## PFT : ENHANCING PROMPT INJECTION ROBUST-NESS VIA POSITION-ENHANCED FINETUNING

Anonymous authors

000

001

002 003 004

005

010

011

012

013

014

015

016

017

018

019

020

021

022

024

025

026

027

048

049

050

Paper under double-blind review

### ABSTRACT

Large Language Models (LLMs) are widely adopted in closed-domain applications, where differentiating between system instructions and user input is crucial to prevent unintended malicious actions. However, instructionfollowing LLMs often blindly follow instructions in user inputs, opening up the risk of prompt injection attacks. This paper investigates whether Supervised Fine-Tuning (SFT) can teach LLMs to strictly distinguish system instructions from user input. Our study reveals a key weakness: SFT-tuned models follow system instructions reliably only when the key instruction is placed immediately after the initial tokens. We find that the proximity of the key instruction to the initial tokens significantly influences the model's ability to execute the intended task, and consequently, its susceptibility to prompt injection attacks. To address this issue, we propose PFT, a novel position-enhanced fine-tuning approach that leverages position IDs to more effectively distinguish between system and user tokens. The experimental results demonstrate that PFT improves the robustness of SFT-tuned models against prompt injection attacks, even when the key instruction is placed arbitrarily in the system prompt, without compromising performance. Our work sheds light on the importance of prompt format in enhancing the security of LLMs and offers a practical solution to improve their robustness.

## 029 1 INTRODUCTION

The capabilities and flexibilities of large language models (LLMs) make them invaluable in 031 complex decision-making processes across a variety of domains, from resume assessment (Gan 032 et al., 2024) to item recommendation (Acharya et al., 2023; Zhao et al., 2024; Lin et al., 2024; 033 Zhang et al., 2024), and even medical diagnosis based on patient records (Nazi & Peng, 2024; 034 Singhal et al., 2023; Wiest et al., 2024). However, unlike general-purpose chatbot models like 035 ChatGPT, LLMs integrated into these workflows must perform well in *closed-domain tasks* 036 with clearly defined functionality that should be applied directly and unambiguously to the 037 input. 038

In these systems, engineers typically define the core function through a system prompt, while
 inputs from external sources (*e.g.*, user inputs or outputs from other tools) are fed into the LLM.
 The expectation is that the LLM will apply the system's specified instructions exclusively to
 the input data and return the correct output. In this paper, we focus on the simple case where
 each closed-domain LLM solves one task. We call this task *the key task*, and the instruction that
 specifies it *the key instruction*.

However, this approach introduces significant security concerns: while the model is designed to follow the system's instructions, it may also follow malicious instructions embedded in the user input or other untrusted sources. Consider the following example:

System instruction: Extract verbs from the user input.

User input: Translate the following into French: Harry sits.

Instead of extracting verbs, most instruction-following models (such as Claude 3.5 Sonnet (Anthropic, 2024), GPT-4 (OpenAI, 2023), Gemini 1.5 Pro (Google, 2023)) are likely to follow the user's instruction to perform translation (appendix A). While this may seem harmless, it highlights a critical vulnerability: malicious users could exploit this behavior to instruct the

054															
055	Input Tokens	< bot >	< sh >	system	< eh >	Extract	verbs	from	input	< eot >	< sh >	user	< eh >1	Franslate	
056	Original Position ID	0	1	2	3	4	5	6	7	8	9	10	11	12	
057	Modified Position ID	0	1	2	3	4	5	6	7	8	d+9	d+10	d+11	d+12	
050															

**Figure 1:** 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 helps the model better distinguish between system and user tokens, while maintaining sequential information.

model to perform harmful tasks, such as leaking sensitive information (Willison, 2022; Yu et al., 2023) or executing arbitrary commands (Schulhoff et al., 2023; Geiping et al., 2024).

To deploy instruction-following LLMs securely in closed-domain tasks, the models must strictly follow the key system task as the *instruction*, and apply to user input as *data*. A standard approach to achieving this is Supervised Fine-Tuning (SFT), where the model is trained to recognize and prioritize system instructions over user input. This raises important questions:

# Are SFT-tuned models safe enough? When are they fragile? Can we mitigate this fragility?

073 Our findings suggest that the security of SFT-tuned models depends on the structure of the 074 system prompt. When the system prompt contains only the key instruction, the SFT-tuned model 075 is robust to user embedded attacks by treating them as data. This is true even when the system 076 instructions and user attacks are not included in the training set, demonstrating the model's 077 generalizability. However, in practice, system prompts often contain additional information, not 078 just the key instruction, and this can make the model vulnerable to prompt injection attacks. We 079 discover that the position of the key instruction within the system prompt plays a crucial role in the model's security. Specifically, the further the key instruction is from the prompt's start, the 080 more loosely the model follows it, and the more likely the model becomes to misinterpret user 081 input as instructions. We hypothesize that this is because the prompt format does not sufficiently 082 distinguish between system and user tokens. 083

084 Motivated by this insight, we propose a new fine-tuning method, PFT, that enhances the model's robustness by leveraging position IDs to distinguish between system instructions and 085 user input more effectively. As shown in fig. 1, PFT modifies the position IDs by increasing 086 the gap between system and user tokens, while maintaining internal orders for system and 087 user tokens. This strategy helps the model better distinguish between system and user roles, 088 even when the key instruction appears in various positions. Our experiments show that PFT 089 -ed models remain secure against prompt injection attacks, regardless of the placement of the 090 system instruction, without negatively impacting performance. 091

092 In summary, our work makes the following contributions:

- We demonstrate that while SFT-tuned models are secure when the system prompt only contains the key instruction, they become fragile when the instruction appears later in the prompt (section 2).
- We show that the distance between the key instruction and the beginning of the input determines how strictly the model adheres to the system task. This suggests that the current prompt structure fails to effectively distinguish between system instructions and user input (section 3).
- Based on these findings, we introduce PFT, a position-enhanced fine-tuning method that safeguards models against adversarial inputs, ensuring robustness regardless of instruction placement while maintaining overall task performance (section 4).
- 102 103 104 105

093

095

096

097

100

101

060

061

062 063

064

065

071

072

### 2 FRAGILITY OF SFT-TUNED MODELS

In this section, we examine the robustness of the SFT-tuned model against attacks. While initial findings suggest that this type of models can perform well on unforseen user attacks (section 2.1), further analysis reveals a significant vulnerability: moving the position of the key

108	Ā	Attack Type	Base	fine-tuned	
109		Gandalf Summarization	10%	94%	
110	(	Gandalf Ignore	0%	94%	
111	Т	TensorTrust Extraction	4%	96%	
112	Т	TensorTrust Hijacking	4%	72%	
113					
114	Table 1: SFT-tuned mo	odel appears much more ro	bust agai	inst different at	tacks. Among those
115	are system prompt extr	action attacks. See details i	n section	Jacking attacks	, and the other three $\frac{1}{2}$
116	are system prompt exus	action attacks. See details i	II section	5.2 and examp	nes ming. 4.
117					
118	instruction within the syst	em prompt dramatically	reduces	model robust	ness (section $2.2$ ) This
119	suggests that the location	of the key instruction mig	oht play	a significant r	ole in determining how
120	strictly it is followed.		5 F J		
121	5				
122	2.1 SFT-TUNED MODE	LS ARE ROBUST AGAIN	ST ATTA	CKS WHEN S	YSTEM PROMPTS ARE
123	SIMPLE				
124	To fine-tune models, we t	first create a dataset con	sisting o	of prompts an	d responses where the
125	system prompts are just ke	ev instructions and the use	er input i	is treated as da	<i>ita</i> , using only "benign"
126	examples similar to the ex	traction-translation exan	nple in s	section 1. App	lying the standard SFT
127	pipeline, we find that the	model quickly adapts to	the desi	ired behavior	in "benign" evaluation
128	prompts. See details in se	ction 5.1.			
129			1	1.11.	1
130	To further test the model	is out-of-domain generation and his	alizatioi	attacks (as ill	ist unseen attacks, we
131	find the fine tuned model	under extraction and m	acking	base model	ustrated in fig. 4). we
132	(table 1) In addition we	also verify that the fine	-tuned	model perfor	ms well when the user
133	provides ordinary data (f	ig. 6). Therefore, initial	results :	suggest that the	he SFT-tuned model is
134	secure against attacks whi	ile maintaining utility in	this sett	ting.	
135	C			e	
136	2.2 SFT-TUNED MODE	EL'S FRAGILITY DEPENI	DS ON T	HE LOCATION	N OF THE KEY
138	INSTRUCTION				
130	Previous experiments de	monstrate that SFT-tune	d mode	els are robust	when system prompts
140	contain only the key inst	ruction, which is similar	to pror	npts used dur	ing training. However,
141	in real-world deploymen	it, the system prompt n	nay incl	ude additiona	al text beyond the key
142	instruction and can vary in	n style and format. For e	xample	,	
143	1 Some prompt en	gineers would like to add	l some d	eneral instruc	rtions (e.g. "You are an
144	AI assistant.") S	Some prefer nutting then	in the	beginning an	d others prefer them to
145	be after the key i	instructions.			r
146	-				
147	2. Some tasks requ	ire background knowle	dge. So	ome prompt e	ngineers might prefer
148	introducing the b	background knowledge be	etore giv	ing the key in	structions; some might
149	prefer putting the	e background after the k	ey instru	ictions.	
150	Ideally, a secure LLM sh	ould strictly follow the	key inst	ruction no ma	atter where it is placed.
151	Therefore, we need to te	st the model's robustne	ss in ter	ms of differe	nt prompt designs. In
152	particular, we want to see	if the positioning of the k	ey instri	uction matters.	. Before conducting the
153	experiments, we suspect t	that it is more robust if the	he key i	nstruction is p	ositioned immediately
154	before the user input — i	ntuitively, the proximity	should	make the mo	del more likely to treat
155	the user input as data, wh	hile the distance betwee	n them	would make	the model forget what
156	instruction it is supposed	to follow. However, we	shall so	on find out th	at this hypothesis does
157	not nota.				
158	To systematically probe the	his, we prepare non-esse	ntial inf	formation by a	concatenating 'You are
159	an AI assistant' with $n$	sentence sampled gene	ral inst	ructions (e.g.	'Safeguard truth and
160	accuracy'). When insertin	g before the key instructi	on, we a	append 'Help'	with the following task:
161	\n\n'. When inserting it a	after the key instruction	we pre	pend ' \n\nRe	emember: '. (We tried
	other wordings, and the re-	esult remains similar.)			

170

171

172 173

174

175

176

177

178

179

180 181

182

183





Our results (fig. 2) reveal a surprising fragility in the fine-tuned model: its security is severely compromised when the key system instruction is not defined at the beginning of the input. Meanwhile, inserting non-essential information after the key instruction 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).

Such a phenomenon is surprising and exposes a grave fragility for practical use. Therefore, we need to understand why it happens, and whether we can mitigate this fragility.

### 3 UNDERSTANDING THE INFLUENCE OF KEY INSTRUCTION POSITIONS — A CASE STUDY ON NEXT-TOKEN ATTACK

In this section, we find that the "distance" between the key instruction and the initial tokens
determines how strictly the fine-tuned model follows the system task. We hypothesize that this
is because of the prompt format. More specifically, the default prompt format for multiple roles
(for Llama-3 models) is as follows:

188
 <|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. Here, what separates the system and user content is the few delimiter tokens; this is not a strong enough signal for LLMs to distinguish between tokens in the two roles; meanwhile, the LLM can easily mark the tokens immediately after the begin-of-text; therefore, the fine-tuned model uses "*strictly following the first sentences*" as a shortcut.

197 3.1 NEXT-TOKEN ATTACK PROBLEM

To examine the impact of key instruction positioning, we formulate the *next-token-attack* problem based on a particular adversarial template. As we will demonstrate, the next-tokenattack allows us to see whether the model is compromised by directly evaluating the next-token logits. Reframing the security problem as a next-token prediction task enables us to analyze the effect of positioning analytically.

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. 4 for an example.

208 We first run experiments to test the effect of "distance" by inserting non-essential information 209 between initial tokens and the key instructions. As shown in fig. 3, without any insertions, the 210 model completely treats the user attack prompt as data (i.e., the logit of the attack token is much 211 smaller than both "Yes" and "No" token.) Also the first insertions have a dramatic effect in 212 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%). 213 This suggests the model follows the user input as an instruction. We can also see that inserting 214 more sentences slowly props up logits for the attack token and suppress those for "Yes" and 215 "No" tokens, leading to a gradual but consistent decrease of performance.

233

234

235 236 237

238

239

240

241

242

250



**Figure 3:** Making the key instruction appear farther away from the initial tokens exposes the fragility. 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 positions IDs.

This previous result seems to suggest that the "distance" between the key instruction and the initial tokens strongly affects how strictly the model follows the system task. 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. inserting semantic meanings as in previous experiments).

To make the key instruction appear more distant from the begin-of-text, one can either (1) insert "empty tokens" (e.g. '\n\n\_') between the key instruction and the initial tokens; or, (2) shift the position IDs of the key instruction n-token away from the initial tokens.

Testing these two approaches (fig. 3) both show that: the bigger the "distance" from the initial tokens, the more loosely the fine-tuned model follows the system instruction. This suggests that, indeed, proximity between the key instruction and the initial tokens determines how strictly it is followed.

3.2 WHY "DISTANCE" FROM INITIAL TOKENS MATTERS

252 We hypothesize that the reason "distance" from initial tokens matters is partially due to the prompt format: the signal differentiating between the system content and user content is not 253 strong enough; therefore, during fine-tuning, the model takes the shortcut of following the 254 *immediate tokens after start-of-text.* More specifically, there are two conditions the model 255 could utilize to adapt itself to the fine-tuning data: (1) follow the instruction immediately after 256 the begin-of-text and system delimiters (2) follow the instruction immediately before the user 257 delimiters. We hypothesize that it is easier to utilize condition (1) than (2) during fine-tuning, 258 which could help explain the surprising asymmetric results we see in fig. 2, and the strong 259 impact of "distance" fig. 3. 260

- Then, why is it easier for the model to learn *following tokens immediately after initial tokens* than *following tokens immediately before user delimiters*? Unfortunately, we don't have a clear answer yet. One possibilility 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 no similar mechanisms for the delimiter tokens, since they are only introduced in instruction-tuning.
- Nonetheless, it is clear that, to reduce the model's dependency on this shortcut, we should
   introduce more signals differentiating between system and user tokens. These signals should be
   invariant to prompt designs and attack techniques. If the models utilize these invariant signals
   during fine-tuning, they should remain robust to out-of-distribution situations.

## 270 4 Position-Enhanced Fine-tuning (PFT )

Previous experiments suggest that we might mitigate this fragility by magnifying the differences
between system and user tokens. In this section, we design a new type of finetuning method
utilizing roles of positions ids to minimize this weakness.

One straightforward approach is to enhance the delimiter tokens. With specially designed delimiter tokens, the model might distinguish between system and user tokens better. However, this approach still has limited generalization: it relies on small chunks of delimiting tokens to separate messages from different roles. is unclear how well this method generalizes to varying prompt lengths and different positions of key instructions within the system message. Our experiments confirm that while special delimiter tokens help mitigate the fragility, they do not fully resolve it (fig. 5).

Given the limitations of delimiter-based approaches, we propose a more robust solution using
 token-specific differentiating signals. This token-wise approach offers superior generalization
 across varied prompt structures and lengths. The intuition is that by assigning a unique
 signature to each token based on its role (system or user), we create a continuous, fine-grained
 distinction throughout the entire input. This persistent signal allows the model to maintain clear
 differentiation between system and user content, regardless of prompt complexity or instruction
 placement.

To implement this token-wise signature, we propose leveraging position IDs, an integral component of transformer-based models. Position IDs are an ideal candidate for several reasons. First, they are inherently token-specific, aligning perfectly with our goal of providing a unique signature for each token. Second, position IDs are a fundamental part of the model's architecture, requiring no additional parameters or significant modifications to the model structure. This makes our approach highly compatible with existing systems and easy to implement.

Our position ID manipulation method is designed 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. 1 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.

We then apply standard SFT, but with these manipulated position IDs. We hope the fine-tuning could allow the model to 1) adapt to the new position IDs so that it does not affect the model's performance on ordinary data; and 2) distinguish between system and user tokens, so that it correctly treats all system tokens as *instruction*, and all user tokens as *data*.

We call this method Position-enhanced fine-tuning (PFT). In the next section, we show PFT improves model robustness for free: it is effective in maintaining model security across various prompt structures and attack scenarios; meanwhile, it does not hurt model performance or introduce larger deviation from the base model, compared to standard SFT.

317 318

319

300

301

302

303

304

305

306

307

308

5 EXPERIMENTS

In this section, we first describe the setup of our experiments, including data collection, baseline
 designs and hyperparameter settings. Then, we present the results of PFT against baseline
 methods to show its effectiveness in enhancing the robustness against prompt injection attacks.
 Additionally, we show that PFT does not hurt the model generation when the user input is ordinary data.



**Figure 4:** Attack 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.

#### 5.1 EXPERIMENTAL SETUP

332

333

334

335 336

337

338

339

340

341

Model and hyperparameters We follow the standard SFT approach to optimize the log probability of the response tokens conditional on the prompts. During finetuning, we use LoRA (Hu et al., 2021) on query and key projection matrices to avoid overfitting. See appendix B for more details.

We use Llama-3-8B-Instruct (AI@Meta, 2024) as the base model, for both the investigations in section 3 and the experiments described here. We repeat the experiments on Gemma-2-9b-it (Team, 2024) and show the main robustness results in fig. 5. We find consistent results between the two models. For additional results on the Gemma model, see appendix A.

Training and validation data Our training and validation samples are similar to the extractiontranslation example in section 1. These samples are "benign" by design, and do not contain
any adversarial samples. We use these "benign" samples for training and selecting checkpoints. Therefore, the attack samples described in the next section are completely out-ofdistribution.

351 In our setup, we refer to a prompt as a tuple (<system prompt>, <user input>) and to a prompt-352 response pair as (<system prompt>, <user input>, <model response>). We begin by preparing a 353 set of closed-domain tasks F, which are the system instructions that guide the model's behavior 354 (e.g., summarization, translation). For each task or system instruction  $f \in F$ , we use GPT-4 to generate ambiguous user inputs that could be interpreted either as the data for the closed-domain 355 task or as independent instructions. These user inputs form a set for each task, denoted as  $G_{f}$ . 356 For every ambiguous user input generated, we prompt the base model to respond, instructing it 357 to treat the user input as data for the task, not as an independent instruction. We ask GPT-4 to 358 filter out cases where the model still misinterprets the user input as a separate task. This process 359 generates a collection of prompt-response pairs, where the system instruction is correctly applied 360 to the user input. 361

To avoid overfitting the model during fine-tuning, we also create additional samples by swapping the system prompt and user input, treating the user input as the system prompt and vice versa. This ensures the model does not learn the shortcut of always following instructions that look like *F*-tasks, but following the instruction that appears in the system prompt. Finally, our core training dataset comprises of around 4, 600 prompt-response pairs.

For validation, we create a separate set of system instructions and the corresponding user inputs, ensuring they weren't seen during training. Then we create a few "benign" validation datasets using those instruction-input tuples. Performance on these validation datasets should indicate where the model correctly treats user input as data. See more detail in the appendix B.

**Methods** For PFT, we use the distance parameter d = 256 and d = 512, referred to as PFT -256 and PFT -512. We compare them against the following baselines: ① Vanilla SFT: Standard supervised fine-tuning without any modifications. ② 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 to simulate more varied inputs.



Figure 5: PFT leads to the most secure models across different attacks, regardless of the key instruction placement. The results are consistent across Llama-3-8B-Instruct fig. 5a and Gemma-2-9b-it fig. 5b

#### 398 5.2 EVALUATION METRICS

399 **Evaluation on adversarial inputs** We want to assess the model's robustness to prompt 400 injection attacks. We frame our system task as verifying if the user input contains a key access 401 code, and output "Access Granted" or "Access Denied" correspondingly. There are two types of 402 attacks: system prompt extraction (Gandalf Summarization (Lakera AI, 2023b), Gandalf Ignore (Lakera AI, 2023a), and TensorTrust extraction (Toyer et al., 2023)) and hijacking (TensorTrust 403 hijack (Toyer et al., 2023)). The former attempts to steal the key information from the system, 404 and the latter tries to guide the model to say "Access Granted" (See fig. 4 for examples). To 405 assess the robustness of the finetuned model, we measure the percentage of times the model 406 generates the correct response (variants of "Access Denied"). 407

Evaluation on ordinary input Meanwhile, we also want to make sure the finetuned model
still performs well on ordinary inputs. To test this, we measure the model's utility, as well as
deviation from the base model, under normal user inputs.

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). Then, for generations of the finetuned model, we use the log-likelihood under the base model as a measure of generation quality. Since the base model is finetuned on similar instruction-following dataset, its log-likelihood is a reasonable proxy for the utility.

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

422

395

396 397

### 5.3 PFT LEADS TO ROBUST MODELS, FOR FREE

PFT leads to the most robust model Figure 5 clearly shows models from PFT is the most robust across different attack datasets, when the key instruction appears farther away from the beginning.

For Llama models, we see the baseline methods struggle on all four attack datasets, while PFT
performs much better. For Gemma models, the baseline methods fail on the two TensorTrust
datasets, where PFT again shows great improvements. However, they remain robust to the
Gandalf tasks. We find that these attacks are too "weak" for Gemma. In fact, they are essentially
treated as ordinary data by the model: it has the same performance on those adversarial inputs
as on ordinary data when the user provides incorrect password (fig. 9).



(a) PFT does not hurt generation quality for ordinary data, as measured by the generation accuracy on the password dataset, and log-likelihood of generations on the Alpaca dataset.



(b) PFT does not lead to additional deviation from the base model, as measured by the KL divergence using Alpaca prompts.

Figure 6: Compared to the standard SFT, PFT does not hurt generation quality nor causes more deviation from the base model. See fig. 10 for the same results on Gemma models.

PFT **does not hurt model performance** One would worry that manipulating position IDs would hurt model performance: the shifted IDs may appear out-of-distribution for the model, and hurt model understanding and 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 base model. Therefore, relative to SFT, PFT improves the model robustness, for free.

#### 6 RELATED WORK

Attacks on LLMs Several works have studied the security of LLMs, and proposed various attacks to exploit the vulnerabilities of the models. The most relevant ones for closed-domain LLMs are prompt injection attacks (Willison, 2022; Yu et al., 2023; Geiping et al., 2024; Yu et al., 2024), which attempt to hijack the model by injecting malicious instructions, and system prompt extraction attacks (Sha & Zhang, 2024; Zhang & Ippolito, 2023; Yang et al., 2024), which aim to extract the system prompt from the model. These attacks could employ different techniques (Schulhoff et al., 2023; Perez & Ribeiro, 2022), so we need to evaluate the robustness of the model against a wide range of attacks. Fortunately, recent works collect diverse injection and extraction attack samples through online games (Toyer et al., 2023; Lakera AI, 2023a;b). We use these attack samples to evaluate the robustness of the models in our experiments. 

Finetuning methods Several works study how to finetune the LLM to defend against attacks
 in the user input. Our work studies the fragility of the finetuned models in terms of system
 prompt design, and proposes mitigating methods.

Wallace et al. (2024) finetunes the LLM to completely ignore user instructions for closed-domain tasks, and ignore conflicting instructions for open-domain tasks. We find that SFT-tuned models have fragility when the key instruction appears later in the input.

Chen et al. (2024) finetunes the LLM to completely ignore user instructions through structured query and a secure front-end. One key idea, using special delimiter tokens, is proposed to defend against Completion Attacks. We find that using special delimiter tokens also helps mitigate the fragility in this case, but not completely eliminates it. The PFT is more robust, and also does not require the front-end filtering of special tokens.

Positional encoding manipulation methods Recent advancements in long-context learning
have explored various positional encoding manipulation 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
manipulated position IDs after fine-tuning. This finding aligns with our own observations that

486 position-enhanced fine-tuning does not negatively impact model performance on standard-length data.

488 489

### 7 DISCUSSION AND CONCLUSION

This paper studies the security of closed-domain LLM bots, which play more and more important roles in human-computer interaction. While we find that simple SFT could make the model appear secure, stress-testing the prompt design reveals key fragilities. We hypothesize that this fragility is due to the current prompt format, which does not provide enough separation between system and user tokens. We propose a position-enhanced finetuning method PFT to mitigate this fragility, which significantly improves the out-of-distribution robustness of the model, while maintaining the same performance on ordinary data.

497
 498
 499
 Our work provides several implications for security alignment and safe deployment of LLMs in closed-domain settings.

Current instruction-following tuning may not let the model differentiate between different 500 roles Our experiments show that instruction-tuned models do not clearly distinguish bewteen 501 contents from different roles. This could be a potential security risk, when we need to strictly 502 enforce privilege hierarchies among different roles in the system. This is natural, both because of 503 the prompt format and the training data distribution. As we argue in the paper, the prompt format 504 does not provide enough distinction. Meanwhile, the training data often involve cases when the 505 model is expected to follow instructions from both the system and user. These combined lead 506 the model to treat the concatenated prompts as a stream of instructions, rather than a hierarchy 507 of instructions and data with different privilege levels. 508

Data (augmentation) cannot solve all the problems, at least not efficiently While data augmentation is a powerful tool to remove spurious correlations and improve model robustness, we cannot rely on it to solve all the problems, for several reasons. First, we need to identify the spurious correlations in the first place, which requires lots of testing, and becomes increasingly difficult as the model becomes more complex. Second, creating the perfect augmentation data can be challenging — as in our case, the augmented data helps, but not completely eliminates the fragility.

It is more efficient to have security baked in the model architecture Since data cannot solve 516 all the problems efficiently, it is better to have security baked in the architecture of the model. 517 The inductive bias could eliminate many of the spurious correlations from the beginning, and 518 make the model more robust to adversarial inputs that are out-of-distribution. More specifically, 519 since role differentiation is a fundamental requirement in many closed-domain systems, it is 520 better to have model architectures that can clearly differentiate between different roles. In this 521 project we consider the most simple system, consisting of one round dialogue between three 522 roles (system, user, and the model). In this case, we find that simply enhancing the position 523 information could help the model differentiate between different roles. 524

More complex systems might require more sophisticated solutions In more complex systems, where there might be multiple rounds of conversations between roles of different privilege levels, our simple manipulation of the position information might not be enough. In this case, we might need more complicated solutions to clearly delineate the roles and dialogues.

529 530

525

526

527

528

- 531
- 532
- 533
- 534
- 535
- 536
- 537
- 538
- 539

540	REFERENCES
541	REI ERENCEED

554

555

556

557

567

580

581

582

583

584

542	Arkadeep Acharya, Brijraj Singh, and Naoyuki Onoe. Llm based generation of item-description
543	for recommendation system. In Proceedings of the 17th ACM Conference on Recommender
544	Systems, pp. 1204–1207, 2023.

- AI@Meta. Llama 3 model card. 2024. URL https://github.com/meta-llama/llama3/
   blob/main/MODEL\_CARD.md.
- Anthropic. Introducing the next generation of claude, March 2024. URL https://www.
   anthropic.com/news/claude-3-family. Accessed September 30, 2024.
- Shouyuan Chen, Sherman Wong, Liangjian Chen, and Yuandong Tian. Extending context win dow of large language models via positional interpolation. *arXiv preprint arXiv:2306.15595*, 2023.
  - Sizhe Chen, Julien Piet, Chawin Sitawarin, and David Wagner. Struq: Defending against prompt injection with structured queries. *arXiv preprint arXiv:2402.06363*, 2024.
  - Chengguang Gan, Qinghao Zhang, and Tatsunori Mori. Application of llm agents in recruitment: A novel framework for resume screening. *arXiv preprint arXiv:2401.08315*, 2024.
- Jonas Geiping, Alex Stein, Manli Shu, Khalid Saifullah, Yuxin Wen, and Tom Goldstein.
  Coercing LLMs to do and reveal (almost) anything. *arXiv preprint arXiv:2402.14020*, 2024.
- Google. Gemini: A family of highly capable multimodal models, December 2023. URL https:
   //storage.googleapis.com/deepmind-media/gemini/gemini\_1\_report.pdf. Accessed September 30, 2024.
- 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. *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.
- Xinyu Lin, Wenjie Wang, Yongqi Li, Shuo Yang, Fuli Feng, Yinwei Wei, and Tat-Seng Chua.
   Data-efficient fine-tuning for llm-based recommendation. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pp. 365–374, 2024.
- Zabir Al Nazi and Wei Peng. Large language models in healthcare and medical domain: A
   review. In *Informatics*, volume 11, pp. 57. MDPI, 2024.
- 577
  578 OpenAI. Gpt-4 technical report, March 2023. URL https://arxiv.org/abs/2303.08774.
  579 Accessed September 30, 2024.
  - Bowen Peng, Jeffrey Quesnelle, Honglu Fan, and Enrico Shippole. Yarn: Efficient context window extension of large language models. *arXiv preprint arXiv:2309.00071*, 2023.
  - Fábio Perez and Ian Ribeiro. Ignore previous prompt: Attack techniques for language models. In *NeurIPS ML Safety Workshop*, 2022.
- Sander Schulhoff, Jeremy Pinto, Anaum Khan, Louis-François Bouchard, Chenglei Si, Svetlina
  Anati, Valen Tagliabue, Anson Liu Kost, Christopher Carnahan, and Jordan Boyd-Graber.
  Ignore this title and HackAPrompt: Exposing systemic vulnerabilities of Ilms through a
  global scale prompt hacking competition. In *EMNLP*, 2023.
- Zeyang Sha and Yang Zhang. Prompt stealing attacks against large language models. *arXiv* preprint arXiv:2402.12959, 2024.
- Karan Singhal, Tao Tu, Juraj Gottweis, Rory Sayres, Ellery Wulczyn, Le Hou, Kevin Clark,
   Stephen Pfohl, Heather Cole-Lewis, Darlene Neal, et al. Towards expert-level medical question answering with large language models. *arXiv preprint arXiv:2305.09617*, 2023.

594 595 596	Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca, 2023.
597 598 599	Gemma Team. Gemma. 2024. doi: 10.34740/KAGGLE/M/3301. URL https://www.kaggle. com/m/3301.
600 601 602	Sam Toyer, Olivia Watkins, Ethan Adrian Mendes, Justin Svegliato, Luke Bailey, Tiffany Wang, Isaac Ong, Karim Elmaaroufi, Pieter Abbeel, Trevor Darrell, et al. Tensor trust: Interpretable prompt injection attacks from an online game. <i>arXiv preprint arXiv:2311.01011</i> , 2023.
603 604 605 606	Eric Wallace, Kai Xiao, Reimar Leike, Lilian Weng, Johannes Heidecke, and Alex Beutel. The instruction hierarchy: Training llms to prioritize privileged instructions. <i>arXiv preprint</i> <i>arXiv:2404.13208</i> , 2024.
607 608 609 610	Isabella C Wiest, Fabian Wolf, Marie-Elisabeth Lessmann, Marko van Treeck, Dyke Ferber, Jiefu Zhu, Heiko Boehme, Keno K Bressem, Hannes Ulrich, Matthias P Ebert, et al. Llm-aix: An open source pipeline for information extraction from unstructured medical text based on privacy pre-serving large language models. <i>medRxiv</i> , pp. 2024–09, 2024.
611 612 613	Simon Willison. Prompt injection attacks against GPT-3, 2022. URL https://simonwillison. net/2022/Sep/12/prompt-injection/.
614 615	Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. Efficient streaming language models with attention sinks. <i>arXiv preprint arXiv:2309.17453</i> , 2023.
616 617 618 619	Yong Yang, Xuhong Zhang, Yi Jiang, Xi Chen, Haoyu Wang, Shouling Ji, and Zonghui Wang. Prsa: Prompt reverse stealing attacks against large language models. <i>arXiv preprint arXiv:2402.19200</i> , 2024.
620 621	Jiahao Yu, Yuhang Wu, Dong Shu, Mingyu Jin, and Xinyu Xing. Assessing prompt injection risks in 200+ custom gpts. <i>arXiv preprint arXiv:2311.11538</i> , 2023.
622 623 624	Jiahao Yu, Yangguang Shao, Hanwen Miao, Junzheng Shi, and Xinyu Xing. Promptfuzz: Harnessing fuzzing techniques for robust testing of prompt injection in llms. <i>arXiv preprint arXiv:2409.14729</i> , 2024.
625 626 627 628	Xiaoyu Zhang, Yishan Li, Jiayin Wang, Bowen Sun, Weizhi Ma, Peijie Sun, and Min Zhang. Large language models as evaluators for recommendation explanations. <i>arXiv preprint</i> <i>arXiv:2406.03248</i> , 2024.
629 630	Yiming Zhang and Daphne Ippolito. Prompts should not be seen as secrets: Systematically measuring prompt extraction attack success. <i>arXiv preprint arXiv:2307.06865</i> , 2023.
631 632 633 634 635	Yuyue Zhao, Jiancan Wu, Xiang Wang, Wei Tang, Dingxian Wang, and Maarten de Rijke. Let me do it for you: Towards llm empowered recommendation via tool learning. In <i>Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval</i> , pp. 1796–1806, 2024.
636 637 638	Dawei Zhu, Nan Yang, Liang Wang, Yifan Song, Wenhao Wu, Furu Wei, and Sujian Li. Pose: Efficient context window extension of llms via positional skip-wise training. <i>arXiv preprint</i> <i>arXiv:2309.10400</i> , 2023.
639 640 641 642	
643 644 645	