 0.00 \end{aligned}$	0.95 ± 0.06 0.85 ± 0.25 1.00 ± 0.00 0.95 ± 0.08	0.99 ± 0.01 0.92 ± 0.14 1.00 ± 0.00 0.98 ± 0.03
GPQA	Loss ZLIB Mink++ Ref	0.97 ± 0.05 0.95 ± 0.04 1.00 ± 0.00 1.00 ± 0.00	$\begin{array}{c} 0.81 \pm 0.26 \\ 0.45 \pm 0.22 \\ 1.00 \pm 0.01 \\ 0.97 \pm 0.04 \end{array}$	0.94 ± 0.08 0.77 ± 0.15 1.00 ± 0.00 0.99 ± 0.01	0.98 ± 0.03 0.97 ± 0.02 1.00 ± 0.00 1.00 ± 0.00	$\begin{array}{c} 0.93 \pm 0.10 \\ 0.57 \pm 0.24 \\ 0.99 \pm 0.01 \\ 0.97 \pm 0.05 \end{array}$	0.95 ± 0.07 0.83 ± 0.13 1.00 ± 0.00 0.99 ± 0.02
SimpleQA	Loss ZLIB Mink++ Ref	0.98 ± 0.03 0.98 ± 0.03 0.98 ± 0.03 0.98 ± 0.04	0.81 ± 0.34 0.80 ± 0.33 0.81 ± 0.32 0.78 ± 0.36	0.90 ± 0.25 0.90 ± 0.23 0.91 ± 0.21 0.90 ± 0.25	0.99 ± 0.03 0.99 ± 0.03 0.99 ± 0.03 0.98 ± 0.03	0.81 ± 0.32 0.80 ± 0.33 0.82 ± 0.31 0.79 ± 0.36	$\begin{array}{c} 0.92 \pm 0.20 \\ 0.91 \pm 0.20 \\ 0.92 \pm 0.21 \\ 0.90 \pm 0.24 \end{array}$
RCV1	Loss ZLIB Mink++ Ref	0.99 ± 0.01 0.93 ± 0.07 1.00 ± 0.00 0.99 ± 0.01	$\begin{array}{c} 0.86 \pm 0.21 \\ 0.71 \pm 0.26 \\ 0.97 \pm 0.05 \\ 0.77 \pm 0.28 \end{array}$	0.97 ± 0.06 0.81 ± 0.18 0.99 ± 0.01 0.99 ± 0.03	$\begin{array}{c} 0.99 \pm 0.02 \\ 0.94 \pm 0.08 \\ 1.00 \pm 0.01 \\ 0.99 \pm 0.01 \end{array}$	$\begin{array}{c} 0.85 \pm 0.24 \\ 0.73 \pm 0.28 \\ 0.96 \pm 0.06 \\ 0.80 \pm 0.28 \end{array}$	0.96 ± 0.09 0.83 ± 0.19 0.99 ± 0.02 0.98 ± 0.05
PubMedQA	Loss ZLIB Mink++ Ref	$\begin{array}{c} 0.81 \pm 0.07 \\ 0.77 \pm 0.07 \\ 0.90 \pm 0.02 \\ 1.00 \pm 0.00 \end{array}$	$\begin{array}{c} 0.16 \pm 0.11 \\ 0.10 \pm 0.05 \\ 0.31 \pm 0.08 \\ 0.98 \pm 0.02 \end{array}$	$\begin{array}{c} 0.36 \pm 0.15 \\ 0.30 \pm 0.13 \\ 0.57 \pm 0.08 \\ 1.00 \pm 0.00 \end{array}$	$\begin{array}{c} 0.95 \pm 0.03 \\ 0.88 \pm 0.05 \\ 0.99 \pm 0.01 \\ 1.00 \pm 0.00 \end{array}$	$\begin{array}{c} 0.51 \pm 0.21 \\ 0.32 \pm 0.17 \\ 0.93 \pm 0.07 \\ 1.00 \pm 0.00 \end{array}$	0.75 ± 0.14 0.57 ± 0.15 0.98 ± 0.02 1.00 ± 0.00

5.1.2 Multiple Active Domains

As discussed earlier in § 3, we explore three methods of combining knowledge from multiple domains the user has access to: (a) activating all the domain-specific LoRA modules (Activate), (b) merging the LoRA modules (Merge), and (c) training separate LoRA modules on the union of domains and using those for inference (*Union*). Table 3 reports the UGI (Δ_U) for these approaches when the user has access to two domains for all the data sets. We note that WMDP and GPQA have only three security domains, and hence activating any two domains always lead to overlap when calculating Δ_U as per 4.2. For these data sets, we calculate Δ_U on the non-overlapping data. Activate is computationally inexpensive but suffers from considerable utility loss when compared to the previous case of single domain. This is due to the high interference across the multiple domains in the activation space, which is a known issue in the multi-task learning literature [49, 39, 31]. The utility loss suppresses the absolute Δ_U in our experiments. As can be seen in Figure 3, *Merge* reduces the cross-domain interference, but still suffers from utility loss. Interestingly Merge achieves even lower Δ_U than Activate when combining two domains, as shown in Table 3. Although it quickly outperforms Activate when the user has access to more than two domains, the utility loss due to model merging interference [38, 43, 46, 50] also results in progressive degradation of Δ_U (see Figure 3). Union retains Δ_U even beyond four domains, and hence is the best choice when combining knowledge from several domains. But this comes at the cost of more training-time computation since new domain-specific modules have to be trained for the union of domains, and there could be potential combinatorial blow-up of the number of such combinations. As with the single active domain case,

Table 3: Utility Gap Index (Δ_U) for models with different approaches of combining domains when
user has access to two domains. All reported values are $mean \pm std$ across domains.

	Metric		Llama-3.1-8B			Mistral-0.1-7B	
		Activate	Merge	Union	Activate	Merge	Union
WMDP	$\Delta_{bleurt} \ \Delta_{bert}$	0.09 ± 0.01 0.05 ± 0.01	0.07 ± 0.02 0.03 ± 0.01	0.11 ± 0.02 0.06 ± 0.01	0.10 ± 0.02 0.05 ± 0.01	0.08 ± 0.03 0.04 ± 0.02	0.14 ± 0.03 0.07 ± 0.02
_ ≥	Δ_{acc}	0.27 ± 0.07	0.21 ± 0.09	0.34 ± 0.11	0.32 ± 0.04	0.25 ± 0.07	0.49 ± 0.09
GPQA	$egin{array}{l} \Delta_{bleu} \ \Delta_{bleurt} \ \Delta_{bert} \ \Delta_{acc} \end{array}$	$\begin{array}{c} 0.15 \pm 0.06 \\ 0.10 \pm 0.02 \\ 0.07 \pm 0.02 \\ 0.09 \pm 0.04 \end{array}$	$\begin{array}{c} 0.11 \pm 0.06 \\ 0.06 \pm 0.02 \\ 0.04 \pm 0.03 \\ 0.05 \pm 0.02 \end{array}$	$\begin{array}{c} 0.51 \pm 0.07 \\ 0.26 \pm 0.03 \\ 0.18 \pm 0.02 \\ 0.31 \pm 0.08 \end{array}$	$\begin{array}{c} 0.24 \pm 0.10 \\ 0.14 \pm 0.06 \\ 0.11 \pm 0.04 \\ 0.16 \pm 0.07 \end{array}$	$\begin{array}{c} 0.17 \pm 0.10 \\ 0.10 \pm 0.04 \\ 0.08 \pm 0.03 \\ 0.08 \pm 0.07 \end{array}$	$\begin{array}{c} 0.62 \pm 0.02 \\ 0.32 \pm 0.02 \\ 0.21 \pm 0.02 \\ 0.51 \pm 0.04 \end{array}$
SimpleQA	$egin{array}{l} \Delta_{bleu} \ \Delta_{bleurt} \ \Delta_{bert} \ \Delta_{acc} \end{array}$	0.26 ± 0.09 0.16 ± 0.05 0.07 ± 0.03 0.20 ± 0.07	$\begin{array}{c} 0.23 \pm 0.09 \\ 0.12 \pm 0.04 \\ 0.05 \pm 0.02 \\ 0.18 \pm 0.07 \end{array}$	$\begin{array}{c} 0.61 \pm 0.03 \\ 0.32 \pm 0.04 \\ 0.14 \pm 0.02 \\ 0.59 \pm 0.05 \end{array}$	0.30 ± 0.13 0.19 ± 0.05 0.08 ± 0.03 0.27 ± 0.09	$\begin{array}{c} 0.25 \pm 0.04 \\ 0.14 \pm 0.02 \\ 0.06 \pm 0.01 \\ 0.21 \pm 0.03 \end{array}$	$\begin{array}{c} 0.61 \pm 0.08 \\ 0.33 \pm 0.05 \\ 0.14 \pm 0.03 \\ 0.62 \pm 0.09 \end{array}$
RCV1	$egin{array}{l} \Delta_{bleu} \ \Delta_{bleurt} \ \Delta_{bert} \ \Delta_{acc} \end{array}$	0.05 ± 0.03 0.11 ± 0.01 0.08 ± 0.01 0.03 ± 0.01	$\begin{array}{c} 0.04 \pm 0.02 \\ 0.07 \pm 0.03 \\ 0.06 \pm 0.02 \\ 0.04 \pm 0.04 \end{array}$	$\begin{array}{c} 0.16 \pm 0.09 \\ 0.22 \pm 0.08 \\ 0.16 \pm 0.04 \\ 0.24 \pm 0.14 \end{array}$	$\begin{array}{c} 0.04 \pm 0.02 \\ 0.08 \pm 0.01 \\ 0.06 \pm 0.01 \\ 0.02 \pm 0.02 \end{array}$	0.01 ± 0.03 0.03 ± 0.04 0.03 ± 0.05 0.01 ± 0.03	$\begin{array}{c} 0.19 \pm 0.10 \\ 0.22 \pm 0.08 \\ 0.18 \pm 0.06 \\ 0.26 \pm 0.15 \end{array}$
PubMedQA	$egin{array}{c} \Delta_{bleu} \ \Delta_{bleurt} \ \Delta_{bert} \ \Delta_{acc} \end{array}$	0.01 ± 0.00 0.01 ± 0.00 0.01 ± 0.00	0.00 ± 0.00 0.00 ± 0.00 0.00 ± 0.00	0.01 ± 0.00 0.01 ± 0.00 0.01 ± 0.00	0.01 ± 0.00 0.01 ± 0.00 0.01 ± 0.00	0.00 ± 0.00 0.00 ± 0.00 0.00 ± 0.00	0.01 ± 0.01 0.01 ± 0.01 0.01 ± 0.00

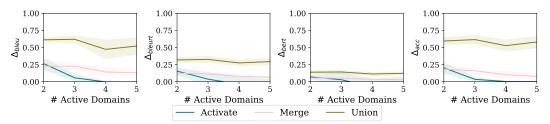


Figure 3: Utility Gap Index, Δ_U ($mean \pm std$) for Llama-3.1-8B models fine-tuned on SimpleQA when user has access to multiple security domains.

we observe close to zero utility gap on PubMedQA as the domains share the same data distribution. We observe similar results for Mistral-0.1-7B model (see Figure 9 in the appendix).

The DDI results for a two-domain setting appear in Table 4 (See Table 9 and Table 10 in the appendix for complete results). As we can see, we achieve high DDI values (e.g., close to $\alpha=1.0$ for auc-roc). In other words, an auditor can *almost perfectly* identify which domain is in effect, even when the corresponding utility gap (Δ_U) is far below 1.0 (Figure 3). *Union* consistently attains the highest DDI, followed by *Activate* and then *Merge* mirroring the trend observed with Δ_U . *Union*'s superiority however comes at the cost of greater tuning-time computation. *Union*'s near-perfect distinguishability mirrors the effect of model performance (with increasing domains) on Δ_U (see Figure 3). Crucially, the high DDI values confirm that even when Δ_U drops due to model generalization or degradation due to activation space or parameter interference, access control remains uncompromised; DDI therefore provides the more sensitive indicator of enforcement efficacy.

6 Conclusion and Discussion

We presented a comprehensive treatment of the access control problem on fine-tuned LLMs that includes novel formalism, empirical evaluation metrics, access control enforcement mechanisms, and evaluation of the mechanisms as well as the proposed metrics. We formalized a new class of LLMs called *Permissioned LLMs (PermLLM)* whose access control enforcement can be verified both theoretically and empirically using the formal tools provided in our work.

Table 4: DDI values (auc-roc) for models with different approaches of combining domains when user has access to two domains. All reported values are $mean \pm std$ across domains

	Metric		Llama-3.1-8B			Mistral-0.1-7B	
		Activate	Merge	Union	Activate	Merge	Union
WMDP	Loss ZLIB Mink++ Ref	0.98 ± 0.02 0.92 ± 0.08 0.90 ± 0.05 1.00 ± 0.00	0.93 ± 0.05 0.86 ± 0.09 0.94 ± 0.04 0.99 ± 0.00	0.99 ± 0.02 0.97 ± 0.05 1.00 ± 0.00 1.00 ± 0.00	0.99 ± 0.02 0.93 ± 0.09 0.96 ± 0.02 1.00 ± 0.00	0.95 ± 0.04 0.87 ± 0.09 0.94 ± 0.03 0.99 ± 0.00	0.99 ± 0.01 0.98 ± 0.03 1.00 ± 0.00 1.00 ± 0.00
GPQA	Loss ZLIB Mink++ Ref	$\begin{array}{c} 0.99 \pm 0.01 \\ 0.90 \pm 0.06 \\ 0.95 \pm 0.06 \\ 1.00 \pm 0.00 \end{array}$	$\begin{array}{c} 0.93 \pm 0.02 \\ 0.82 \pm 0.07 \\ 0.97 \pm 0.03 \\ 0.99 \pm 0.01 \end{array}$	$\begin{aligned} 1.00 &\pm 0.00 \\ 0.99 &\pm 0.01 \\ 1.00 &\pm 0.00 \\ 1.00 &\pm 0.00 \end{aligned}$	$\begin{array}{c} 0.99 \pm 0.01 \\ 0.93 \pm 0.08 \\ 0.98 \pm 0.02 \\ 1.00 \pm 0.00 \end{array}$	$\begin{array}{c} 0.96 \pm 0.04 \\ 0.86 \pm 0.09 \\ 0.98 \pm 0.01 \\ 0.99 \pm 0.02 \end{array}$	$\begin{aligned} 1.00 &\pm 0.00 \\ 0.99 &\pm 0.01 \\ 1.00 &\pm 0.00 \\ 1.00 &\pm 0.00 \end{aligned}$
SimpleQA	Loss ZLIB Mink++ Ref	0.96 ± 0.03 0.94 ± 0.04 0.85 ± 0.10 0.96 ± 0.03	0.95 ± 0.03 0.93 ± 0.04 0.92 ± 0.03 0.96 ± 0.04	$\begin{array}{c} 0.97 \pm 0.04 \\ 0.97 \pm 0.04 \\ 0.97 \pm 0.03 \\ 0.97 \pm 0.04 \end{array}$	0.97 ± 0.03 0.97 ± 0.03 0.92 ± 0.04 0.98 ± 0.03	0.96 ± 0.02 0.95 ± 0.03 0.93 ± 0.05 0.98 ± 0.03	0.97 ± 0.04 0.97 ± 0.04 0.97 ± 0.04 0.96 ± 0.04
RCV1	Loss ZLIB Mink++ Ref	$\begin{array}{c} 0.96 \pm 0.02 \\ 0.82 \pm 0.02 \\ 0.80 \pm 0.13 \\ 0.97 \pm 0.01 \end{array}$	$\begin{array}{c} 0.90 \pm 0.01 \\ 0.72 \pm 0.02 \\ 0.84 \pm 0.07 \\ 0.95 \pm 0.00 \end{array}$	0.98 ± 0.00 0.90 ± 0.05 0.99 ± 0.00 0.98 ± 0.01	$\begin{array}{c} 0.93 \pm 0.04 \\ 0.82 \pm 0.05 \\ 0.69 \pm 0.25 \\ 0.96 \pm 0.02 \end{array}$	$\begin{array}{c} 0.85 \pm 0.01 \\ 0.69 \pm 0.03 \\ 0.70 \pm 0.16 \\ 0.94 \pm 0.01 \end{array}$	0.98 ± 0.01 0.90 ± 0.05 0.99 ± 0.00 0.98 ± 0.00
PubMedQA	Loss ZLIB Mink++ Ref	0.84 ± 0.05 0.78 ± 0.05 0.85 ± 0.10 0.99 ± 0.02	0.66 ± 0.02 0.61 ± 0.01 0.79 ± 0.05 0.99 ± 0.01	0.79 ± 0.04 0.72 ± 0.04 0.94 ± 0.02 1.00 ± 0.00	0.78 ± 0.11 0.73 ± 0.10 0.72 ± 0.24 0.93 ± 0.10	0.65 ± 0.01 0.61 ± 0.01 0.74 ± 0.06 0.98 ± 0.00	0.81 ± 0.05 0.74 ± 0.05 0.95 ± 0.02 1.00 ± 0.00

Limitations. Our approach does not support deep hierarchy of domains with arbitrary overlaps. Another issue we observe is with the scalability beyond a handful of domains. This either leads to severe degradation of utility (as in the case of *Activate*) or it becomes compute-intensive (for *Union*). We leave this exploration for future work. We also note some limitations in the experiments that we do not expect to change our key claims. First, we only run one model fine-tuning per parameter setting due to the computation overhead. Second, we use the default value for LoRA rank as our preliminary experiments with different ranks suggested limited impact on model utility. For our formalism in § 2, we assume that adversaries do not tamper with their credentials or domain access, otherwise they can gain arbitrary domain information. This is enforced by the enclosing system via authentication. Finally, we note that our DDI metrics are only as good as the best possible MIA.

Related Work. Access control problems in agentic systems can manifest in interesting ways, such as context hijacking [2], and may require further constraining the purview of individual agent contexts. Retrieval Augmented Generation (RAG) systems [23, 32, 51] are also susceptible to the access control problem. However, the access control needs to be enforced in the information retrieval engine of the system [4, 14] and is beyond our work's scope (although we do provide a formalism for access control in RAG-based systems in Appendix A). One may draw some parallels between our formalism of response relevance and access advantage metric with prior works on *indistinguishability* [1, 3, 9, 13] in security and privacy. The mechanisms in this lineage of works are singularly focused on eliminating distinguishability between the effects of different data on computations. In contrast, PermLLM's objective is to maximize domain separation, which implies maximization of distinguishability.

Broader Impacts. Our work aims to bolster the security and privacy of individual's data by enforcing strict access control, such that only people with prior authorization can get access to the information. Our work is applicable to healthcare, finance, and more broadly, enterprise / governance applications that deal with sensitive data of individuals.

7 Acknowledgements

We would like to thank Pallika Kanani and Dan Roth for fruitful discussions that motivated the problem setup of access control in large language model based applications. We would also like to thank David Evans and the anonymous reviewers for providing helpful feedback on the paper.

References

- [1] Afonso Arriaga, Manuel Barbosa, and Pooya Farshim. Private functional encryption: Indistinguishability-based definitions and constructions from obfuscation. Cryptology ePrint Archive, Paper 2016/018, 2016. URL https://eprint.iacr.org/2016/018.
- [2] Eugene Bagdasarian, Ren Yi, Sahra Ghalebikesabi, Peter Kairouz, Marco Gruteser, Sewoong Oh, Borja Balle, and Daniel Ramage. Airgapagent: Protecting privacy-conscious conversational agents, 2024. URL https://arxiv.org/abs/2405.05175.
- [3] Jens-Matthias Bohli and Andreas Pashalidis. Relations among privacy notions. *ACM Trans. Inf. Syst. Secur.*, 14(1), 2011. URL https://doi.org/10.1145/1952982.1952986.
- [4] Stefan Büttcher and Charles L. A. Clarke. A security model for full-text file system search in multi-user environments. In *Proceedings of the 4th Conference on USENIX Conference on File and Storage Technologies Volume 4*, page 13, USA, 2005. USENIX Association.
- [5] Nicholas Carlini, Florian Tramer, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom Brown, Dawn Song, Ulfar Erlingsson, et al. Extracting training data from large language models. In *30th USENIX Security Symposium (USENIX Security 21)*, pages 2633–2650, 2021.
- [6] Nicholas Carlini, Steve Chien, Milad Nasr, Shuang Song, Andreas Terzis, and Florian Tramer. Membership inference attacks from first principles. In 2022 IEEE symposium on security and privacy (SP), pages 1897–1914. IEEE, 2022.
- [7] Shih-Han Chan. Encrypted prompt: Securing Ilm applications against unauthorized actions, 2025. URL https://arxiv.org/abs/2503.23250.
- [8] Yang Chen, Ethan Mendes, Sauvik Das, Wei Xu, and Alan Ritter. Can language models be instructed to protect personal information?, 2023. URL https://arxiv.org/abs/2310. 02224.
- [9] Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. Calibrating noise to sensitivity in private data analysis. In *Proceedings of the Third Conference on Theory of Cryptography*, TCC'06, pages 265–284, 2006.
- [10] David F. Ferraiolo, John F. Barkley, and D. Richard Kuhn. A role-based access control model and reference implementation within a corporate intranet. ACM Trans. Inf. Syst. Secur., 2(1): 34–64, 1999. ISSN 1094-9224. URL https://doi.org/10.1145/300830.300834.
- [11] David F. Ferraiolo, D. Richard Kuhn, and Ramaswamy Chandramouli. Role-based access control. *Information Security and Privacy Series*, 2007.
- [12] William Fleshman, Aleem Khan, Marc Marone, and Benjamin Van Durme. Adapterswap: Continuous training of llms with data removal and access-control guarantees, 2025. URL https://arxiv.org/abs/2404.08417.
- [13] Shafi Goldwasser, Silvio Micali, and Ronald L. Rivest. A digital signature scheme secure against adaptive chosen-message attacks. SIAM Journal on Computing, 17(2):281–308, 1988.
- [14] Pawam Goyal. Private information retrieval with access control, 2023.
- [15] Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd of models, 2024. URL https://arxiv.org/abs/2407.21783.
- [16] Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin de Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning for nlp, 2019. URL https://arxiv.org/abs/1902.00751.
- [17] Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models, 2021. URL https://arxiv.org/abs/2106.09685.

- [18] Bargav Jayaraman, Lingxiao Wang, Katherine Knipmeyer, Quanquan Gu, and David Evans. Revisiting membership inference under realistic assumptions. *Proceedings on Privacy Enhancing Technologies*, 2021.
- [19] Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. Mistral 7b, 2023. URL https://arxiv.org/abs/2310.06825.
- [20] Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William Cohen, and Xinghua Lu. PubMedQA: A dataset for biomedical research question answering. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 2567–2577. Association for Computational Linguistics, 2019. doi: 10.18653/v1/D19-1259. URL https://aclanthology.org/D19-1259/.
- [21] Jean Kaddour, Joshua Harris, Maximilian Mozes, Herbie Bradley, Roberta Raileanu, and Robert McHardy. Challenges and applications of large language models, 2023. URL https://arxiv.org/abs/2307.10169.
- [22] David D Lewis, Yiming Yang, Tony G Rose, and Fan Li. Rcv1: A new benchmark collection for text categorization research. *Journal of machine learning research*, 5:361–397, 2004.
- [23] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. Retrieval-augmented generation for knowledge-intensive nlp tasks. In Proceedings of the 34th International Conference on Neural Information Processing Systems, NIPS '20, Red Hook, NY, USA, 2020. Curran Associates Inc. ISBN 9781713829546.
- [24] Nathaniel Li, Alexander Pan, Anjali Gopal, Summer Yue, Daniel Berrios, Alice Gatti, Justin D Li, Ann-Kathrin Dombrowski, Shashwat Goel, Long Phan, et al. The WMDP benchmark: Measuring and reducing malicious use with unlearning. In *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pages 28525–28550. PMLR, 2024. URL https://proceedings.mlr.press/v235/li24bc.html.
- [25] Qin Liu, Fei Wang, Chaowei Xiao, and Muhao Chen. Sudolm: Learning access control of parametric knowledge with authorization alignment, 2025. URL https://arxiv.org/ abs/2410.14676.
- [26] Yi Liu, Gelei Deng, Zhengzi Xu, Yuekang Li, Yaowen Zheng, Ying Zhang, Lida Zhao, Tianwei Zhang, Kailong Wang, and Yang Liu. Jailbreaking chatgpt via prompt engineering: An empirical study, 2024. URL https://arxiv.org/abs/2305.13860.
- [27] Anay Mehrotra, Manolis Zampetakis, Paul Kassianik, Blaine Nelson, Hyrum Anderson, Yaron Singer, and Amin Karbasi. Tree of attacks: Jailbreaking black-box llms automatically. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang, editors, Advances in Neural Information Processing Systems, volume 37, pages 61065-61105. Curran Associates, Inc., 2024. URL https://proceedings.neurips.cc/paper_files/paper/2024/file/70702e8cbb4890b4a467b984ae59828a-Paper-Conference.pdf.
- [28] Fatemehsadat Mireshghallah, Kartik Goyal, Archit Uniyal, Taylor Berg-Kirkpatrick, and Reza Shokri. Quantifying privacy risks of masked language models using membership inference attacks. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 8332–8347, 2022.
- [29] Hamid Mozaffari and Virendra J Marathe. Semantic membership inference attack against large language models. *arXiv preprint arXiv:2406.10218*, 2024.
- [30] NIST SP Joint Task Force. Security and privacy controls for information systems and organizations, *nist sp* 800-53 rev. 5, 2020. URL https://csrc.nist.gov/pubs/sp/800/53/r5/upd1/final.

- [31] Oleksiy Ostapenko, Zhan Su, Edoardo Maria Ponti, Laurent Charlin, Nicolas Le Roux, Matheus Pereira, Lucas Caccia, and Alessandro Sordoni. Towards modular llms by building and reusing a library of loras, 2024. URL https://arxiv.org/abs/2405.11157.
- [32] OWASP GenAI Security Project. Llm08:2025 vector and embedding weaknesses, 2025. URL https://genai.owasp.org/llmrisk/llm082025-vector-and-embedding-weaknesses/#:~:text=1. %20Permission%20and%20access%20control.
- [33] David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R. Bowman. Gpqa: A graduate-level google-proof q&a benchmark, 2023. URL https://arxiv.org/abs/2311.12022.
- [34] Shang Shang, Zhongjiang Yao, Yepeng Yao, Liya Su, Zijing Fan, Xiaodan Zhang, and Zhengwei Jiang. Intentobfuscator: A jailbreaking method via confusing llm with prompts. In *ESORICS*, pages 146–165, 2024. URL https://doi.org/10.1007/978-3-031-70903-6_8.
- [35] William F. Shen, Xinchi Qiu, Meghdad Kurmanji, Alex Iacob, Lorenzo Sani, Yihong Chen, Nicola Cancedda, and Nicholas D. Lane. Lunar: Llm unlearning via neural activation redirection, 2025. URL https://arxiv.org/abs/2502.07218.
- [36] Weijia Shi, Anirudh Ajith, Mengzhou Xia, Yangsibo Huang, Daogao Liu, Terra Blevins, Danqi Chen, and Luke Zettlemoyer. Detecting pretraining data from large language models. *arXiv* preprint arXiv:2310.16789, 2023.
- [37] Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. Membership inference attacks against machine learning models. In 2017 IEEE symposium on security and privacy (SP), pages 3–18. IEEE, 2017.
- [38] George Stoica, Pratik Ramesh, Boglarka Ecsedi, Leshem Choshen, and Judy Hoffman. Model merging with SVD to tie the Knots, 2024. URL https://arxiv.org/abs/2410. 19735.
- [39] Tu Vu, Brian Lester, Noah Constant, Rami Al-Rfou, and Daniel Cer. Spot: Better frozen model adaptation through soft prompt transfer, 2022. URL https://arxiv.org/abs/2110. 07904.
- [40] Alexander Wei, Nika Haghtalab, and Jacob Steinhardt. Jailbroken: How does llm safety training fail? In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, Advances in Neural Information Processing Systems, volume 36, pages 80079–80110. Curran Associates, Inc., 2023. URL https://proceedings.neurips.cc/paper_files/paper/2023/file/fd6613131889a4b656206c50a8bd7790-Paper-Conference.pdf.
- [41] Jason Wei, Nguyen Karina, Hyung Won Chung, Yunxin Joy Jiao, Spencer Papay, Amelia Glaese, John Schulman, and William Fedus. Measuring short-form factuality in large language models, 2024. URL https://arxiv.org/abs/2411.04368.
- [42] Lingling Xu, Haoran Xie, Si-Zhao Joe Qin, Xiaohui Tao, and Fu Lee Wang. Parameter-efficient fine-tuning methods for pretrained language models: A critical review and assessment, 2023. URL https://arxiv.org/abs/2312.12148.
- [43] Prateek Yadav, Derek Tam, Leshem Choshen, Colin Raffel, and Mohit Bansal. Ties-merging: Resolving interference when merging models, 2023. URL https://arxiv.org/abs/2306.01708.
- [44] Samuel Yeom, Irene Giacomelli, Matt Fredrikson, and Somesh Jha. Privacy risk in machine learning: Analyzing the connection to overfitting, 2018.
- [45] Fangcong Yin, Xi Ye, and Greg Durrett. Lofit: Localized fine-tuning on LLM representations. In Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems, 2024.

- [46] Le Yu, Bowen Yu, Haiyang Yu, Fei Huang, and Yongbin Li. Language models are super mario: Absorbing abilities from homologous models as a free lunch, 2024. URL https://arxiv.org/abs/2311.03099.
- [47] Jingyang Zhang, Jingwei Sun, Eric Yeats, Yang Ouyang, Martin Kuo, Jianyi Zhang, Hao Yang, and Hai Li. Min-k%++: Improved baseline for detecting pre-training data from large language models. *arXiv preprint arXiv:2404.02936*, 2024.
- [48] Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen Yang, Yushuo Chen, Zhipeng Chen, Jinhao Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu, Peiyu Liu, Jian-Yun Nie, and Ji-Rong Wen. A survey of large language models, 2025. URL https://arxiv.org/abs/2303.18223.
- [49] Xiangyun Zhao, Haoxiang Li, Xiaohui Shen, Xiaodan Liang, and Ying Wu. A modulation module for multi-task learning with applications in image retrieval. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018.
- [50] Ziyu Zhao, Tao Shen, Didi Zhu, Zexi Li, Jing Su, Xuwu Wang, Kun Kuang, and Fei Wu. Merging loras like playing lego: Pushing the modularity of lora to extremes through rank-wise clustering, 2024. URL https://arxiv.org/abs/2409.16167.
- [51] Hongbin Zhong, Matthew Lentz, Nina Narodytska, Adriana Szekeres, and Kexin Rong. Honeybee: Efficient role-based access control for vector databases via dynamic partitioning, 2025. URL https://arxiv.org/abs/2505.01538.

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: The key claims of the paper mentioned in the abstract and the introduction are accurately reflected in the rest of the paper.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the
 contributions made in the paper and important assumptions and limitations. A No or
 NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: We discuss the key limitations of our work in § 6. We also clearly mention the underlying assumptions for our formalism of access control in § 2. Finally, in the experiments with multiple active domains in § 5.1.2, we find overlap when activating any two domains for WMDP and GPQA data sets to calculate the access advantage, as these data sets have only three security domains. For these experiments, we calculated the access advantage metrics on the non-overlapping data.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [Yes]

Justification: We provide the proof of correctness of our access control mechanisms and the underlying assumptions in Appendix B.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and crossreferenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: We provide a detailed description of our data sets and model fine-tuning hyperparameters in Appendix D.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
- (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
- (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
- (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
- (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility.

In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [No]

Justification: Due to our organizational policies regarding intellectual property, we can not currently open-source the code and the fine-tuned models. However, we provide complete citations for the public data sets we use and also describe the model fine-tuning hyperparameter settings in detail in Appendix D. We use the publicly available pretrained models for our fine-tuning, which are also properly cited in the paper.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (https://nips.cc/ public/quides/CodeSubmissionPolicy) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https: //nips.cc/public/quides/CodeSubmissionPolicy) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: We describe the training and test data set splits, and the model fine-tuning hyperparameters in Appendix D.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: We provide mean \pm standard deviations for all our results reported in the experiments. These are calculated across the different domains for each data set. Only the training and validation set losses are not reported with $mean \pm std$ as we only fine-tune a single model per setting due to the computational overhead.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: We discuss the compute resources used for our experiments in Appendix D.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?

Answer: [Yes]

Justification: We strictly adhere to the code of ethics. We made sure to preserve complete anonymity in the submission.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: While we see no negative societal impacts of our work, we discuss the broader impacts of our work in § 6.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: We do not intend to release any of our fine-tuned models. We do not show any of the data set records in the paper, even though the data sets we use are publicly available under permissible licenses.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do
 not require this, but we encourage authors to take this into account and make a best
 faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: We mention the licenses for the data sets and the pretrained models we used in our experiments in Appendix D.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: We do not release any models or new data sets as a part of this submission. The core contributions are the set of mechanisms to ensure access control in LLMs, and the evaluation metrics to audit the effectiveness of access control.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects. Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects. Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent)
 may be required for any human subjects research. If you obtained IRB approval, you
 should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [NA]

Justification: We only use pretrained LLMs for our experimental evaluation of fine-tuning with access control over various data sets. LLMs were not used for writing the paper, nor do they replace any human intellectual efforts, and as such do not contribute to the novelty of the paper.

Guidelines:

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (https://neurips.cc/Conferences/2025/LLM) for what should or should not be described.

A Formalizing Access Control for Retrieval Augmented Generation

For Retrieval Augmented Generation (RAG), we assume a pre-trained LLM f that is used in applications without additional fine-tuning. Instead, we augment f with a retriever engine R to give us a retrieval augmented LLM f_R .

Each query q_c to f_R is accompanied by a context c, retrieved by R, that enhances f_R 's response to the query. Let R retrieve contexts from the context database C, i.e. $c \in C$. Furthermore, we have $C = \bigcup_{i=1}^n C_{s_i} \sim C_{s_i}$, where each C_{s_i} is a collection of contexts belonging to security domain s_i .

For this discussion, we define \mathcal{M} as an access control mechanism that dictates the mapping of every $C_{s_i} \subseteq C$ to the security domain s_i . We say that a RAG system that uses contexts from the context database C is permissioned (PermRAG), if it uses retriever $R_C^{\mathcal{M}}$, which in turn uses the access control mechanism \mathcal{M} to retrieve context $c \in C_{s_i}$ from a selected security domain s_i . Intuitively, given a security domain s_i , R uses \mathcal{M} to retrieve context $c \in C_{s_i}$. One can trivially generalize this definition of PermRAG to work with subsets of security domains instead of a singleton security domain s_i .

For PermRAG, we assume an identical enclosing system setting as in PermLLM (see § 2): Given a user u the enclosing system determines u's access credentials $cred_u$ and calls authenticate ($cred_u$) that takes user credentials $cred_u$ and maps them to a subset of security domains S_u that u can access. User u cannot arbitrarily change S_u . Each of user u's subsequent query q to f_R is annotated with S_u . The retriever $R_C^{\mathcal{M}}$ of f_R uses access control mechanism \mathcal{M} to retrieve a context $c \in C_{S_u}$.

Definition A.1 (Relevant Response for PermRAG). Given a PermRAG f_R , with retriever R_C^M , a query q from user u, and S_u the security domains u has access to, $r = f_R(q)$ is the response by f_R to query q. Response r is said to be relevant to S_u (i.e. $r = r_{S_u}$) if retriever R_C^M uses a context $c \in C_{S_u}$ to augment the query for r.

To empirically quantify response relevance, we can use the same response relevance score, $relv(f_R(q), S_u)$ that quantifies the information gained on data in the security domains q's user u has access to (this is the same set of security domains that mapping \mathcal{M} gives for u for the retriever $R_C^{\mathcal{M}}$, i.e. S_u). Here $R_C^{\mathcal{M}}$ retrieves the query context $c \in C$ using mapping \mathcal{M} ; c is then augmented to the query q. We restrict the domain of relv to the real number interval [0,1], where 1 is the best expected score for relevance. Similar to PermLLM, we define $access\ advantage$ for PermRAG as follows:

Definition A.2 (Access Advantage for PermRAG). Given PermRAG f_R that uses retriever $R_C^{\mathcal{M}}$ which in turn uses the context database C containing data from domains $\mathbb{S} = \bigcup_{i=1}^n \{s_i\}$, with access control mechanism \mathcal{M} , a subset of security domains $S_u \subseteq \mathbb{S}$, context $c \in C$ that is augmented to query q, f_R achieves α -access advantage w.r.t. S_u if:

$$\mathbb{E}_{q \sim \mathcal{D}_{S_u}, S_v \subseteq \mathbb{S}; S_u \cap S_v = \phi} [relv(f_R(q), S_u) \ominus relv(f_R(q), S_v)] \ge \alpha$$

where relv() is the response relevance score on the corresponding security domain subset $(S_u \text{ or } S_v)$, \ominus is a "difference" operator specific to the access control assessment method (e.g. subtraction), and α is an advantage threshold that lies in the range [0,1].

B Formal Access Control Enforcement in PermLLM Mechanisms

We now present formal proofs for correct access control enforcement in our PermLLM mechanisms presented in § 3: *Activate*, *Merge*, and *Union*.

Refreshing the formalism from § 2, we consider a universe of n different security domains $\mathbb{S} = \bigcup_{i=1}^n \{s_i\}$, and a training data set consisting of data from these domains $D = \bigcup_{i=1}^n D_{s_i} \sim \mathcal{D}_{s_i}$ (here D_{s_i} is a data set sampled from data distribution \mathcal{D}_{s_i} of domain s_i). Let f_D be the LLM tuned using data set D. Let W be the set of f_D 's parameters. Model tuning changes values of a subset of W. Let security domain s_i affect, per the meaning of affect in § 2, a subset of parameters $W_{s_i} \subseteq W$. Thus data from D_{s_i} is used to change parameters W_{s_i} during model fine-tuning. Let \mathcal{M} be the access control mechanism that dictates the mapping of security domain s_i to parameters W_{s_i} via the affects relation.

Consider a set of LoRA adapters [17] $l_1, l_2, ..., l_m$. Each adapter l_i comprises parameters W_{l_i} , such that $W_{l_i} \cap W_{l_j} = \phi, \forall i \neq j$. Let i be the adapter Id for adapter l_i . Let $f_D^{\mathcal{M}}$ by the PermLLM that uses mapping \mathcal{M} of security domains to parameters during tuning and testing. Let $\mathcal{F}^{\mathcal{M}}$ be the system enclosing $f_D^{\mathcal{M}}$ that performs the mapping from user credentials $cred_u$ to sets of security domains S_u for each user u. We make two assumptions about $\mathcal{F}^{\mathcal{M}}$: (i) $\mathcal{F}^{\mathcal{M}}$ can correctly determine and maintain the security domains S_u a user u has access to; and (ii) S_u remains opaque to the user and any other adversary and as a result, cannot be tampered with by any user or adversary.

We assume that both fine-tuning and testing are mediated through $\mathcal{F}^{\mathcal{M}}$. During fine-tuning, the dataset D is passed to $\mathcal{F}^{\mathcal{M}}$. $\mathcal{F}^{\mathcal{M}}$ extracts information about the security domains $s_1,..,s_n$ covered by D. For settings where users have access to multiple security domains, the list of security domain combinations that users have access to is also passed on to $\mathcal{F}^{\mathcal{M}}$. $\mathcal{F}^{\mathcal{M}}$ does the mapping between security domain groups and LoRA adapters differently for each of our PermLLM mechanisms:

- Activate $\mathcal{F}^{\mathcal{M}}$ maps each security domain s_i to a unique LoRA adapter l_i . For fine-tuning of $f_D^{\mathcal{M}}$, minibatches sampled for each s_i are routed to the corresponding LoRA adapter l_i , the other LoRA adapters are deactivated for the sampled mini-batch.
- Merge Security domain-LoRA adapter mappings and fine-tuning of $f_D^{\mathcal{M}}$ proceeds identically to that in *Activate*. However, after the fine-tuning is done, the security domain groups are used to merged LoRA adapters. These new LoRA adapters are added to the set of LoRA adapters in $f_D^{\mathcal{M}}$. The mapping \mathcal{M} is also updated with the new mappings between security domain groups and LoRA adapters.
- Union Datasets corresponding to the security domain groups are used to fine-tuning unique LoRA adapters. \mathcal{M} is also updated with these new security domain group-LoRA adapter mappings.

At the end of fine-tuning, \mathcal{M} will have a mapping between each security domain group S_u (for each respective user u) and each LoRA adapter in mechanisms Merge and Union. More formally,

Lemma B.1. In Merge and Union, after fine-tuning, for every user u that has access to $S_u \subseteq \mathbb{S}$, $\exists l_{S_u}$, where l_{S_u} is a LoRA adapter, S_u affects parameters $W_{l_{S_u}}$, and $W_{l_{S_u}}$ is not affected by any other security domains in S.

In case of *Activate*, S_u is used at inference time to activate the LoRA adapters l_{s_i} , where $s_i \in S_u$. More formally,

Lemma B.2. In Activate, after fine-tuning, for every user u that has access to $S_u \subseteq \mathbb{S}$, $\forall s_i \in S_u$, s_i affects parameters $W_{l_{s_i}}$, and $W_{l_{s_i}}$ is not affected by any other security domain $s_j \in S_u$, $i \neq j$, or $s_k \in \mathbb{S} \setminus S_u$.

At inference time, user u sends query q to $\mathcal{F}^{\mathcal{M}}$. $\mathcal{F}^{\mathcal{M}}$ first determines u's security domains S_u , and then passes q and S_u to $f_D^{\mathcal{M}}$, which then activates the LoRA adapter/s corresponding to S_u : l_{S_u} in case of Merge and Union , and l_{s_i} , where $s_i \in S_u$, in case of $\mathit{Activate}$. Our assumptions about accessibility of S_u to the user or adversary ensure that the adversary cannot tamper with S_u within the scope of $\mathcal{F}^{\mathcal{M}}$.

Theorem B.3. Given any query q from any user u, the response $r = f_D^{\mathcal{M}}(q)$ is relevant to S_u for \mathcal{M} in Activate, Merge, or Union.

Proof. From Lemmas Theorem B.1 and Theorem B.2, through the fine-tuning process S_u affects parameters $W_{l_{S_u}}$ in Merge and Union, and parameters $W_{l_{s_i}}, \forall s_i \in S_u$ in Activate. At inference time, these same parameters (along with the pretrained model's parameters) are used to generate response $r = f_D^{\mathcal{M}}(q)$. By implication, the parameters affected by S_u are used to generated r. Hence r is relevant to S_u , i.e. $r = r_{S_u}$.

Since the above response relevance condition applies for all responses $r = f_D^{\mathcal{M}}(q)$ on all queries q by all users u, we say that Activate, Merge, and Union correctly enforce parameter separation and hence correctly enforce access control for all users u.

C Audit Games

We formalize black-box games that capture: (i) the distinguishability of security domain-specific responses for DDI, and (ii) the utility disparity induced by access restrictions for UGI. Intuitively, in these auditing games, we measure how *effectively* an external auditor can conclude if the access control mechanism is correctly implemented by verifying if the correct domain adapter is activated for a query. This effectiveness is directly correlated with the access advantage score for the target security domain(s). Higher access advantage score denotes *stronger* access control enforcement. A perfectly separated system provides the auditor with an access advantage score of 1.0.

We consider the same threat setting and auditor privileges for our adversarial games between auditor \mathcal{A} and system \mathcal{S} enclosing the PermLLM $f_{\mathcal{D}}^{\mathcal{M}}$ as described in § 2.3.

Game 1: Domain Distinguishability. This game assesses whether the auditor can effectively conclude if the correct security domains were used based on the generated responses. The primary motivation of this game is to measure the distinguishability of different security domains' distributions.

- 1. Auditor \mathcal{A} chooses security domain set S_u and emulates user u. \mathcal{A} sends user credentials $cred_u$ and query $q \sim \mathcal{D}_{S_u}$ to system \mathcal{S} . \mathcal{S} verifies the user credential $cred_u$ and sends back the model response $f_D^{\mathcal{M}}(q)$ to \mathcal{A} .
- 2. \mathcal{A} chooses security domain set S_v such that $S_v \cap S_u = \phi$ and emulates user v. \mathcal{A} sends user credentials $cred_v$ and the same query $q \sim \mathcal{D}_{S_u}$ to \mathcal{S} . \mathcal{S} verifies the user credential $cred_v$ and sends back the model response $f_D^{\mathcal{M}}(q)$ to \mathcal{A} .
- 3. \mathcal{A} sends the models responses and domain information to membership inference oracle O to obtain domain distinguishability score $m(O(f_D^{\mathcal{M}}(q)|S_u, f_D^{\mathcal{M}}(q)|S_v))$, where $m(\cdot)$ is a membership inference metric (e.g., AUC-ROC or TPR@1%FPR) in the [0,1] range.
- 4. \mathcal{A} concludes the access control mechanism is correctly implemented if the domain distinguishability score $m(O(f_D^{\mathcal{M}}(q)|S_u, f_D^{\mathcal{M}}(q)|S_v)) \geq \alpha$.

Note that we can change the above game to distinguish members $(q \sim \mathcal{D}_{S_u})$ and non-members $(q \sim \mathcal{D}_{S_v})$ for the target domain set S_u , similar to prior MIA setups, which is what we do in our experiments in § 5.

Game 2: Utility Gap Evaluation. The second game evaluates how distinctly the responses from two different security domains impact the utility perceived by users. The rationale behind this game is to confirm that enforced access controls result in meaningful variations in response quality.

- 1. Auditor \mathcal{A} chooses security domain set S_u and emulates user u. \mathcal{A} sends user credentials $cred_u$ and query $q \sim \mathcal{D}_{S_u}$ to system \mathcal{S} . \mathcal{S} verifies the user credential $cred_u$ and sends back the model response $f_D^{\mathcal{M}}(q)$ to \mathcal{A} .
- 2. \mathcal{A} chooses security domain set S_v such that $S_v \cap S_u = \phi$ and emulates user v. \mathcal{A} sends user credentials $cred_v$ and the same query $q \sim \mathcal{D}_{S_u}$ to \mathcal{S} . \mathcal{S} verifies the user credential $cred_v$ and sends back the model response $f_D^{\mathcal{M}}(q)$ to \mathcal{A} .
- 3. Given a utility function $U(\cdot)$ (e.g., BLEURT or task accuracy) that outputs values in [0,1] range, \mathcal{A} concludes the access control mechanism is correctly implemented if the utility gap score $|U(f_D^{\mathcal{M}}(q)|S_u) U(f_D^{\mathcal{M}}(q)|S_v)| \ge \alpha$.

We aggregate the utility gaps from this game across all domain set pairs to obtain our UGI metric.

D Detailed Experiment Setup

D.1 Models

For our instantiation of PermLLM, we fine-tune Llama-3.1-8B[15] and Mistral-0.1-7B[19] pretrained models on four datasets covering multiple distinct security domains (henceforth called *domains*), where we fine-tune a separate LoRA adapter for each domain. To compare our PermLLM, we train two additional models with full fine-tuning and LoRA fine-tuning respectively on entire training data. Note that these models are only used for utility baselines as they do not provide access control. For all the LoRA adapters, we use 64 rank and 0.1 dropout. We use AdamW optimizer with 0.1 weight decay

Table 5: Data Set Details. Generalization Loss Gap (i.e., gap between model's loss on training and test sets) for all models are reported after fine-tuning for 5 epochs on each data set.

Data Set	Data Se	Data Set Size		a-3.1-8B	Loss Gap	Mistral-0.1-7B Loss Gap		
(# Domains)	Train	Test	Full FT	LoRA	PermLLM	Full FT	LoRA	PermLLM
WMDP (3)	2936	732	1.96	0.52	1.15	1.36	0.65	1.07
GPQA (3)	360	88	2.51	1.06	1.04	1.58	0.61	1.09
SimpleQA (10)	4089	1018	2.91	0.96	1.49	1.87	0.90	1.25
RCV1 (4)	45622	22811	4.07	0.35	0.83	2.48	0.37	0.74
PubMedQA (10)	200000	11269	3.53	0.07	0.36	2.56	0.07	0.35

to fine-tuned all the models for 5 epochs with 300 warmup steps, 2 batch size and 5×10^{-4} learning rate (except for Mistral-0.1-7B full fine-tuning that uses a learning rate of 5×10^{-5}). We performed grid search over multiple learning rates and warmup steps and found these values to give the best results. For all our experiments, we use 8 H100 GPUs (with 80GB VRAM per GPU), 4 workers per GPU, and 384 GB RAM. One epoch of fine-tuning took from few minutes (for our smallest data set: GPQA) to a couple of hours (for our largest data set: RCV1). Mistral-0.1-7B is released under Apache 2.0 license, and Llama-3.1-8B is released under Llama 3.1 Community License.

D.2 Data Sets

For our experiments, we require data sets that consist of multiple distinct domains and are possibly not seen by the pretrained models. We use five different data sets, namely, WMDP [24], GPQA [33], SimpleQA [41], RCV1 [22], and PubMedQA [20]. While the first three data sets were collected after the pretraining cutoff dates for Llama-3.1-8B and Mistral-0.1-7B, RCV1 is an older data set and hence we do not know if it was used in pretraining. However, we observe a high initial training loss on this data set, thereby indicating that it was either not used in pretraining or was catastrophically forgotten by the models, allowing for a gradual reduction in training loss during our fine-tuning (see Figure 7). Table 5 shows the data set details, along with the generalization gap (test loss - train loss) for different approaches of fine-tuning the models on these data sets. See Figure 4, Figure 5, Figure 6, Figure 7, and Figure 8 for complete training and test loss trajectories across different data sets.

WMDP. Weapons of Mass Destruction Proxy (WMDP) [24] is a data set consisting of multi-choice question—answer pairs spanning three domains: biological weapons (*bio*), chemical weapons (*chem*) and cyber-warfare weapons (*cyber*). We do 4:1 split of the data set to obtain training and test sets. The training set consists of 2936 question—answer pairs where 1019 are from *bio*, 327 are from *chem* and the remaining 1590 are from *cyber*. The test set size is 732 records, consisting of 254 *bio*, 81 *chem* and 397 *cyber* records. The largest record from this data set consists of 1934 tokens (tokenized using Llama3 tokenizer). This data set is released under MIT License.

GPQA. Graduate-Level Google-Proof Q&A Benchmark (GPQA) [33] data set consists of general question—answer pairs from three domains: *biology*, *chemistry* and *physics*. We do 4:1 split of the data set to obtain training and test sets. The training set consists of 360 question—answer pairs where 63 are from *biology*, 147 are from *chemistry* and the remaining 150 are from *physics*. The test set size is 88 records, consisting of 15 *biology*, 36 *chemistry* and 37 *physics* records. The largest record from this data set consists of 911 tokens (tokenized using Llama3 tokenizer). This data set is released under MIT License.

SimpleQA. SimpleQA [41] is a factuality benchmark that measures the ability for language models to answer short, fact-seeking questions. It consists of general question–answer pairs from ten domains: art, geography, history, music, other, politics, science and technology, sports, tv shows, and video games. We do 4:1 split of the data set to obtain training and test sets. The training set consists of 4089 question–answer pairs divided across all ten domains. The test set size is 1018 records spanning across all ten domains. The largest record from this data set consists of 156 tokens (tokenized using Llama3 tokenizer). This data set is released under MIT License.

RCV1. RCV1 [22] is a benchmark dataset on text categorization. It is a collection of newswire articles produced by Reuters between 1996 and 1997. It contains 804,414 manually labeled newswire

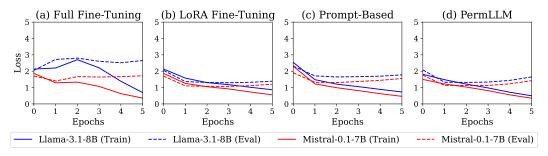


Figure 4: Comparing model loss on WMDP data set.

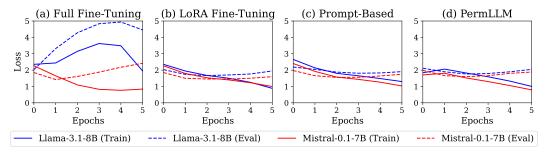


Figure 5: Comparing model loss on GPQA data set.

documents, broadly categorized with respect to three categories: *industries*, *topics* and *regions*. We took a subset of this data set and created four non-overlapping domains using *topics*: commercial (*CCAT*), economic (*ECAT*), governance (*GCAT*), and mechanical (*MCAT*). We then did 2:1 split of the subset to obtain training and test sets. The training set consists of 45622 question–answer pairs where 23822 are from *CCAT*, 7460 are from *GCAT*, 3370 are from *ECAT* and the remaining 10970 are from *MCAT*. The test set size is 22811 records, consisting of 11911 *CCAT*, 3730 *GCAT*, 1685 *ECAT*, and 5485 *MCAT* records. The largest record from this data set consists of 1199 tokens (tokenized using Llama3 tokenizer). This data set is released under CC BY 4.0 License.

PubMedQA. PubMedQA [20] contains approximately 200K medical articles formatted as (Context + Question + Answer). We encoded these articles using the GTE sentence encoder and applied k-means clustering to the resulting embeddings to derive 10 non-overlapping security domains. While clustering enforces semantic similarity within each domain and dissimilarity across domains, the underlying data distribution remains the same, since all samples originate from the same dataset. The largest record from this data set consists of 1614 tokens (tokenized using Llama3 tokenizer). This data set is released under MIT License.

D.3 Model Utility Evaluation

We use four metrics to evaluate the utility of the model generations: Bleurt Score (bluert), Bert F1-Score (bert), Sacrebleu Score (bleu) and Verbatim Accuracy (acc). These metrics measure how

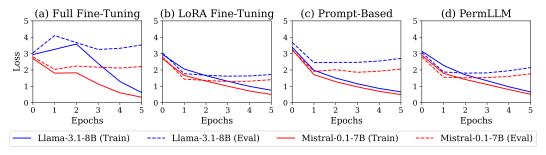


Figure 6: Comparing model loss on SimpleQA data set.

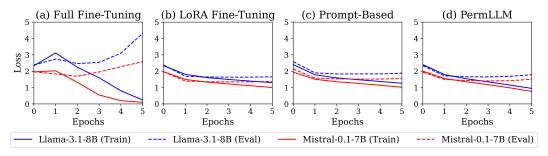


Figure 7: Comparing model loss on RCV1 data set.

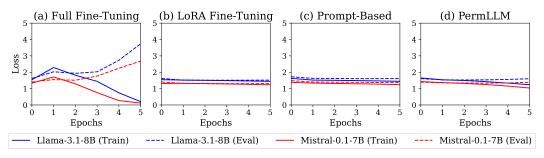


Figure 8: Comparing model loss on PubMedQA data set.

similar the generated text is to the ground truth. bleurt and bert measure the semantic similarity, bleu measures the fraction of common n-grams, and acc gives a binary decision of whether the generated text verbatim matches the ground truth. All the metrics lie in a [0,1] range, where values close to 1 indicate high model utility.

We check the utility of *Activate* to determine if tuning different LoRA adapters for each security domain leads to acceptable model utility. To that end, we show in Table 6 the utility of Llama-3.1-8B models fine-tuned on different data sets with the three approaches: full fine-tuning, LoRA fine-tuning and our PermLLM. We do not report the *bleu* score for WMDP as it is a multi-choice question-answering task where model only has to generate a single token. *bleu* requires generating at least four tokens. Our approach achieves similar or better utility on the training set compared to the LoRA approach. On the test set, our approach achieves similar utility to LoRA for most of the data sets, except for SimpleQA where LoRA performs better. This is because SimpleQA has more domains (10 in total), thus each of our individual domain adapter sees only a fraction of data of what LoRA approach's adapter sees (given that SimpleQA is already a small data set). We expect the performance of our domain-specific adapters to increase as the data set size increases. Full fine-tuning is highly sensitive to training hyper-parameters, and as a result it either completely overfits on training set to achieve high utility (e.g., on SimpleQA and RCV1), or it underfits and achieves low utility (e.g., on WMDP and GPOA). We observe similar results for Mistral-0.1-7B models (see Table 7).

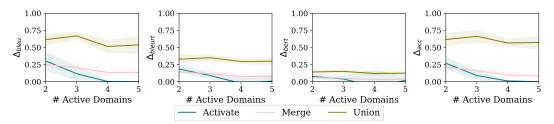


Figure 9: Utility Gap Index, Δ_U (mean \pm std) for Mistral-0.1-7B models fine-tuned on SimpleQA when user has access to multiple security domains.

Table 6: Utility comparison of Llama-3.1-8B models trained with different approaches. All reported values are $mean \pm std$ across domains.

	Metric	Full Fine	e-Tuning	LoRA Fi	ne-Tuning	Pern	nLLM
		Train	Test	Train	Test	Train	Test
ЭР	bleurt	0.74 ± 0.06	0.74 ± 0.06	0.90 ± 0.08	0.85 ± 0.08	0.92 ± 0.08	0.82 ± 0.06
WMDP	bert	0.89 ± 0.03	0.89 ± 0.03	0.96 ± 0.03	0.94 ± 0.03	0.97 ± 0.03	0.93 ± 0.03
<u> </u>	acc	0.26 ± 0.07	0.27 ± 0.07	0.76 ± 0.20	0.60 ± 0.20	0.84 ± 0.22	0.49 ± 0.15
_	bleu	0.26 ± 0.02	0.05 ± 0.03	0.45 ± 0.12	0.10 ± 0.05	0.39 ± 0.20	0.10 ± 0.04
GPQA	bleurt	0.53 ± 0.05	0.39 ± 0.05	0.64 ± 0.09	0.46 ± 0.07	0.62 ± 0.11	0.47 ± 0.07
5	bert	0.67 ± 0.06	0.59 ± 0.05	0.77 ± 0.08	0.67 ± 0.05	0.75 ± 0.09	0.67 ± 0.05
	acc	0.24 ± 0.06	0.02 ± 0.03	0.32 ± 0.05	0.05 ± 0.05	0.31 ± 0.09	0.04 ± 0.05
Ą	bleu	0.80 ± 0.06	0.34 ± 0.11	0.65 ± 0.06	0.29 ± 0.08	0.67 ± 0.10	0.09 ± 0.04
SimpleQA	bleurt	0.86 ± 0.03	0.58 ± 0.05	0.80 ± 0.02	0.61 ± 0.02	0.82 ± 0.04	0.53 ± 0.04
np	bert	0.96 ± 0.01	0.84 ± 0.02	0.94 ± 0.01	0.86 ± 0.01	0.95 ± 0.02	0.82 ± 0.03
Sir	acc	0.68 ± 0.10	0.20 ± 0.12	0.52 ± 0.07	0.17 ± 0.07	0.55 ± 0.13	0.02 ± 0.02
	bleu	0.75 ± 0.08	0.14 ± 0.08	0.22 ± 0.10	0.16 ± 0.08	0.27 ± 0.10	0.16 ± 0.08
RCV1	bleurt	0.88 ± 0.04	0.46 ± 0.12	0.57 ± 0.13	0.49 ± 0.11	0.62 ± 0.13	0.50 ± 0.12
\lesssim	bert	0.94 ± 0.03	0.67 ± 0.09	0.75 ± 0.08	0.70 ± 0.07	0.78 ± 0.08	0.70 ± 0.08
-	acc	0.78 ± 0.06	0.16 ± 0.10	0.27 ± 0.14	0.17 ± 0.10	0.31 ± 0.15	0.18 ± 0.10
- V	bleu	0.71 ± 0.05	0.07 ± 0.01	0.09 ± 0.01	0.09 ± 0.01	0.10 ± 0.02	0.09 ± 0.01
PubMedQA	bleurt	0.77 ± 0.03	0.38 ± 0.01	0.40 ± 0.01	0.40 ± 0.01	0.42 ± 0.01	0.40 ± 0.02
ρWe	bert	0.90 ± 0.02	0.64 ± 0.02	0.68 ± 0.01	0.68 ± 0.01	0.69 ± 0.02	0.67 ± 0.02
Pu	acc	-	-	-	-	-	-

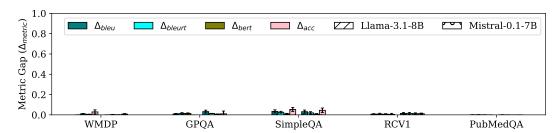


Figure 10: Utility Gap Index, Δ_U ($mean \pm std$) for prompt-based access control baseline when user has access to one security domain.

E Prompt-Based Access Control

Recent works [8, 25] have proposed enforcing some form of access control in system prompts, however we note that they do not provide absolute access control and are vulnerable to jailbreaking prompts. Regardless, we implement prompt-based access control as a baseline where each query is tagged with a prompt prefix (e.g., "use domain 1") and the rest of the fine-tuning pipeline is similar to LoRA fine-tuning. We add the relevant prompt prefixes during both model fine-tuning and inference. The models fine-tuned with prompt-based access control achieve similar training and test loss to that of LoRA fine-tuning across all the data sets, as shown in Figure 4, Figure 5, Figure 6, Figure 7, and Figure 8. However, this baseline fails to provide any meaningful access control, even when a user has access to only one security domain as shown in Figure 10 and Table 8. As shown in the figure and table, the utility gap index is close to zero and DDI scores are close to random guessing across all the data sets for both Llama and Mistral models fine-tuned with prompt-based access control. The reason is that the prompt prefix for different domains only differ in one or two tokens and hence the model tends to ignore this difference and continues generating responses even for domains the user has no access to. Exploring different prompt structures might lead to better access control but is beyond the scope of this work. We observe a similar trend when the user has access to multiple security domains as shown in Figure 11 for SimpleQA data set.

Table 7: Utility comparison of Mistral-0.1-7B models trained with different approaches. All reported values are $mean \pm std$ across domains.

	Metric	Full Fine	e-Tuning	LoRA Fi	ne-Tuning	Perm	nLLM
		Train	Test	Train	Test	Train	Test
WMDP	$egin{aligned} bleurt \ acc \end{aligned}$	0.95 ± 0.01 0.98 ± 0.01 0.88 ± 0.04	0.82 ± 0.03 0.92 ± 0.02 0.46 ± 0.14	0.96 ± 0.02 0.99 ± 0.01 0.92 ± 0.07	0.87 ± 0.03 0.94 ± 0.02 0.60 ± 0.09	0.96 ± 0.01 0.99 ± 0.01 0.93 ± 0.04	0.86 ± 0.03 0.94 ± 0.02 0.58 ± 0.11
GPQA	$bleu \ bleurt \ bert \ acc$	$\begin{array}{c} 0.46 \pm 0.03 \\ 0.65 \pm 0.04 \\ 0.75 \pm 0.05 \\ 0.38 \pm 0.04 \end{array}$	$\begin{array}{c} 0.06 \pm 0.05 \\ 0.42 \pm 0.08 \\ 0.62 \pm 0.07 \\ 0.04 \pm 0.05 \end{array}$	0.35 ± 0.08 0.59 ± 0.09 0.73 ± 0.08 0.24 ± 0.04	$\begin{array}{c} 0.11 \pm 0.07 \\ 0.47 \pm 0.06 \\ 0.68 \pm 0.05 \\ 0.05 \pm 0.06 \end{array}$	0.55 ± 0.18 0.67 ± 0.09 0.79 ± 0.08 0.40 ± 0.09	$\begin{array}{c} 0.13 \pm 0.06 \\ 0.47 \pm 0.08 \\ 0.66 \pm 0.09 \\ 0.08 \pm 0.02 \end{array}$
SimpleQA	$egin{array}{l} bleu \ bleurt \ bert \ acc \end{array}$	0.94 ± 0.02 0.94 ± 0.01 0.99 ± 0.01 0.91 ± 0.04	$\begin{array}{c} 0.36 \pm 0.11 \\ 0.60 \pm 0.04 \\ 0.85 \pm 0.02 \\ 0.23 \pm 0.12 \end{array}$	0.73 ± 0.06 0.84 ± 0.03 0.96 ± 0.01 0.62 ± 0.08	0.34 ± 0.09 0.62 ± 0.03 0.87 ± 0.01 0.20 ± 0.10	0.70 ± 0.13 0.83 ± 0.06 0.95 ± 0.03 0.60 ± 0.16	$\begin{array}{c} 0.10 \pm 0.04 \\ 0.52 \pm 0.04 \\ 0.82 \pm 0.03 \\ 0.03 \pm 0.02 \end{array}$
RCV1	$egin{array}{c} bleu \ bleurt \ bert \ cc \end{array}$	0.92 ± 0.06 0.93 ± 0.02 0.98 ± 0.02 0.92 ± 0.03	$\begin{array}{c} 0.17 \pm 0.09 \\ 0.48 \pm 0.12 \\ 0.69 \pm 0.08 \\ 0.19 \pm 0.11 \end{array}$	0.28 ± 0.13 0.60 ± 0.13 0.78 ± 0.09 0.31 ± 0.15	$\begin{array}{c} 0.20 \pm 0.10 \\ 0.51 \pm 0.12 \\ 0.71 \pm 0.08 \\ 0.20 \pm 0.11 \end{array}$	$\begin{array}{c} 0.37 \pm 0.14 \\ 0.66 \pm 0.12 \\ 0.81 \pm 0.08 \\ 0.38 \pm 0.17 \end{array}$	$\begin{array}{c} 0.19 \pm 0.09 \\ 0.50 \pm 0.12 \\ 0.71 \pm 0.08 \\ 0.19 \pm 0.10 \end{array}$
PubMedQA	$egin{array}{c} bleu \ bleurt \ bert \ acc \end{array}$	0.75 ± 0.04 0.80 ± 0.03 0.92 ± 0.01	0.08 ± 0.01 0.39 ± 0.01 0.65 ± 0.02	0.09 ± 0.01 0.41 ± 0.01 0.69 ± 0.01	0.08 ± 0.01 0.41 ± 0.01 0.68 ± 0.02	0.11 ± 0.02 0.43 ± 0.02 0.70 ± 0.02	0.08 ± 0.01 0.41 ± 0.01 0.68 ± 0.01
1.00			.00	1.00		1.00	

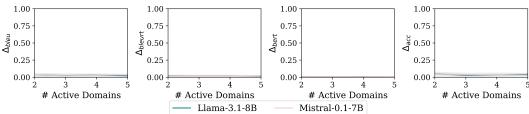


Figure 11: Utility Gap Index, Δ_U ($mean \pm std$) for prompt-based access control baseline on different models fine-tuned on SimpleQA when user has access to multiple security domains.

F MIAs against LLMs

In Section 4, we defined the Domain Distinguishability Index (DDI) as the average success rate of an adversary playing the Domain Distinguishability game over all domain set pairs. That game is implemented with *membership inference attacks* (MIAs) [44, 5, 29, 36, 47]: the auditor compares a *member* set drawn from the active domain's training data with a *non-member* set drawn from some other domain, and tries to tell them apart. The better this separation, the larger the DDI. Here, in this section, we expand on the MIA toolbox that underpins DDI—detailing evaluation metrics and the specific attacks we deploy against LLMs. More generally, an MIA for an LLM f assigns a *membership score* A(x, f) to a candidate text x. Thresholding this score at ε declares x a member (if $A(x, f) \ge \varepsilon$) or a non-member (if $A(x, f) < \varepsilon$).

F.1 Metrics

We employ two complementary metrics to quantify the success of our membership inference attacks, as used by prior MIA works [18, 6, 28]:

(1) Attack ROC curves: The Receiver Operating Characteristic (ROC) curve illustrates the trade-off between the True Positive Rate (TPR) and the False Positive Rate (FPR) for the attacks. The FPR measures the proportion of non-member samples that are incorrectly classified as members, while the TPR represents the proportion of member samples that are correctly identified as members. We report

Table 8: DDI values for prompt-based access control baseline when user has access to one security domain.

	MIA	auc-roc	Llama-3.1-8B tpr@1%fpr	tpr@5%fpr	auc-roc	Mistral-0.1-7B tpr@1%fpr	tpr@5%fpr
WMDP	Loss ZLIB Mink Mink++ Ref	$ \begin{vmatrix} 0.53 \pm 0.02 \\ 0.52 \pm 0.01 \\ 0.53 \pm 0.03 \\ 0.55 \pm 0.06 \\ 0.53 \pm 0.02 \end{vmatrix} $	$\begin{array}{c} 0.02 \pm 0.01 \\ 0.01 \pm 0.01 \\ 0.02 \pm 0.02 \\ 0.02 \pm 0.03 \\ 0.02 \pm 0.01 \end{array}$	$\begin{array}{c} 0.06 \pm 0.01 \\ 0.06 \pm 0.01 \\ 0.08 \pm 0.03 \\ 0.08 \pm 0.05 \\ 0.06 \pm 0.00 \end{array}$	$ \begin{array}{c} 0.54 \pm 0.02 \\ 0.52 \pm 0.01 \\ 0.53 \pm 0.01 \\ 0.52 \pm 0.03 \\ 0.53 \pm 0.01 \end{array} $	$\begin{array}{c} 0.02 \pm 0.01 \\ 0.01 \pm 0.01 \\ 0.02 \pm 0.01 \\ 0.01 \pm 0.00 \\ 0.01 \pm 0.00 \\ 0.01 \pm 0.01 \end{array}$	$\begin{array}{c} 0.07 \pm 0.02 \\ 0.06 \pm 0.01 \\ 0.06 \pm 0.01 \\ 0.05 \pm 0.01 \\ 0.06 \pm 0.01 \end{array}$
GPQA	Loss ZLIB Mink Mink++ Ref	$ \begin{vmatrix} 0.55 \pm 0.02 \\ 0.54 \pm 0.02 \\ 0.57 \pm 0.05 \\ 0.54 \pm 0.10 \\ 0.57 \pm 0.02 \end{vmatrix} $	$\begin{array}{c} 0.02 \pm 0.00 \\ 0.02 \pm 0.00 \\ 0.02 \pm 0.01 \\ 0.05 \pm 0.06 \\ 0.04 \pm 0.02 \end{array}$	$\begin{array}{c} 0.06 \pm 0.10 \\ 0.07 \pm 0.01 \\ 0.12 \pm 0.02 \\ 0.12 \pm 0.08 \\ 0.13 \pm 0.05 \end{array}$	$\begin{array}{c} 0.56 \pm 0.03 \\ 0.54 \pm 0.02 \\ 0.59 \pm 0.07 \\ 0.55 \pm 0.12 \\ 0.56 \pm 0.03 \end{array}$	$\begin{array}{c} 0.03 \pm 0.01 \\ 0.03 \pm 0.01 \\ 0.05 \pm 0.06 \\ 0.06 \pm 0.06 \\ 0.03 \pm 0.02 \end{array}$	$\begin{array}{c} 0.12 \pm 0.04 \\ 0.08 \pm 0.02 \\ 0.13 \pm 0.07 \\ 0.12 \pm 0.09 \\ 0.13 \pm 0.07 \end{array}$
SimpleQA	Loss ZLIB Mink Mink++ Ref	$ \begin{vmatrix} 0.53 \pm 0.25 \\ 0.52 \pm 0.16 \\ 0.52 \pm 0.28 \\ 0.50 \pm 0.43 \\ 0.53 \pm 0.22 \end{vmatrix} $	$\begin{array}{c} 0.08 \pm 0.14 \\ 0.04 \pm 0.05 \\ 0.09 \pm 0.15 \\ 0.31 \pm 0.40 \\ 0.03 \pm 0.03 \end{array}$	$\begin{array}{c} 0.16 \pm 0.20 \\ 0.09 \pm 0.09 \\ 0.17 \pm 0.22 \\ 0.36 \pm 0.43 \\ 0.11 \pm 0.12 \end{array}$	$ \begin{array}{c} 0.55 \pm 0.22 \\ 0.53 \pm 0.14 \\ 0.55 \pm 0.22 \\ 0.52 \pm 0.35 \\ 0.54 \pm 0.15 \end{array} $	$\begin{array}{c} 0.09 \pm 0.15 \\ 0.03 \pm 0.03 \\ 0.09 \pm 0.15 \\ 0.22 \pm 0.33 \\ 0.04 \pm 0.06 \end{array}$	$\begin{array}{c} 0.16 \pm 0.20 \\ 0.09 \pm 0.08 \\ 0.17 \pm 0.20 \\ 0.28 \pm 0.35 \\ 0.09 \pm 0.09 \end{array}$
RCV1	Loss ZLIB Mink Mink++ Ref	$ \begin{vmatrix} 0.50 \pm 0.02 \\ 0.50 \pm 0.01 \\ 0.50 \pm 0.04 \\ 0.50 \pm 0.05 \\ 0.50 \pm 0.01 \end{vmatrix} $	$\begin{array}{c} 0.01 \pm 0.00 \\ 0.01 \pm 0.00 \\ 0.01 \pm 0.00 \\ 0.01 \pm 0.00 \\ 0.01 \pm 0.01 \\ 0.01 \pm 0.01 \end{array}$	$\begin{array}{c} 0.05 \pm 0.01 \\ 0.05 \pm 0.02 \\ 0.05 \pm 0.01 \\ 0.05 \pm 0.01 \\ 0.05 \pm 0.01 \\ 0.05 \pm 0.01 \end{array}$	$\begin{array}{c} 0.50 \pm 0.01 \\ 0.50 \pm 0.00 \\ 0.50 \pm 0.01 \\ 0.50 \pm 0.01 \\ 0.50 \pm 0.04 \\ 0.50 \pm 0.01 \end{array}$	$\begin{array}{c} 0.01 \pm 0.00 \\ 0.01 \pm 0.00 \\ 0.01 \pm 0.02 \\ 0.01 \pm 0.00 \\ 0.01 \pm 0.00 \\ 0.01 \pm 0.00 \end{array}$	$\begin{array}{c} 0.05 \pm 0.00 \\ 0.05 \pm 0.00 \\ 0.05 \pm 0.01 \\ 0.05 \pm 0.01 \\ 0.05 \pm 0.01 \\ 00.5 \pm 0.00 \end{array}$
PubMedQA	Loss ZLIB Mink Mink++ Ref	$ \begin{vmatrix} 0.50 \pm 0.00 \\ 0.50 \pm 0.00 \\ 0.50 \pm 0.01 \\ 0.50 \pm 0.01 \\ 0.50 \pm 0.03 \end{vmatrix} $	$\begin{array}{c} 0.01 \pm 0.00 \\ 0.01 \pm 0.00 \end{array}$	$\begin{array}{c} 0.05 \pm 0.00 \\ 0.05 \pm 0.01 \end{array}$	$\begin{array}{c} 0.50 \pm 0.00 \\ 0.50 \pm 0.00 \\ 0.50 \pm 0.01 \\ 0.50 \pm 0.01 \\ 0.50 \pm 0.03 \end{array}$	$\begin{array}{c} 0.01 \pm 0.00 \\ 0.01 \pm 0.00 \end{array}$	$\begin{array}{c} 0.05 \pm 0.00 \\ 0.05 \pm 0.00 \\ 0.05 \pm 0.01 \\ 0.05 \pm 0.00 \\ 0.05 \pm 0.00 \\ 0.05 \pm 0.01 \end{array}$

the Area Under the ROC Curve (AUC-ROC) as an aggregate metric to assess the overall success of the attacks. AUC-ROC is a threshold-independent metric, and it shows the probability that a positive instance (member) has higher score than a negative instance (non-member).

(2) Attack TPR at low FPR: This metric is crucial for determining the effectiveness of an attack at confidently identifying members of the training dataset without falsely classifying non-members as members. We focus on low FPR thresholds, specifically 1%, and 5%. For instance, the TPR at an FPR of 1% is calculated by setting the detection threshold so that only 1% of non-member samples are predicted as members.

F.2 Existing MIAs

LOSS [44]: The LOSS method utilizes the loss value of model f(.) for the given text x as the membership score; a lower loss suggests that the text was seen during training, so $A(x, f) = \ell(f, x)$.

Ref [5]: Calculating membership scores based solely on loss values often results in high false negative rates. To improve this, a difficulty calibration method can be employed to account for the intrinsic complexity of x. For example, repetitive or common phrases typically yield low loss values. One method of calibrating this input complexity is by using another LLM, Ref(.), assumed to be trained on a similar data distribution. The membership score is then defined as the difference in loss values between the target and reference models, $A(x, f) = \ell(x, f) - \ell(x, Ref)$. In our evaluations, we used the base models (i.e., Llama-3.1-8B and Mistral-0.1-7B) before any fine-tuning as the reference models.

Zlib [5]: Another method to calibrate the difficulty of a sample is by using its zlib compression size, where more complex sentences have higher compression sizes. The membership score is then calculated by normalizing the loss value by the zlib compression size, $A(x, f) = \frac{\ell(x, f)}{7lib(x)}$.

Table 9: DDI values for models (with base model Llama-3.1-8B) with different approaches of combining domains when user has access to two domains. All reported values are $mean \pm std$ across domains

	MIA	auc-roc	Activate tpr@1%fpr	tpr@5%fpr	auc-roc	Merge tpr@1%fpr	tpr@5%fpr	auc-roc	Union tpr@1%fpr	tpr@5%fpr
WMDP	Loss ZLIB Mink Mink++ Ref	$\begin{array}{c} 0.98 \pm 0.02 \\ 0.92 \pm 0.08 \\ 0.99 \pm 0.01 \\ 0.90 \pm 0.05 \\ 1.00 \pm 0.00 \end{array}$	$\begin{array}{c} 0.77 \pm 0.22 \\ 0.60 \pm 0.27 \\ 0.88 \pm 0.08 \\ 0.62 \pm 0.21 \\ 0.98 \pm 0.02 \end{array}$	$\begin{array}{c} 0.87 \pm 0.13 \\ 0.67 \pm 0.28 \\ 0.93 \pm 0.04 \\ 0.71 \pm 0.16 \\ 0.99 \pm 0.01 \end{array}$	$\begin{array}{c} 0.93 \pm 0.05 \\ 0.86 \pm 0.09 \\ 0.96 \pm 0.02 \\ 0.94 \pm 0.04 \\ 0.99 \pm 0.00 \end{array}$	$\begin{array}{c} 0.53 \pm 0.25 \\ 0.38 \pm 0.21 \\ 0.65 \pm 0.19 \\ 0.65 \pm 0.21 \\ 0.81 \pm 0.05 \end{array}$	$\begin{array}{c} 0.67 \pm 0.21 \\ 0.50 \pm 0.26 \\ 0.78 \pm 0.12 \\ 0.80 \pm 0.15 \\ 0.91 \pm 0.02 \end{array}$	$\begin{array}{c} 0.99 \pm 0.02 \\ 0.97 \pm 0.05 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$	$\begin{array}{c} 0.90 \pm 0.14 \\ 0.77 \pm 0.31 \\ 0.94 \pm 0.08 \\ 1.00 \pm 0.00 \\ 0.98 \pm 0.02 \end{array}$	$\begin{array}{c} 0.94 \pm 0.09 \\ 0.80 \pm 0.28 \\ 0.99 \pm 0.01 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$
GPQA	Loss ZLIB Mink Mink++ Ref	$\begin{array}{c} 0.99 \pm 0.01 \\ 0.90 \pm 0.06 \\ 0.99 \pm 0.01 \\ 0.95 \pm 0.06 \\ 1.00 \pm 0.00 \end{array}$	$\begin{array}{c} 0.81 \pm 0.09 \\ 0.38 \pm 0.26 \\ 0.92 \pm 0.11 \\ 0.82 \pm 0.10 \\ 0.99 \pm 0.01 \end{array}$	$\begin{array}{c} 0.93 \pm 0.05 \\ 0.63 \pm 0.22 \\ 0.97 \pm 0.04 \\ 0.85 \pm 0.13 \\ 0.99 \pm 0.01 \end{array}$	$\begin{array}{c} 0.93 \pm 0.02 \\ 0.82 \pm 0.07 \\ 0.96 \pm 0.01 \\ 0.97 \pm 0.03 \\ 0.99 \pm 0.01 \end{array}$	$\begin{array}{c} 0.38 \pm 0.14 \\ 0.26 \pm 0.17 \\ 0.69 \pm 0.07 \\ 0.75 \pm 0.13 \\ 0.87 \pm 0.12 \end{array}$	$\begin{array}{c} 0.72 \pm 0.03 \\ 0.44 \pm 0.16 \\ 0.80 \pm 0.07 \\ 0.88 \pm 0.10 \\ 0.93 \pm 0.09 \end{array}$	$\begin{array}{c} 1.00 \pm 0.00 \\ 0.99 \pm 0.01 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$	$\begin{array}{c} 0.97 \pm 0.04 \\ 0.79 \pm 0.30 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$	$\begin{array}{c} 0.99 \pm 0.01 \\ 0.96 \pm 0.05 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$
SimpleQA	Loss ZLIB Mink Mink++ Ref	$\begin{array}{c} 0.96 \pm 0.03 \\ 0.94 \pm 0.04 \\ 0.94 \pm 0.06 \\ 0.85 \pm 0.10 \\ 0.96 \pm 0.03 \end{array}$	$\begin{array}{c} 0.42 \pm 0.32 \\ 0.35 \pm 0.28 \\ 0.41 \pm 0.33 \\ 0.25 \pm 0.19 \\ 0.37 \pm 0.35 \end{array}$	$\begin{array}{c} 0.73 \pm 0.26 \\ 0.66 \pm 0.23 \\ 0.68 \pm 0.27 \\ 0.57 \pm 0.16 \\ 0.73 \pm 0.30 \end{array}$	$\begin{array}{c} 0.95 \pm 0.03 \\ 0.93 \pm 0.04 \\ 0.94 \pm 0.03 \\ 0.92 \pm 0.03 \\ 0.96 \pm 0.04 \end{array}$	$\begin{array}{c} 0.47 \pm 0.28 \\ 0.41 \pm 0.24 \\ 0.47 \pm 0.22 \\ 0.34 \pm 0.16 \\ 0.43 \pm 0.40 \end{array}$	$\begin{array}{c} 0.74 \pm 0.21 \\ 0.67 \pm 0.17 \\ 0.71 \pm 0.18 \\ 0.62 \pm 0.13 \\ 0.69 \pm 0.35 \end{array}$	$ \begin{array}{c} 0.97 \pm 0.04 \\ 0.97 \pm 0.04 \\ 0.98 \pm 0.03 \\ 0.97 \pm 0.03 \\ 0.97 \pm 0.04 \end{array} $	0.62 ± 0.38 0.61 ± 0.38 0.57 ± 0.38 0.57 ± 0.37 0.58 ± 0.42	$\begin{array}{c} 0.83 \pm 0.29 \\ 0.82 \pm 0.29 \\ 0.84 \pm 0.25 \\ 0.85 \pm 0.24 \\ 0.79 \pm 0.31 \end{array}$
RCV1	Loss ZLIB Mink Mink++ Ref	$\begin{array}{c} 0.96 \pm 0.02 \\ 0.82 \pm 0.02 \\ 0.97 \pm 0.02 \\ 0.80 \pm 0.13 \\ 0.97 \pm 0.01 \end{array}$	$\begin{array}{c} 0.40 \pm 0.09 \\ 0.27 \pm 0.07 \\ 0.60 \pm 0.14 \\ 0.32 \pm 0.19 \\ 0.50 \pm 0.09 \end{array}$	$\begin{array}{c} 0.76 \pm 0.15 \\ 0.46 \pm 0.06 \\ 0.87 \pm 0.08 \\ 0.49 \pm 0.24 \\ 0.86 \pm 0.09 \end{array}$	$\begin{array}{c} 0.90 \pm 0.01 \\ 0.72 \pm 0.02 \\ 0.92 \pm 0.02 \\ 0.84 \pm 0.07 \\ 0.95 \pm 0.00 \end{array}$	$\begin{array}{c} 0.24 \pm 0.05 \\ 0.11 \pm 0.03 \\ 0.29 \pm 0.04 \\ 0.28 \pm 0.22 \\ 0.26 \pm 0.07 \end{array}$	$\begin{array}{c} 0.52 \pm 0.07 \\ 0.28 \pm 0.03 \\ 0.65 \pm 0.08 \\ 0.52 \pm 0.19 \\ 0.63 \pm 0.05 \end{array}$	$\begin{array}{c} 0.98 \pm 0.00 \\ 0.90 \pm 0.05 \\ 0.99 \pm 0.00 \\ 0.99 \pm 0.00 \\ 0.98 \pm 0.01 \end{array}$	$\begin{array}{c} 0.55 \pm 0.23 \\ 0.52 \pm 0.20 \\ 0.80 \pm 0.08 \\ 0.90 \pm 0.05 \\ 0.50 \pm 0.31 \end{array}$	$\begin{array}{c} 0.94 \pm 0.01 \\ 0.67 \pm 0.13 \\ 0.97 \pm 0.01 \\ 0.98 \pm 0.00 \\ 0.95 \pm 0.02 \end{array}$
PubMedQA	Loss ZLIB Mink Mink++ Ref	$\begin{array}{c} 0.84 \pm 0.05 \\ 0.78 \pm 0.05 \\ 0.91 \pm 0.05 \\ 0.85 \pm 0.10 \\ 0.99 \pm 0.02 \end{array}$	$\begin{array}{c} 0.21 \pm 0.11 \\ 0.13 \pm 0.07 \\ 0.39 \pm 0.17 \\ 0.36 \pm 0.22 \\ 0.84 \pm 0.16 \end{array}$	$\begin{array}{c} 0.44 \pm 0.12 \\ 0.33 \pm 0.09 \\ 0.63 \pm 0.17 \\ 0.55 \pm 0.23 \\ 0.93 \pm 0.08 \end{array}$	$ \begin{vmatrix} 0.66 \pm 0.02 \\ 0.61 \pm 0.01 \\ 0.73 \pm 0.02 \\ 0.79 \pm 0.05 \\ 0.99 \pm 0.01 \end{vmatrix} $	$\begin{array}{c} 0.04 \pm 0.01 \\ 0.03 \pm 0.01 \\ 0.07 \pm 0.02 \\ 0.12 \pm 0.05 \\ 0.75 \pm 0.18 \end{array}$	$\begin{array}{c} 0.14 \pm 0.01 \\ 0.12 \pm 0.01 \\ 0.22 \pm 0.03 \\ 0.33 \pm 0.08 \\ 0.96 \pm 0.04 \end{array}$	$ \begin{array}{c} 0.79 \pm 0.04 \\ 0.72 \pm 0.04 \\ 0.87 \pm 0.03 \\ 0.94 \pm 0.02 \\ 1.00 \pm 0.00 \\ \end{array} $	$\begin{array}{c} 0.14 \pm 0.05 \\ 0.09 \pm 0.04 \\ 0.25 \pm 0.11 \\ 0.44 \pm 0.17 \\ 0.97 \pm 0.07 \end{array}$	$\begin{array}{c} 0.32 \pm 0.09 \\ 0.25 \pm 0.07 \\ 0.49 \pm 0.11 \\ 0.72 \pm 0.12 \\ 1.00 \pm 0.00 \end{array}$

Table 10: DDI values for models (with base model Mistral-0.1-7B) with different approaches of combining domains when user has access to two domains. All reported values are $mean \pm std$ across domains.

	MIA	auc-roc	Activate tpr@1%fpr	tpr@5%fpr	auc-roc	Merge tpr@1%fpr	tpr@5%fpr	auc-roc	Union tpr@1%fpr	tpr@5%fpr
WMDP	Loss ZLIB Mink Mink++ Ref	$\begin{array}{c} 0.99 \pm 0.02 \\ 0.93 \pm 0.09 \\ 0.99 \pm 0.01 \\ 0.96 \pm 0.02 \\ 1.00 \pm 0.00 \end{array}$	$\begin{array}{c} 0.85 \pm 0.21 \\ 0.69 \pm 0.30 \\ 0.89 \pm 0.14 \\ 0.77 \pm 0.04 \\ 1.00 \pm 0.00 \end{array}$	$\begin{array}{c} 0.92 \pm 0.11 \\ 0.74 \pm 0.30 \\ 0.95 \pm 0.07 \\ 0.86 \pm 0.04 \\ 1.00 \pm 0.00 \end{array}$	$ \begin{vmatrix} 0.95 \pm 0.04 \\ 0.87 \pm 0.09 \\ 0.96 \pm 0.03 \\ 0.94 \pm 0.03 \\ 0.99 \pm 0.00 \end{vmatrix} $	$\begin{array}{c} 0.62 \pm 0.21 \\ 0.47 \pm 0.26 \\ 0.73 \pm 0.11 \\ 0.58 \pm 0.03 \\ 0.86 \pm 0.09 \end{array}$	$\begin{array}{c} 0.73 \pm 0.19 \\ 0.58 \pm 0.29 \\ 0.83 \pm 0.12 \\ 0.80 \pm 0.05 \\ 0.96 \pm 0.02 \end{array}$	$ \begin{vmatrix} 0.99 \pm 0.01 \\ 0.98 \pm 0.03 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{vmatrix} $	$\begin{array}{c} 0.93 \pm 0.10 \\ 0.83 \pm 0.23 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$	$\begin{array}{c} 0.96 \pm 0.06 \\ 0.88 \pm 0.16 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$
GPQA	Loss ZLIB Mink Mink++ Ref	$\begin{array}{c} 0.99 \pm 0.01 \\ 0.93 \pm 0.08 \\ 1.00 \pm 0.00 \\ 0.98 \pm 0.02 \\ 1.00 \pm 0.00 \end{array}$	$\begin{array}{c} 0.83 \pm 0.18 \\ 0.50 \pm 0.35 \\ 0.94 \pm 0.07 \\ 0.80 \pm 0.14 \\ 1.00 \pm 0.00 \end{array}$	$\begin{array}{c} 0.95 \pm 0.06 \\ 0.74 \pm 0.32 \\ 0.98 \pm 0.02 \\ 0.92 \pm 0.06 \\ 1.00 \pm 0.00 \end{array}$	$ \begin{vmatrix} 0.96 \pm 0.04 \\ 0.86 \pm 0.09 \\ 0.98 \pm 0.02 \\ 0.98 \pm 0.01 \\ 0.99 \pm 0.02 \end{vmatrix} $	$\begin{array}{c} 0.55 \pm 0.24 \\ 0.33 \pm 0.23 \\ 0.74 \pm 0.14 \\ 0.75 \pm 0.13 \\ 0.84 \pm 0.23 \end{array}$	$\begin{array}{c} 0.87 \pm 0.06 \\ 0.56 \pm 0.21 \\ 0.87 \pm 0.12 \\ 0.89 \pm 0.07 \\ 0.97 \pm 0.04 \end{array}$	$ \begin{vmatrix} 1.00 \pm 0.00 \\ 0.99 \pm 0.01 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{vmatrix} $	$\begin{array}{c} 0.97 \pm 0.04 \\ 0.88 \pm 0.17 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \\ 0.97 \pm 0.04 \end{array}$	$\begin{array}{c} 0.98 \pm 0.02 \\ 0.97 \pm 0.04 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$
SimpleQA	Loss ZLIB Mink Mink++ Ref	$\begin{array}{c} 0.97 \pm 0.03 \\ 0.97 \pm 0.03 \\ 0.97 \pm 0.03 \\ 0.97 \pm 0.03 \\ 0.92 \pm 0.04 \\ 0.98 \pm 0.03 \end{array}$	$\begin{array}{c} 0.58 \pm 0.33 \\ 0.51 \pm 0.32 \\ 0.51 \pm 0.34 \\ 0.46 \pm 0.21 \\ 0.65 \pm 0.39 \end{array}$	$\begin{array}{c} 0.82 \pm 0.27 \\ 0.78 \pm 0.28 \\ 0.83 \pm 0.24 \\ 0.68 \pm 0.21 \\ 0.86 \pm 0.27 \end{array}$	$ \begin{vmatrix} 0.96 \pm 0.02 \\ 0.95 \pm 0.03 \\ 0.96 \pm 0.02 \\ 0.93 \pm 0.05 \\ 0.98 \pm 0.03 \end{vmatrix} $	$\begin{array}{c} 0.49 \pm 0.24 \\ 0.44 \pm 0.23 \\ 0.49 \pm 0.24 \\ 0.45 \pm 0.28 \\ 0.64 \pm 0.34 \end{array}$	$\begin{array}{c} 0.79 \pm 0.17 \\ 0.72 \pm 0.19 \\ 0.77 \pm 0.18 \\ 0.73 \pm 0.19 \\ 0.85 \pm 0.25 \end{array}$	$ \begin{vmatrix} 0.97 \pm 0.04 \\ 0.96 \pm 0.04 \end{vmatrix} $	$\begin{array}{c} 0.50 \pm 0.42 \\ 0.51 \pm 0.42 \\ 0.51 \pm 0.41 \\ 0.50 \pm 0.41 \\ 0.48 \pm 0.43 \end{array}$	$\begin{array}{c} 0.76 \pm 0.31 \\ 0.75 \pm 0.31 \\ 0.79 \pm 0.27 \\ 0.76 \pm 0.29 \\ 0.73 \pm 0.34 \end{array}$
RCVI	Loss ZLIB Mink Mink++ Ref	$\begin{array}{c} 0.93 \pm 0.04 \\ 0.82 \pm 0.05 \\ 0.93 \pm 0.05 \\ 0.69 \pm 0.25 \\ 0.96 \pm 0.02 \end{array}$	$\begin{array}{c} 0.39 \pm 0.23 \\ 0.30 \pm 0.10 \\ 0.44 \pm 0.24 \\ 0.27 \pm 0.20 \\ 0.35 \pm 0.12 \end{array}$	$\begin{array}{c} 0.62 \pm 0.23 \\ 0.50 \pm 0.08 \\ 0.68 \pm 0.23 \\ 0.45 \pm 0.33 \\ 0.71 \pm 0.18 \end{array}$	$ \begin{vmatrix} 0.85 \pm 0.01 \\ 0.69 \pm 0.03 \\ 0.85 \pm 0.02 \\ 0.70 \pm 0.16 \\ 0.94 \pm 0.01 \end{vmatrix} $	$\begin{array}{c} 0.14 \pm 0.03 \\ 0.10 \pm 0.04 \\ 0.16 \pm 0.03 \\ 0.18 \pm 0.13 \\ 0.15 \pm 0.05 \end{array}$	$\begin{array}{c} 0.35 \pm 0.02 \\ 0.26 \pm 0.06 \\ 0.40 \pm 0.04 \\ 0.35 \pm 0.21 \\ 0.52 \pm 0.08 \end{array}$	$ \begin{vmatrix} 0.98 \pm 0.01 \\ 0.90 \pm 0.05 \\ 0.99 \pm 0.00 \\ 0.99 \pm 0.00 \\ 0.98 \pm 0.00 \end{vmatrix} $	$\begin{array}{c} 0.53 \pm 0.22 \\ 0.48 \pm 0.23 \\ 0.73 \pm 0.12 \\ 0.89 \pm 0.03 \\ 0.45 \pm 0.25 \end{array}$	$\begin{array}{c} 0.92 \pm 0.01 \\ 0.67 \pm 0.14 \\ 0.97 \pm 0.01 \\ 0.98 \pm 0.00 \\ 0.97 \pm 0.00 \end{array}$
PubMedQA	Loss ZLIB Mink Mink++ Ref	$\begin{array}{c} 0.78 \pm 0.11 \\ 0.73 \pm 0.10 \\ 0.83 \pm 0.14 \\ 0.72 \pm 0.24 \\ 0.93 \pm 0.10 \\ \end{array}$	$\begin{array}{c} 0.19 \pm 0.15 \\ 0.13 \pm 0.10 \\ 0.30 \pm 0.26 \\ 0.29 \pm 0.31 \\ 0.69 \pm 0.30 \end{array}$	$\begin{array}{c} 0.38 \pm 0.19 \\ 0.28 \pm 0.14 \\ 0.49 \pm 0.27 \\ 0.44 \pm 0.33 \\ 0.80 \pm 0.25 \end{array}$	$ \begin{vmatrix} 0.65 \pm 0.01 \\ 0.61 \pm 0.01 \\ 0.70 \pm 0.02 \\ 0.74 \pm 0.06 \\ 0.98 \pm 0.00 \end{vmatrix} $	$\begin{array}{c} 0.04 \pm 0.01 \\ 0.03 \pm 0.01 \\ 0.05 \pm 0.02 \\ 0.08 \pm 0.04 \\ 0.59 \pm 0.09 \end{array}$	$\begin{array}{c} 0.13 \pm 0.02 \\ 0.11 \pm 0.01 \\ 0.18 \pm 0.02 \\ 0.23 \pm 0.07 \\ 0.93 \pm 0.04 \end{array}$	$ \begin{vmatrix} 0.81 \pm 0.05 \\ 0.74 \pm 0.05 \\ 0.89 \pm 0.04 \\ 0.95 \pm 0.02 \\ 1.00 \pm 0.00 \end{vmatrix} $	$\begin{array}{c} 0.16 \pm 0.09 \\ 0.11 \pm 0.04 \\ 0.26 \pm 0.11 \\ 0.53 \pm 0.19 \\ 0.97 \pm 0.05 \end{array}$	$\begin{array}{c} 0.36 \pm 0.12 \\ 0.27 \pm 0.18 \\ 0.53 \pm 0.12 \\ 0.76 \pm 0.13 \\ 1.00 \pm 0.00 \end{array}$

Min-K [36]: This attack hypothesizes that non-member samples often have more tokens assigned lower likelihoods. It first calculates the likelihood of each token as Min-K% $_{\text{token}}(x_t) = \log p(x_t|x_{< t})$, for each token x_t given the prefix $x_{< t}$. The membership score is then calculated by averaging over the lowest K% of tokens with lower likelihood, $A(x,f) = \frac{1}{|\min - k\%|} \sum_{x_i \in min - k\%} \text{Min-K}\%_{\text{token}}(x_t)$.

Min-K++ [47]: This method improves on Min-K by utilizing the insight that maximum likelihood training optimizes the Hessian trace of likelihood over the training data. It calculates a normalized score for each token x_t given the prefix $x_{< t}$ as Min-K%++ $_{\text{token}}(x_t) = \frac{\log p(x_t|x_{< t}) - \mu_{x_{< t}}}{\sigma_{x_{< t}}}$, where $\mu_{x_{< t}}$ is the mean log probability of the next token across the vocabulary, and $\sigma_{x_{< t}}$ is the standard deviation. The membership score is then aggregated by averaging the scores of the lowest K% tokens, $A(x,f) = \frac{1}{|\min - k\% + 1|} \sum_{x_t \in min - k\%} \text{Min-K}\% + +_{\text{token}}(x_t)$.