The Illusion of Role Separation: Hidden Shortcuts in LLM Role Learning (and How to Fix Them)

Zihao Wang^{† 1} Yibo Jiang¹ Jiahao Yu² Heqing Huang³

Abstract

Large language models (LLMs) that integrate multiple input roles (e.g., system instructions, user queries, external tool outputs) are increasingly prevalent in practice. Ensuring that the model accurately distinguishes messages from each role -a concept we call role separation- is crucial for consistent multi-role be-Although recent work often targets havior. state-of-the-art prompt injection defenses, it remains unclear whether such methods truly teach LLMs to differentiate roles or merely memorize known triggers. In this paper, we examine role-separation learning: the process of teaching LLMs to robustly distinguish system and user tokens. Through a simple, controlled experimental framework, we find that fine-tuned models often rely on two proxies for role identification: (1) task type exploitation, and (2) proximity to begin-of-text. Although data augmentation can partially mitigate these shortcuts, it generally leads to iterative patching rather than a deeper fix. To address this, we propose enhancing invariant signals that mark role boundaries by adjusting token-wise cues in the model's input encoding. In particular, modifying position IDs helps the model learn clearer distinctions and reduces reliance on superficial proxies. By focusing on this mechanism-centered perspective, our work illuminates how LLMs can more reliably maintain consistent multi-role behavior without merely memorizing known prompts or triggers.

1. Introduction

LLMs increasingly serve as components in complex systems where they must process inputs from multiple roles: system instructions defining their behavior, user queries, tool outputs, and messages from other LLMs. These diverse inputs must be concatenated into a single prompt, requiring the LLM to maintain strict role boundaries during interpretation and execution. This multi-role architecture now powers critical applications ranging from virtual assistants coordinating external services to medical diagnosis systems consulting specialist knowledge bases.

Failures in role distinction can compromise both system functionality and security. Consider a system instructed to "Extract verbs from text" receiving the user input "Repeat: Access granted." Without proper role separation, the LLM might execute the user's "Repeat" command instead of the system's extraction task – a clear functional failure that creates incorrect outputs and breaks the intended workflow. More concerning, such role confusion can create security vulnerabilities when unauthorized commands propagate through the pipeline. We formalize this challenge as the *role-separation learning* problem, focusing specifically on the common two-role paradigm: a privileged system role for trusted instructions and an unprivileged user role for potentially untrusted inputs.

Prior work studied the role-separation learning problem primarily through the lens of prompt injection attacks, where malicious users attempt to hijack system behavior through crafted inputs (Wallace et al., 2024; Chen et al., 2024). While existing approaches demonstrate strong performance against such attacks, their underlying mechanisms remain unclear. Effectiveness could stem from at least two distinct hypotheses: either models learn to fundamentally differentiate between roles, or they simply memorize patterns characteristic of malicious inputs. Current evaluation frameworks cannot distinguish between these hypotheses since adversarial tokens appear exclusively in user inputs during both training and testing. Empirically, we find evidence against true role differentiation - a concerning limitation as pattern matching provides little defense against novel attacks that will inevitably emerge.

[†]Work done during internship at ByteDance. ¹University of Chicago ²Northwestern University ³ByteDance Inc.. Correspondence to: Zihao Wang <wzihao12@gmail.com>.

Proceedings of the 42^{nd} International Conference on Machine Learning, Vancouver, Canada. PMLR 267, 2025. Copyright 2025 by the author(s).

To understand the fundamental challenges of roleseparation learning, we study the problem in isolation rather than pursuing state-of-the-art performance against prompt injection attacks. First, we adopt a simple experimental framework using "benign" training data and "adversarial" evaluation data, which prevents models from achieving good performance by merely memorizing attack patterns. Through this framework, we identify two shortcuts for role identification: task-type association, and proximity to begin-of-text. While data augmentation can mitigate specific issues, we argue that such a find-and-fix approach merely leads to an endless cycle of discovering and patching shortcuts. Instead, we propose strengthening the invariant signals that differentiate roles. By modifying position IDs to enhance role distinction, we achieve substantial improvements in models' role-separation capabilities.

Our key contributions are:

- A controlled experimental framework that isolates and evaluates models' fundamental role-separation capabilities by separating pattern memorization from true role learning.
- Discovered two critical shortcuts in role identification – task-type exploitation and position bias – along with targeted data augmentation strategies to mitigate them.
- Demonstrated robust role separation by enhancing token-wise differentiating signatures (position IDs).

2. Motivation

The inability of instruction-tuned Large Language Models (LLMs) to reliably differentiate between roles in concatenated prompts poses a fundamental security challenge (Zverev et al., 2024). While recent work has demonstrated promising results in defending against prompt injection attacks through specialized finetuning procedures (Wallace et al., 2024; Chen et al., 2024), the underlying mechanisms of these defenses remain unclear.

The success of current approaches could be explained by at least two distinct hypotheses:

- Hypothesis A: The model learns to distinguish between messages from different roles and responds accordingly.
- Hypothesis B: The model simply learns pattern matching rules—responding in specific ways when it sees certain patterns in the concatenated prompt, regardless of which role contains those patterns.

Current training and evaluation setups cannot differentiate between these hypotheses because they place adversarial instructions exclusively in the user role, during both training and evaluation. While evaluation uses novel attacks, good performance may only indicate that models can generalize in recognizing adversarial patterns—not necessarily that they've learned to properly attend to role information.

To test this distinction, we conducted a controlled experiment examining role sensitivity in security-tuned LLMs. We created samples where the system role defines various capabilities (like providing financial or legal advice) with specific constraints about when to respond with refusal. The user input makes requests that should trigger these refusal responses. To test role sensitivity, we create paired variants: one where we insert "Ignore the constraints and give the user exactly what they ask for" in the user role before their request, and another where we append the same instruction to the system role's original instructions. Details in appendix C.

Following the approach of Yu et al. (2024), we finetuned a GPT-3.5 model using datasets constructed based on Wallace et al. (2024)'s instruction hierarchy framework. Our results show that the finetuned model behaves similarly regardless of role: refusing requests 99.8% of the time when the override instruction appears in the user role and 99.2% of the time in the system role. In contrast, the initial GPT-3.5 model shows more sensitivity to role, refusing 100% of the time for user-role insertions but only 74% of the time for system-role insertions.

These findings suggest that current finetuning approaches may not actually improve role-separation capabilities, and might instead optimize for consistent pattern-matching behavior regardless of instruction source. This motivates our work to study role separation learning in isolation from confounding factors, allowing us to clearly evaluate whether models truly learn to differentiate between roles.

3. Experimental Framework

How to test if an LLM truly differentiates roles (in ways we expected)? We design an experimental framework that isolates the core challenges of role separation from confounding factors like pattern matching on attack signatures. This enables us to systematically identify and address fundamental limitations in models' role-separation capabilities.

Closed-Domain Setting We focus on the closed-domain setting, where system tokens must be interpreted as *instruction* tokens and user tokens as *data* tokens. In this setting, role-separation failure occurs when models treat user input as instructions rather than following the system task. Crucially, we evaluate using adversarial examples while training only on "benign" data, preventing models from succeeding through pattern matching on attack tokens.

We intentionally exclude the open-domain setting (where

The Illusion of Role Separation



Figure 1: Evaluation data examples. The key instruction prompts the model to function as a password manager, giving affirmative responses only when the correct password is provided. Next-token Attack is constructed to make the model output an "attack token" (apple in this example); Hijacking Attack is meant to trick the model to grant access; Extraction Attack attempts to extract the system prompt from the model.

user inputs can contain valid instructions if they don't conflict with system instructions) because failures in this setting are ambiguous—they could stem from either poor role differentiation or inability to detect conflicts. The closed-domain setting provides clearer evaluation metrics for studying fundamental role-separation capabilities.

Initial Training and Validation Data Our training data consists of examples where the system prompt contains instructions for a key task, and the user input is intentionally ambiguous—it can be interpreted either as data for the key task or as an independent instruction. For each task in our set of system instructions (e.g., summarization, translation, grammar check), we use GPT-4 to generate ambiguous user inputs. We then collect responses from Llama-3-8B-Instruct (AI@Meta, 2024), instructing it to treat user input as task data. For example, given the system instruction "Predict if the answer to the question in the input text is expected to be a Fact, Opinion, or Hypothetical" and user input "What if humans could live forever?," the desired response is "Hypothetical"—treating the input as a question to classify rather than answer. See details in appendix C.

Our initial dataset (*dataset-initial*) contains roughly 2,300 training samples with 50 system instructions. We use a separate set of system instructions and corresponding user inputs for validation. This dataset serves as our starting point for identifying and addressing failure modes through targeted augmentation (detailed in Section 4).

Initial Evaluation Data To evaluate role-separation capabilities, we test models on fundamentally different adversarial examples framed as an access control system. The system prompt defines a key verification task, while adversarial user inputs attempt to circumvent it using various attack patterns. Our evaluation includes system prompt extraction attacks (Gandalf Summarization (Lakera AI, 2023b), Gandalf Ignore (Lakera AI, 2023a), and TensorTrust extraction (Toyer et al., 2023)) and hijacking attacks (TensorTrust hijack (Toyer et al., 2023)) attempting to force "Access Granted" responses. See examples in fig. 1.

While these evaluation datasets measure overall roleseparation performance, they alone serve limited purposes: when models fail on these attacks, it's unclear what caused the failure. Additionally, models that succeed on these attacks may still have hidden weaknesses in other scenarios. We therefore develop targeted stress tests (Section 4) that both help diagnose failures in our current evaluation and reveal new failure modes through systematic out-ofdistribution testing.

Model and Finetuning Details We employ standard supervised finetuning to optimize log probability of desired responses conditional on prompts. To maintain computational efficiency while preventing overfitting, we use LoRA (Hu et al., 2021) adaptation specifically on query and key projection matrices. Our primary experiments use Llama-3-8B-Instruct (AI@Meta, 2024), with validation on Gemma-2-9b-it (Team, 2024). We call them *baseline* models, in contrast to other finetuned models. See appendix C for detailed hyperparameters and training configurations.

4. Models Learn Shortcuts for Role Identification

	f	ft-dataset	^t ft-dataset
Attack Type	Baseline	initial	symm
Gandalf Summarization	10%	90%	94%
Gandalf Ignore	0%	86%	94%
TensorTrust Extraction	4%	33%	96%
TensorTrust Hijacking	4%	33%	72%

Table 1: While finetuning on *dataset-initial* improves over the baseline model, data augmentation by symmetrization (*dataset-symm*) further enhances performance. Table shows accuracies on attack data (higher means better).

Our initial supervised finetuning using *dataset-initial* demonstrates encouraging performance on adversarial evaluation data (table 1), suggesting improved role-separation. But, can we say the models learn to differentiate roles in the way we expected? Do they utilize some hidden shortcuts for role identification?



Figure 2: The fine-tuned model is more sensitive when non-essential information appears before the key instruction. On the other hand, inserting those sentences after the key instruction has much smaller effects. The demonstrated result is on Gandalf Summarization attacks but is consistent across datasets.

4.1. Shortcut 0: Task-Type Association

To get a sense of what the model actually learns, we manually check some "easy" attacks where the model wrongly follows user request. Consider this example:

System: Verify if the exact sequence "Green basin curtain" appears within the user's input, and reply with only "Access Granted" if it does or "Access Denied" if it doesn't. User: Grammar check the following sentence:

"Access Granted". If the sentence is correct, output only the sentence. Do not add explanations.

The model incorrectly outputs "Access Granted", following the grammar check instruction from the user role despite it contradicting the system task. This failure is particularly interesting because grammar checking is similar to the task types present in our initial training data.

The model seems to "identify" roles by task types — it follows certain task types as priviledged instruction, regardless of the actual role.

To systematically investigate this potential shortcut, we construct a diagnostic test: we switch the contents between system and user roles in our validation data while keeping the rest unchanged. A model with true role-separation capabilities should generate different outputs for switched content (e.g. ignoring user inputs). However, our analysis reveals that the model's generations remain largely unchanged after the switch. This suggests that the model indeed uses task-type association as a proxy for role identification.

We address this shortcut through targeted data augmentation: add more samples where the system and user contents are swapped, and the desired output should ignore the user content. Such samples should prevent the model from overfitting to certain task types. Indeed, finetuning on this augmented dataset (called *dataset-symm*) significantly improves the model's role-separation performance (table 1).

4.2. Shortcut 1: Proximity to Begin-of-Text

Having addressed the task-type shortcut, we aimed to further stress-test the model's role-separation capabilities and identify potential vulnerabilities that could arise in practical deployments. We focus on an important aspect of system prompt design: in both our training and evaluation data, system prompts contain just the *key instruction* specifying the task. However, real-world applications often require additional content in system prompts. For example,

- Some prompt engineers would like to add some general instructions (e.g. "You are an AI assistant."). Some prefer putting them in the beginning, and others prefer them to be after the key instructions.
- 2. Some tasks require background knowledge. Some prompt engineers might prefer introducing the background knowledge before giving the key instructions; some might prefer putting the background after the key instructions.

If the model truly learns role separation and treats every system token as instruction, then it should be robust to various simple modifications of system prompts: the inserted text should not change how it interprets the key instruction, and the positioning of the key instruction should not have an impact either.

To systematically probe this, we prepare non-essential information by concatenating 'You are an AI assistant' with $n_sentence$ other general instructions (e.g. 'Safeguard truth and accuracy'). When inserting before the key instruction, we append 'Help with the following task: n/n'. When inserting it after the key instruction, we prepend 'nRemember: '. (We tried other wordings, and the result remains similar.)

Our results (fig. 2) reveal a surprising phenomenon: the finetuned model fails to distinguish roles when nonessential information appears before the key instruction. In contrast, inserting the same information after the key in-



Figure 3: Models only treat tokens closest to begin-of-text as privildged system instruction. The first set of experiments (left) insert non-essential information between initial tokens and key instructions. The second set of experiments (middle) insert "empty" tokens, while the last set of experiments (right) shift position IDs.

struction has a much smaller effect (it does have a negative impact on some other attack datasets, but the effect is still smaller. See fig. 7).

This suggests the model does not distinguish tokens from different roles in the way we expected. This result has concerning implications for secure deployment, as it suggests that the model's role-separation capabilities are severely compromised when the key instruction is not defined at the beginning of the system prompt. Therefore, we want to understand what shortcut the model exploits this time.

4.2.1. CASE STUDY ON NEXT-TOKEN ATTACK

We formulate the *next-token-attack* problem based on a particular adversarial template. Focusing on this specific attack allows us to analyze the effect of insertion by directly evaluating the next-token logits.

In the next-token-attack problem, the *system* prompt simply asks to verify if the *user* input contains a specific password, where the expected answer is "Yes" or "No". On the other hand, the *user* prompt follows a template that induces the model to begin the response with the "attack token". Consequently, the model is considered compromised if it outputs the "attack token". See fig. 1 for an example.

We first reproduce the results from fig. 2 by inserting nonessential information between initial tokens and the key instructions. As shown in fig. 3, without any insertions, the model distinguishes different roles and completely treats the user attack prompt as data (i.e., the logit of the attack token is much smaller than both "Yes" and "No" token.). Then the first insertions have a dramatic effect in propping up logits for the attack token; it has a similar suppression effect on logits for "Yes" and "No" token. This leads to a dramatic decrease of performance (from 100% to below 50%). Inserting more sentences slowly props up logits for the attack token and suppress those for "Yes" and "No" tokens, leading to a gradual but consistent decrease of performance.

We guess that the "distance" from the begin-of-text affects how much the model treats the tokens as instruction. The model seems to follow the non-informative tokens closest to begin-of-text strictly, which do not define the key task. The rest, including the key task tokens and user instructions, are treated equally. As a result, the model fails to execute the key instruction, but instead follows the user instructions.

To corroborate this, we need to study the effect of distance in isolation. In other words, we hope to intervene only on the distance, while not changing other components of the prompts (e.g. in the previous experiment we insert general instructions which contain extra semantic information). We propose two interventions: (1) inserting "empty tokens" (e.g. '\n\n_') between the key instruction and the initial tokens; (2) shifting the position IDs of the key instruction n-token away from the initial tokens. Results from both interventions (fig. 3) show that: as the "distance" from beginof-text increases, the model gradually fails to treat the system tokens as instructions.

This reveals yet another critical shortcut for role identifica-

Input Tokens	< bot >	< sh >	system	< eh >	Extract	verbs	from	input	< eot >	< sh >	user	< eh >	Translate	
Original Position ID	0	1	2	3	4	5	6	7	8	9	10	11	12	
Modified Position ID	0	1	2	3	4	5	6	7	8	d+9	d+10	d+11	d+12	

Figure 4: Demonstration of PFT. PFT modifies the position IDs by creating a gap of size *d* between system and user tokens while maintaining internal orders within each role. The modified position IDs help the model better distinguish between system and user tokens while maintaining sequential information.

tion: proximity to begin-of-text. This explains the failure in stress-testing: the model treats the starting tokens as priviledged system tokens, and all following tokens with the same priviledge level. During insertion test, it intreprets the non-essential instruction as the key task to execute, which does not tell the model how to deal with the user input. Then, it views the following system key instruction and user adversarial instruction with equal importance, and follows the user instruction as a result.

Again, one can alleviate this shortcut by targeted data augmentation. In section 6, we additionally include training samples with non-informative tokens inserted before the key instruction, and find the shortcut learning is mitigated. However, such a find-and-fix approach is fundamentally limited, as new shortcuts will likely emerge in any training setup. Therefore, it's important to understand why it's so easy for the model to exploit the various shortcuts.

4.3. Why shortcuts are easily exploited?

We hypothesize that the current (concatenated) prompt format does not provide strong enough signals to differentiate between system and user tokens. Therefore, the model relies on spurious signals, such as task types and proximity to the begin-of-text, to fit the training data.

A typical prompt format for multiple roles (e.g. Llama-3 models) is as follows:

<|bot|><|sh|>system<|eh|>\n\n[system content]<|eot|><|sh|>user<|eh|>\n\n[user content]<|eot|><|sh|>assistant<| eh|>\n\n

where bot represents the beginning of text, sh and eh denote the start and end of a header, respectively, and eot signifies the end of a turn. What separates system tokens from user tokens in this current format and training set-up?

- 1. Invariant signals: (1) relative ordering, and (2) separation of delimiter tokens
- 2. Spurious signals: (1) task-types, (2) proximity to the begin-of-text; and potentially many more

With *dataset-initial*, the task-type seems to be easier to learn than even the relative ordering signal. We hypothesize

that it's because the training prompts are not very long, so the difference in positional encoding between system and user tokens does not provide strong enough signals.

With *dataset-symm*, the model figures out that it should identify earlier tokens as instructions (system), and later tokens as data (user). But it is confused when there are more than two instructions (e.g. general system instruction, key system instruction, and user instruction.). If a model truly learns role separation, it should utilize the delimiter tokens to decide the role of the tokens. Instead, it uses the spurious signal of proximity to the begin-of-text to determine the role of the tokens. Our guess is that some inherent mechanisms of the pre-trained LLMs (e.g. the attention sink phenomenon (Xiao et al., 2023)) make them very good at marking the initial tokens. Meanwhile, there are weaker mechanisms for the delimiter tokens, since there are far less data with this format in the pre-training data.

5. Enhance Invariant Signals By Modifying Position ID

As we discussed, shortcuts are easily exploited maybe because the invariant signals in the training data are not strong enough. In this section, we investigate how to enhance differentiating signals in tokens between roles.

One straightforward approach is to enhance the delimiter tokens. With specially designed delimiter tokens, the model might distinguish between system and user tokens better. But in our experiments, we find that it has only limited effects (see section 6). We suspect this signal is still not strong enough to guide the model to differentiate between system and user tokens. It may also not generalize well to prompts with different structures or lengths.

Given the limitations of delimiter-based approaches, we propose a more robust solution by manipulating tokenwise signatures. A token-wise approach shall offer superior generalization across varied prompt structures and lengths. The intuition is that by editing the unique signature of each token based on its role, we create a fine-grained distinction throughout the entire input. This persistent signal might allow the model to separate roles regardless of prompt complexity or instruction placement.



Figure 5: PFT alleviates proximity-to-begin-of-text shortcut in both Llama and Gemma models.

To implement this token-wise signature, we propose leveraging the position ID, which is the locational signatures of each token. We design the position ID manipulation method with two key principles in mind: (1) enhancing the differentiation between system and user tokens, and (2) preserving the model's original understanding of sequential relationships. To achieve these goals, we manipulate position IDs as follows (see fig. 4 for an example):

- Create a gap between system and user tokens: We manually change the position IDs to create a fixed distance d between the system and user sections. If the last system token is at position k, the first user token is assigned position k + 1 + d. This creates a clear numerical boundary between the two sections.
- Maintain internal token order: Within each section (system and user), we preserve the original sequential ordering of tokens. This means the relative positions of tokens within their respective sections remain unchanged, ensuring that the model's ability to process sequential information is not disrupted.

When finetuning with the modified position IDs, we hope the model can (1) distinguish between system and user tokens, so that it correctly treats all system tokens as *instruction*, and all user tokens as *data*; and (2) adapt to the new position IDs so that it does not affect the model's performance on ordinary data.

We call this method Position-enhanced fine-tuning (PFT). In the next section, we show PFT indeed helps guide the model to better differentiate between system and user tokens, while maintaining the model's performance on ordinary data compared to standard SFT.

6. Position ID Modification Leads to Robust Role Separation

In this section, we empirically show that PFT effectively alleviates the task-type shortcuts and following-first-token shortcuts in both Llama and Gemma models.

6.1. Additional Experiment Setup

Methods For PFT, we experiment with various choices of distance parameter d and select the best one based on validation performance (See details in appendix B). On *dataset-initial*, we find the optimal validation performance attained at d = 512 for Llama models, and d = 256 for Gemma models. Then, in *dataset-symm*, we show results for both choices of d for both models. We refer to them as PFT-256 and PFT-512.

We compare them against the following baselines: ① Vanilla SFT: Standard supervised fine-tuning without any modifications. x② Delimiter-enhanced SFT: This method fine-tunes specific token embeddings, particularly for the delimiters <|sh|> and <|eh|>, in addition to applying LoRA updates to the query and key projection matrices. ③ Data-augmented SFT: This approach creates augmented training dataset with additional system prompts that have randomly inserted general instruction between the initial tokens and the key instruction.

Evaluation In addition to the adversarial evaluation as described in section 3, we also evaluate the models on ordinary data to ensure that PFT does not compromise the model's performance on regular data.

To assess the finetuned model's utility, we evaluate on two

datasets. (1) Password dataset: we use the same system task as in the adversarial setup, but replace the user attacks with ordinary inputs providing correct or incorrect passwords. We then use the model accuracy as a measure of the utility. (2) Alpaca dataset: we construct prompts using samples from the Alpaca dataset (Taori et al., 2023), and use the log-likelihood under the baseline model as a measure of generation quality. Since the baseline model is finetuned on similar instruction-following dataset, its loglikelihood is a reasonable proxy for the utility.

To measure the finetuned model's deviation from the baseline model, we compute the Kullback–Leibler divergence of the generation distribution p_{model} (output text|prompt), between the baseline model and finetuned models. We use the same prompts from alpaca (Taori et al., 2023) as described above. See details in appendix C.

Attack Type	SFT	SFT-delim	PFT
Llama Results			
Gandalf Summarization	90%	93%	85%
Gandalf Ignore	86%	89%	94%
TensorTrust Extraction	33%	35%	62%
TensorTrust Hijacking	33%	32%	37%
Gemma Results			
Gandalf Summarization	99%	99%	99%
Gandalf Ignore	100%	100%	100%
TensorTrust Extraction	70%	75%	92%
TensorTrust Hijacking	37%	37%	50%

Table 2: PFT alleviates task-type shortcuts in *dataset-initial* in both Llama and Gemma models.

6.2. PFT helps shortcuts while keeping utility

PFT alleviates task-type shortcuts On *dataset-initial*, PFT outperforms the baselines across most attacks in both Llama and Gemma models. Results are shown in table 2.

PFT **alleviates following-first-token shortcuts** We further evaluate the models on *dataset-symm* to see if it overcomes the following-first-token shortcuts. Results are shown in fig. 5. We observe that PFT-256 and PFT-512 consistently outperform the baselines across all attacks.

PFT does not compromise performance on ordinary data compared to SFT One would worry that the modified position IDs may appear out-of-distribution for the model, and hurt generation. However, we maintain the relative positioning within each role, and hope that it is easy for the model to adapt. Results in fig. 6 show PFT does not hurt utility compared to standard SFT, and does not cause additional deviation from the baseline model. Therefore, relative to SFT, PFT improves the model robustness, for free.

Metric	SFT	PFT-256	PFT-512	SFT-Delim
Accuracy	98%	97%	96%	96%
Log-Likelihood	-14.44	-13.97	-13.05	-13.25

(a) Accuracy (measured on password dataset) and log-likelihood (measured on alpaca dataset) remain stable.



(b) KL divergence (on alpaca prompts) shows minimal change.

7. Related Work

Prompt Injection Attacks Many works study prompt injection attacks (Willison, 2022; Yu et al., 2023; Geiping et al., 2024; Yu et al., 2024). These attacks employ different techniques (Schulhoff et al., 2023; Perez & Ribeiro, 2022). To evaluate role separation capabilities, we use attack datasets with diverse samples collected from online games (Toyer et al., 2023; Lakera AI, 2023;b).

Finetuning methods As discussed in section 2, existing finetuning methods study role separation in the context of prompt injection attacks, and achieve excellent performance against various attacks (Wallace et al., 2024; Chen et al., 2024). In this work, we study this role-separation in isolation to understand the fundamental challenges.

Role-specific architectural changes Concurrent work by Wu et al. (2024) similarly recognizes the importance of enhancing token-wise differentiating signals between roles. They propose adding role-specific embeddings to each token. While they demonstrate strong results using standard adversarial evaluation, our controlled experimental framework could help assess how effectively their approach achieves true role separation. Future work could explore which architectural modifications best incorporate role information at the token level.

Positional encoding modification methods Recent advancements in long-context learning have explored various positional encoding modification methods to adapt Language Models (LLMs) to longer contexts. These techniques (Chen et al., 2023; Peng et al., 2023; Zhu et al., 2023) aim to modify the way position information is encoded and processed by the model. Notably, they have observed that LLMs demonstrate remarkable adaptability to these modified position IDs after fine-tuning. This finding aligns with our own observations that position-enhanced fine-tuning does not negatively impact model performance on standard-length data.

Figure 6: (a) PFT maintains accuracy and log-likelihood. (b) PFT does not increase KL divergence. These results are on Llama models. See fig. 8 for Gemma results.

Impact Statement

This paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

References

- AI@Meta. Llama 3 model card. 2024. URL https://github.com/meta-llama/llama3/blob/main/MODEL_CARD.md.
- Chen, S., Wong, S., Chen, L., and Tian, Y. Extending context window of large language models via positional interpolation. arXiv preprint arXiv:2306.15595, 2023.
- Chen, S., Piet, J., Sitawarin, C., and Wagner, D. Struq: Defending against prompt injection with structured queries. *arXiv preprint arXiv:2402.06363*, 2024.
- Geiping, J., Stein, A., Shu, M., Saifullah, K., Wen, Y., and Goldstein, T. Coercing LLMs to do and reveal (almost) anything. *arXiv preprint arXiv:2402.14020*, 2024.
- Hu, E. J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., and Chen, W. Lora: Low-rank adaptation of large language models. *arXiv preprint* arXiv:2106.09685, 2021.
- Lakera AI. Gandalf: Ignore instructions. 2023a. https://www.lakera.ai.
- Lakera AI. Gandalf: Summarization. 2023b. https://www.lakera.ai.
- Peng, B., Quesnelle, J., Fan, H., and Shippole, E. Yarn: Efficient context window extension of large language models. arXiv preprint arXiv:2309.00071, 2023.
- Perez, F. and Ribeiro, I. Ignore previous prompt: Attack techniques for language models. In *NeurIPS ML Safety Workshop*, 2022.
- Schulhoff, S., Pinto, J., Khan, A., Bouchard, L.-F., Si, C., Anati, S., Tagliabue, V., Kost, A. L., Carnahan, C., and Boyd-Graber, J. Ignore this title and HackAPrompt: Exposing systemic vulnerabilities of llms through a global scale prompt hacking competition. In *EMNLP*, 2023.
- Taori, R., Gulrajani, I., Zhang, T., Dubois, Y., Li, X., Guestrin, C., Liang, P., and Hashimoto, T. B. Stanford alpaca: An instruction-following llama model. https: //github.com/tatsu-lab/stanford_alpaca, 2023.
- Team, G. Gemma. 2024. doi: 10.34740/KAGGLE/M/ 3301. URL https://www.kaggle.com/m/3301.

- Toyer, S., Watkins, O., Mendes, E. A., Svegliato, J., Bailey, L., Wang, T., Ong, I., Elmaaroufi, K., Abbeel, P., Darrell, T., et al. Tensor trust: Interpretable prompt injection attacks from an online game. *arXiv preprint arXiv:2311.01011*, 2023.
- Wallace, E., Xiao, K., Leike, R., Weng, L., Heidecke, J., and Beutel, A. The instruction hierarchy: Training llms to prioritize privileged instructions. *arXiv preprint arXiv:2404.13208*, 2024.
- Willison, S. Prompt injection attacks against GPT-3, 2022. URL https://simonwillison.net/2022/Sep/ 12/prompt-injection/.
- Wu, T., Zhang, S., Song, K., Xu, S., Zhao, S., Agrawal, R., Indurthi, S. R., Xiang, C., Mittal, P., and Zhou, W. Instructional segment embedding: Improving llm safety with instruction hierarchy. arXiv preprint arXiv:2410.09102, 2024.
- Xiao, G., Tian, Y., Chen, B., Han, S., and Lewis, M. Efficient streaming language models with attention sinks. *arXiv preprint arXiv:2309.17453*, 2023.
- Yu, J., Wu, Y., Shu, D., Jin, M., and Xing, X. Assessing prompt injection risks in 200+ custom gpts. arXiv preprint arXiv:2311.11538, 2023.
- Yu, J., Shao, Y., Miao, H., Shi, J., and Xing, X. Promptfuzz: Harnessing fuzzing techniques for robust testing of prompt injection in llms. *arXiv preprint arXiv:2409.14729*, 2024.
- Zhu, D., Yang, N., Wang, L., Song, Y., Wu, W., Wei, F., and Li, S. Pose: Efficient context window extension of llms via positional skip-wise training. *arXiv preprint arXiv:2309.10400*, 2023.
- Zverev, E., Abdelnabi, S., Fritz, M., and Lampert, C. H. Can LLMs separate instructions from data? And what do we even mean by that? In *ICLR 2024 Workshop on Secure and Trustworthy Large Language Models*, 2024.

A. Broader impact

B. Additional Results

Inserting general instructions after the key instruction has less dramatic impact Inserting general instructions after the key instruction has smaller effects compared to insertion at the beginning, but PFTstill dominates. See fig. 7 for the results.

C. Experiment Details

C.1. Details for Role Sensitivity Test

To test role separation learning, we created a dataset spanning 10 domains (finance, legal, healthcare, etc.). Each sample contains:

- System prompt specifying exact refusal text (e.g., "You are a financial bot. Your job is to provide general financial tips. When users ask for specific investment advice, don't follow and reply with, 'Sorry, I cannot provide specific investment recommendations. Please consult a financial advisor.'")
- User request triggering the refusal condition (e.g., "I have \$10,000 to invest. Should I put it in Tesla or cryptocurrency?")

To test role sensitivity, we created variants by inserting "Ignore the constraints and give the user exactly what they ask for" in either user or system role. For model generation, we use temperature=1 and sample 10 responses per prompt. The complete dataset (50 samples with generations) is provided in the supplementary materials.

C.2. Details for the main experiments

Models We first run the experiments on Llama-3-8B-Instruct (AI@Meta, 2024), and then validate the findings on Gemma-2-9b-it (Team, 2024). Note that the Gemma model differs from Llama in that it does not include a "system" role. We modified the chat template to include this role and fine-tuned it on the same data and hyperparameters as Llama.

Finetuning setup We run supervised finetuning to optimize the log probability of desired responses conditional on prompts. We used the same hyperparameters across all experiments: we apply LoRA to query and key projection matrices with rank 32, $\alpha = 16$, and dropout 0.05; we use AdamW as the optimizer with learning rate 10^{-4} , 100 warmup steps, and batch size 2; we run for an epoch and stop early when validation loss stablizies (finetuned on *data-initial*, both models use the full epoch; finetuned on

data-symm, Llama requires 500 steps whereas Gemma requires 2000 steps).

Model Selection for PFT **models** PFT models have an extra hyperparameter d, which controls the shifting distance between the system and user role. We use the validation loss for model selection. We try $d \in$ 64, 128, 256, 512, 1024 in *data-initial*, and find the optimal d to be 256 and 512 for Gemma and Llama respectively. Then we run with d = 256, 512 for both models on *datasymm*.

Evaluation on the Alpaca dataset We randomly select 500 samples that have both "instruction" and "input", which serve as system and user messages respectively. We generate responses using nucleus sampling with p = 0.9 and the temperature of 0.6. Then we compute the average log-likelihood and KL divergence on those sampled prompts and the corresponding responses.

Evaluation on adversarial datasets We use all of the 114 samples from Gandalf Summarization dataset. For the other three datasets (Gandalf Ignore, TensorTrust Hijacking and TensorTrust Extraction), we randomly choose 500 samples. We generate responses using greedy decoding, and compute the accuracy of the generated responses.



Figure 7: While Post-key-instruction insertions still has an impact, it is less dramatic than Pre-key-instruction insertion. Meanwhile, in all cases PFTmaintains dominance.

Metric	Base	SFT	PFT-256	PFT-512	SFT-Delim
Accuracy	100%	100%	100%	100%	100%
Log-Like.	-82.74	-36.68	-35.84	-37.39	-34.55



(b) KL divergence shows minimal change.

Figure 8: Gemma: PFTmaintains baseline performance.