# 000 D-SPIN: DEFENSIVE SELF-SUPERVISED PROMPT 101 INJECTION

Anonymous authors

Paper under double-blind review

## Abstract

Large Language Models (LLMs) are increasingly used in a variety of important applications, yet their safety and reliability remain major concerns. Various adversarial and jailbreak attacks have been proposed to bypass the safety alignment and cause the model to produce harmful responses. We introduce Self-supervised Prompt INjection (D-SPIN) which can detect and reverse both prompt injection and jailbreak attacks on LLMs. Just by injecting an adaptive defense prompt at inference-time, our method is simple, effective, and compatible with existing safety-aligned models. Our benchmarks demonstrate that our system can reduce the attack success rate by up to 87.9%, while maintaining the performance on benign user requests. In addition, we discuss the situation of an adaptive attacker and show that our method is still resilient against attackers who are aware of our defense.



Figure 1: Self-Supervised Prompt Injection for LLM Defense. Large language models can be jailbroken by an adversarial prompt. In this paper, we show that we can detect and defend against jailbreak and prompt injection attacks by leveraging self-supervised language tasks. The red text shows the adversarial prompt injection to bypass the safe guardrails of LLMs. The blue box shows our defensive method D-SPIN which safeguards the LLMs.

## 1 INTRODUCTION

Large-Language Models (LLMs) have achieved great success in code generation (Rozière et al., 2023), question answering (Touvron et al., 2023), and task planning (Surís et al., 2023). However, they can be jailbroken by adversarial attacks and prompts (Zou et al., 2023; Zeng et al., 2024; Yuan et al., 2023), leading them to exhibit behaviors of manipulation (Carroll et al., 2023; Ji et al., 2024), deception (Park et al., 2023; Hazell, 2023), and explicit content (Kang et al., 2023). This poses an inherent risk as LLMs are widely deployed in many critical applications, like programming.

To align LLMs for safety, prior work focuses on training to make the model more robust. For example, Reinforcement Learning from Human Feedback (RLHF) (Ouyang et al., 2022) trains reward
models based on human ratings, and is used to judge LLM outputs. DPO (Rafailov et al., 2023)
trains model on direct preference optimization. To make train-time defenses robust against attacks,
it is often required to train the model on the generated attacks, which can take significant resources
to gather responses (Casper et al., 2023). In addition, training-time defenses have fundamental difficulties defending against unforeseen attacks and jailbroken prompt injection, due to a mismatched



Figure 2: **Self-supervised Prompt Injections**. Examples of inputs at each stage of the defense. The red text is the adversarial suffix. On the left side, all the malicious examples with the trigger attack are discovered or reversed. The repeat task fails to reproduce the adversarial suffix in red because the token sequence becomes too random. Even though it does repeat the sentence partially, the difference is enough that the detection system will detect it as being malicious. The blue under the 'Reversal' task represent our defensive prompt injections which reverses malicious attacks, but does not impact benign inputs. On the right side, we show the same benign question passing through the layers without issue. Our defense systems is shown to defend against attacks that would otherwise break alignment, while maintaining performance on benign queries.

083

081

074

075

076

077

079

generalization between the smaller safety training data and the broader pre-training data (Wei et al., 2023).

In this paper we introduce a method to defend against these attacks by self-supervised prompt injection, dubbed D-SPIN allowing us to achieve a defense that adapts to each attack online. Just as an attacker finds the right prompt injection to break the LLM, our approach will find the right prompt injection to detect and repair the input. Figure 1 shows our prompt injection defending an adversarial prompt injection.

Our key insight is that the prompts which override the natural guidelines of the LLM also degrade other capabilities of the model. By constructing various self-supervised tasks, we can detect attacks without needing to know the true labels of the input requests or examples. Our constructed selfsupervised tasks also allow us to repair attacked user prompts that slip through detection. This is achieved by adding tokens to the start of the input and taking gradient steps to reduce the perplexity of the resulting string. By using self-supervised prompt injection at test time, we can defend against unseen attacks and even adaptive attacks given the online adaptive property. Figure 2 shows the results of the different tasks rejecting malicious attacks while letting benign requests through.

- Our method is model agnostic and can be compatible with existing defense systems. Our selfsupervised criteria do not depend on malicious or benign labels, thus can be used online at inference time on demand.
- Visualizations, experiments, and qualitative results demonstrate that our self-supervised prompt injection significantly improves the LLM safety upon existing alignment defenses. On Advbench with
  Universal Adversarial Triggers (Zou et al., 2023), we are able to reduce Attack Success Rates (ASR)
  to 12.11% and 0% on Vicuna and Llama-2 respectively. Our method also can detect natural language
  based jailbreak attacks (Wei et al., 2023), adversarial instructions (Wei et al., 2023), and role-play
  (Shen et al., 2024), reducing ASR to 0% on both models, fully protecting them. Our method also is
  shown to be resilient against attackers who are aware of our defense.



Figure 3: Self Supervised Detection of Jailbreak Attacks. We detect jailbreak prompts through self-supervised natural language tasks. Here we show the 'Repeat' task and the 'Interject' task that form the detection layer. We are able to verify if the prompt passes the tests since we know the answer beforehand, either repeating the entire input or mentioning that Paris is the capital of France. Since jailbreak prompts often inject unusual patterns into the text, the model will often underperform on our created self-supervised language tasks compared to a benign user query. Our method uses this difference for detection.

130 131 132

129

## 2 RELATED WORK

- Large-Language Models: LLMs are pre-trained models on a huge corpus of textual data. A mix of these will have open source weights (Touvron et al., 2023), while others are closed off and are only accessible through APIs (Liu et al., 2023b). Studies have shown that various capabilities are learned at the pre-training stage (Zhou et al., 2023). Performance of these models is correlated with parameter size and train time (Hoffmann et al., 2022; Kaplan et al., 2020), however, larger models also exhibit more problematic behavior (Perez et al., 2022; Elazar et al., 2023) and even the fine-tuning of these large models on benign datasets can compromise their safety (Qi et al., 2023).
- 139 Alignment: Alignment is an inbuilt system in recent LLMs. Originally used for text summarization 140 (Stiennon et al., 2022), the main motivation recently has been to align the goals of AI models with 141 human values and morals (Leike et al., 2018; Hendrycks et al., 2022). It can be done through 142 reinforcement learning through human feedback (RLHF) or even through other LLMs as judges 143 (Ziegler et al., 2020; Lee et al., 2023). However, as long as alignment does not remove the harmful 144 behavior entirely, there will exist prompts that can elicit that behavior (Wolf et al., 2023). Wei et al. 145 (2023) showed that two main issues of alignment come from conflicting objectives, and mismatched 146 generalization.
- 147 Adversarial Attacks: Specific modifications to inputs can cause models to output incorrect or ma-148 licious responses. This first gained prominence with Szegedy et al. (2014)'s paper attacking image 149 classification models. For vision tasks, these attacks can also be repaired through natural super-150 vision as well (Mao et al., 2021). For Large-language models, research has been done compiling 151 jailbreak prompts that can cause harmful responses (Gehman et al., 2020; Liu et al., 2023a). These 152 refer to user-created prompts, and usually involve a scenario that encourages the LLM to break away 153 from the preset safety guidelines (Wei et al., 2023). Recent developments such as AutoDAN (Liu et al., 2024) and CodeChameleon (Lv et al., 2024) incorporate LLMs to generalize and automatically 154 generate these malicious prompts for a variety of requests. 155
- Universal Adversarial Triggers (UATs) (Wallace et al., 2021) has been one method to attack these
  models which embeds a sequence of tokens into the user prompt that will cause it to bypass alignment and safety guidelines. In particular, gradient-guided search (Guo et al., 2021) has proven
  remarkably effective in adversarially misleading the model. Zou et al. (2023) show that these attacks are also transferable across models, allowing harmful generation even in closed-box models.
  In response, various defense methods have been suggested, however, many of them (Kumar et al., 2023; Li et al., 2023; Phute et al., 2023) rely on external LLMs to verify the safety of the request.

162 **Defense Methods** exist for guarding against certain forms of Adversarial attacks. One example is 163 Jain et al. (2023) showing a perplexity filter is able to reject certain token attacks that cross a given 164 score. Another approach have been defenses (Cao et al., 2023; Robey et al., 2024) that randomly 165 drop or modify tokens within the input to try and trigger alignment. For shorter prompts, such 166 methods risk dropping important nouns or context, which can classify benign prompts as malicious. These approaches are similar to ones that preprocess the inputs in a way that adds a secret passphrase 167 (Liu et al., 2024) or strictly structures the query Chen et al. (2024). As a layer of defense, some 168 researchers (Inan et al., 2023) proposed adding another LLM to either judge or classify an input as benign before returning the answer. 170

171

## 3 Method

172 173

We begin by describing a series of defenses that use self-supervision to defend against both handcrafted jailbreak attacks and optimization-based adversarial trigger attacks (Zou et al., 2023). We provide two layers of defense, detection and reversal, that complement each other. Both methods rely on self-supervised signals that we construct for language tasks. These signals are shown in Figure 3 being able to accurately distinguish between benign and malicious prompts. Lastly, we consider the scenarios where the attacker can adapt their attack in response to the defense, and empirically show that our methods are still robust.

181

## 182 3.1 JAILBROKEN LANGUAGE MODELS

In a given user interaction, the user inputs the request  $\mathbf{x}$ , and the chatbot responds with response  $\mathbf{y}$ . The LLM generates text outputs by predicting the next token, sampling from a probability distribution conditioned on the input such that  $\mathbf{y} \sim F(\cdot|\mathbf{x})$ .

For **attacks**, the attacker wants to minimize the difference between the response y and their desired goal t. In this case, this would be malicious instructions prefaced by an agreement such as "Sure, here's how". To do this, they can add a set of tokens a which causes the model to respond with their desired output. This results in an attack optimization goal of:

$$\mathbf{a} = \min_{\mathbf{A}} \mathcal{L}_{\text{attack}}(F(\mathbf{x} + \mathbf{a}), \mathbf{t})$$
(1)

Note here that '+' denotes the concatenation of two strings, and that the attack is appended to the end
of the user request. For GCG attack (Zou et al., 2023), this optimization is through gradient descent.
For natural language jailbreaks (Wei et al., 2023), this optimization is through human construction.

196 197

199

200

201

202

203

204

211

212 213

191 192

## 3.2 DETECT JAILBREAKS VIA SELF-SUPERVISION

We find jailbreaks prompt to bypass the LLM alignment contains different structure than benign input. We find jailbreak input often leave a trace, such as degrades other capabilities of the LLM in order to achieve a successful attack. Motivated by this, we propose to construct tasks that we already know the groundtruth or the expected behavior, so that we can detect jailbreak if we observe a violation of those tasks. While there are many self-supervised task in natural language we can construct, here we show self-supervised tasks as defense measures.

**Repeat:** One simple request of LLMs is to ask it to repeat back the input. For most requests, there is no issue at all, and it is only when the input is harmful or it is under attack that the LLM will have trouble doing so. We use Levenshtein Distance " $lev(a_{1:i}, b_{1:j})$ " between the original user request and the newly generated sentence. This distance is then weighted by the average length of the two strings, where s(x) denotes the size of the string, to compensate for differences in request length.

$$\mathcal{L}_{\text{repeater}}(\mathbf{x}') = \frac{2 \operatorname{lev}(\mathbf{x}', F(\mathbf{x}'))}{s(\mathbf{x}') + s(F(\mathbf{x}'))}$$
(2)

Interjection: At the end of the user request, another method is to interject with a separate question
 where we know the correct answer. This utilizes the internal knowledge set of the model, testing if it can answer fluently. We hypothesis that if the LLM has been jailbroken, they also degrade their



Figure 4: Self-supervised loss and ROC Curve for benign and malicious inputs. We show  $\mathcal{L}_{\text{repeater}}$  loss on the left and  $\mathcal{L}_{\text{interject}}$  on the right. For the repeat task, we find the model can repeat benign queries from the TriviaQA task well, but fails to repeat consistently when the query is from Advbench. For the interject task, the model is asked to answer 'What is the capital of France?' and the loss is for predicting whether the next token will be 'Paris'. As seen by the malicious instruction prompts, the loss increases only with the addition of the suffix. The divergent self-supervised loss values between malicious and benign allow us to detect malicious query. For the ROC curves, benign inputs are categorized as negatives, and malicious as positive. Losses for these graphs are calculated using Llama-2, with the best thresholds for the loss being 0.90 and 5.73 respectively.

capability in ansering those questions, like "What is the capital of France?" By calculating softmax for the logits for 'Paris' as the next token, it provides the following loss function:

$$\mathcal{L}_{\text{interject}}(\mathbf{x}', \mathbf{y}) = \frac{e^{P(\mathbf{y}|\mathbf{x}')}}{\sum\limits_{\mathbf{y}\in\mathbf{Y}} e^{P(\mathbf{v}|\mathbf{x}')}}$$
(3)

Each of these loss functions then cluster the inputs that are harmful and benign. Figure 4 shows the different loss functions for the repeater and interjection tasks. By placing a threshold T on each of these self-supervised tasks, we are able to detect when a model has its capabilities attacked.

#### 3.3 REVERSALS

230

231

232

233

234

235

236

237

238 239

240

245

246

247 248

249

251

258

250 As a second layer for defending against adversarial attacks, we can even reverse attacks that slip through detection. To do so, we append additional tokens in front of the request  $\mathbf{x}'$ , and aim to find a 252 series of tokens that will restore the natural alignment of the LLM. This would allow it to work with 253 existing defenses, and further strengthens the models against attacks. Our proposed method, the 254 introduction of defense tokens into the prefix of the user prompt also leaves the core content of the 255 request unchanged. One observation we made with the GCG adversarial triggers (Zou et al., 2023) was that their inclusion would drastically increase the total perplexity of the input. By defining our 256  $\mathcal{L}_{autoreg}$  as the perplexity of the entire user input, we get the following objective: 257

$$\mathbf{d} = \arg\min_{\mathbf{d}} \mathcal{L}_{\text{autoreg}}(\mathbf{d} + \mathbf{x}') \tag{4}$$

259 After every p stages of optimization, we run the completion of the new input request  $\hat{\mathbf{y}} = F(\mathbf{d} + \mathbf{x}')$ 260 and record if the response began with any of the common alignment refusals. To pass through, it 261 must run through all n steps without triggering any of the common denial prefixes. Once the optimal 262 d is found, the final output of the model is  $\hat{\mathbf{y}} = F(\mathbf{d} + \mathbf{x}')$ . 263

264 3.4 TOTAL SYSTEM 265

266 For an attack to be successful, it must pass through all the detection stages, and the generated re-267 sponse will have the defense tokens as the prefix. The multi-layered approach allows each defense system to plug up holes in the others, and adds to the variables that an attack must consider. Due to 268 the compositionality of our self-supervised defense, users can easily remove some tasks to speed up 269 the system.

#### 270 3.5 ADAPTIVE ATTACKER 271

272 Another key advantage of our defense is that it is also robust to adaptive attackers. If the attack 273 ignores our defense strategy, then we will defend it. If the attackers consider our defense strategy and then attack, then they need to optimize to also respect the self-supervised tasks, therefore make their 274 optimization harder due to the additional constraints. Formally, an adaptive attacker will alternate 275 between attacking and defending the same prompt until it converges: 276

$$\mathbf{a} = \arg\min\mathcal{L}_{\text{attack}}(F(\mathbf{x} + \mathbf{a}), t)$$
(5)

(6)

278

282

283

284

285 286 287

288

289

290

291

 $\mathbf{d} = \arg\min_{\mathbf{d}} \mathcal{L}_{autoreg}(\mathbf{d} + \mathbf{x} + \mathbf{a})$ This iterative process is computationally intensive, and has the equivalent result as a constrained optimization problem. We use the Lagrangian penalty method by adding a new penalty to the score function. For each new defense system, we can include it as a new multi-target optimization the attack must consider. With all the detection layers and reversal implemented, the adaptive attack would have a multi-target optimization problem as follows:

$$\mathcal{L}_{\text{adapt}}(x', t, \lambda_r, \lambda_i, \lambda_s, \lambda_p) = \mathcal{L}_{\text{attack}}(x', t) + \lambda_r \mathcal{L}_{\text{repeater}}(x') + \lambda_i \mathcal{L}_{\text{interject}}(x') + \lambda_p L_{\text{autoreg}}(x')$$
(7)

By allocating more optimization budget to satisfying these constraints, the resulting attack is also less efficient. There is a trade-off that the attacker must consider: by optimizing to bypass the detection layers, it also increases the coherence of the input and constraint the input the be benign. If the attacker does not optimize to bypass the detection layers, our defense algorithm will be activated due to the violation of our self-supervised constraints. Thus, adding our self-supervised constraint will result in a lose-lose situation for the attacker, and we will straightly increase the safety.

292 293

295

296 297

298

299

303

304

310

#### 4 **EXPERIMENTS**

4.1 DATASETS

Advbench 'harmful behaviors' contains 520 malicious requests, which we then added adversarial triggers by following Zou et al. (2023)'s method of multi-prompt attack.

300 **TriviaQA** (Joshi et al., 2017) contains a series of trivia questions crafted by humans along with an evidence dataset that contains the answer. We evaluate our performance on the wikipedia validation 301 set with one-shot learning and closed-book performance after reversal. 302

4.2 MODELS

305 Llama 2-chat is a chatbot released by Meta, finetuned for dialogue and aligned through human 306 feedback (Touvron et al., 2023). 307

Vicuna-7b is an open-source chatbot released by LMSYS based on finetuning Llama 2 with con-308 versations from ShareGPT (Zheng et al., 2023). 309

4.3 ATTACKS 311

312 Universal Adversarial Trigger. We implement Greedy Coordinate Gradient-based (GCG) search 313 as done in (Zou et al., 2023). Each malicious request gets appended a suffix, and the suffix gets 314 optimized for up to 500 iterations or until alignment is broken. Then the next request begins with 315 the previous optimized suffix appended. In the case of Vicuna, the suffix injection is able to achieve 316 an Attack Success Rate (ASR) of 100%. 317

Natural Language Jailbreak. Jailbreak Chat is a website for users to share handcrafted prompts 318 that can bypass AI alignment. We take the top 5 upvoted jailbreak prompts and pair them each with 319 30 randomly selected Advbench requests for a total of 150 prompts. These prompts tend to be older 320 attacks, and were published before Llama2 released. 321

Adversarial Instructions. These attacks bypass alignment by circumventing the priorities of the 322 LLM as mentioned in Wei et al. (2023). We test the strongest attack combination presented in their 323 paper including prefix injection and refusal suppression, and pair it with 150 Advbench samples.



Figure 5: **Multi-Benchmark Attack Success Rate (ASR)**. We show performance of our defense against multiple types of attacks and benign inputs (Except for TriviaQA, lower indicates better performance). Our defense is robust and dramatically lowers ASR across all attacks. Even in cases where alignment is effective, D-SPIN is shown to decrease the ASR further.



Figure 6: Adaptive Attack Success Rate. The attacker minimizes the loss with respect to a Lagrangian score function, which balances between attacking the alignment, and bypassing the defenses we created. Then we run the same defenses on the newly generated attack suffixes. We vary the  $\lambda$  penalty score for layers of the defense and measure the attack success rate. While the adaptive attacks do lower the ability of the specific layer, it leads to easier detection by the other layers.

**Multiple Role-Play**. By simulating role playing in conversation, attackers instruct the LLM to act out a role which is able to bypass safety guardrails. This type of attack is similar to 'Do Anyhing Now' (DAN) attacks in (Shen et al., 2024). We pair the 'dev mode v2' (Wei et al., 2023) attack with the same 150 malicious requests from Advbench.

Automated Jailbreaks. Code Chameleon (Lv et al., 2024), AutoDAN (Liu et al., 2024), and ICA
 (Wei et al., 2024), are all automated ways of generating natural language jailbreaks. The attacks
 commonly have a starting attack template that is gradually modified with the attack query, similar
 to GCG-search, to bypass defense systems. We pair each attack with 150 malicious requests from
 Advbench.

4.4 RESULTS

367
 368
 369
 370
 Detection. For our Repeater layer, we inject the string "Repeat only the following sentences exactly without changes: " into the model with the entire input surrounded by quotation marks. We use the best threshold of 0.89 to filter out inputs that failed to correctly repeat.

For our Interjection layer, we chose the question to be "But first, answer only what is the capital of France?". This gets appended to the end of the model, which leaves the system prompt unchanged. On benign cases, the model was able to answer correctly more than 99% of the time, making it a reasonable question to ask. It passes the layer if the generation returns the answer within the first 10 tokens, and in practice, this gave similar results to using the loss threshold.

The ROC curves in figure 4 show that there is a clear delineation in loss between malicious and benign requests when asked to repeat back the entire string. With only the malicious requests, the LLM will identify the harmful content and refuse to repeat the sentence. When the adversarial tokens do appear, the repeated string will usually not include the suffix tokens, leading to a higher
 loss. This is because when predicting the next word, the actual token that follows is too random to
 be reliably chosen.

As for the interjection task, it makes sense that the benign and malicious prompts without the attack have similar losses. This is because the model still has all of its capabilities, and alignment will still work. Only when an attack disturbs the coherency of the input does the loss go up. In our implementation, we classify the request based on whether it is able to generate the token 'Paris' within the first 10 tokens. This results in a False Positive Rate (Benign requests classified as harmful), of less than 1%. While the malicious queries without attacks might be able to answer the question, it also triggers the alignment, which correctly prevents them from generating the 'Paris' token.

Reversal. For each request, we append 5 tokens "!!!!!" to the start of the user message. For 25 steps of optimization, at each stage the reversal calculates the top-256 tokens that best decrease the perplexity of the overall sentence. Then for our batch size of 50, we uniformly pick from one of the 256 tokens and substitute it in. The sequence with the lowest loss is then chosen. At every 5 iterations it checks if the generated output matches one of the common denials for Vicuna or Llama. In both models, the reversal alone is able to reduce ASR by nearly half.

394 **Jailbreaks** present an interesting scenario. Because these are handcrafted prompts, they usually 395 involve role-play and bypass guidelines by misleading the priorities of the model. The sentence 396 structure is coherent, and as a result, they have the same perplexity as benign inputs. This means 397 that our reversal is less effective because minor perturbations do not affect the overall understanding 398 by the model. Most of these prompts also tell the model to 'Ignore all the instructions you got 399 before', which is needed to weaken alignment or safety reminders in the system prompt. In the case 400 of the repeat task, this actually causes the attack to fail that detection layer, because the model no 401 longer repeats the input, and the Levenshtein distance will be above the threshold. As seen in Figure 5, the repeat condition is strong enough for D-SPIN to prevent all attacks. 402

403 Performance on TriviaQA Figure 5 shows that under the benign cases, the overall performance is
 404 not significantly affected. Because of how we chose the thresholds for D-SPIN detection, the benign
 405 cases passed through without any issues. The addition of tokens that lower the perplexity does not
 406 impact how the model answers the question.

407 408 409

4.5 DEFENSE AWARE ADAPTIVE ATTACK

We present the adaptive attack success rates in Figure 6. The same defense is run, but depending on the  $\lambda$  penalty chosen, the attacker can optimize the attack for certain layers. However, this usually comes with a lower ASR in other layers. This is because the attacker has less of a budget to tackle the original goal of breaking alignment. Our defense system still reduces the attack success rate by up to 76%, even under an adaptive attacker, proving that the method is resilient.

415

422

416 4.6 ABLATION STUDY

Computational Overhead With the detection layers, on benign inputs the model goes through two iterations of token generation. As shown in Table 1, the interjection task and the repeat task causes inference to run up to 1.2x and 1.3x longer, respectively. Reversal is 4.5x longer than inference when doing 25 steps of optimization, and using less steps can significantly speed up inference, however it is important that attacks do not succeed for the sake of a faster response time. Compared to the cost

of running the GCG attack, our reversal is still only 1% of the cost.

423 Steps of optimization. There is an inherent 424 trade-off between optimization steps and com-425 putation speed. Looking at Figure 7, the increase 426 in steps does lead to a lower ASR. However, 427 when examining the defense as a whole, the re-428 versal only comes into effect if it slips through the previous detection layers. Under the assump-429 tion that detection will remove a large portion of 430 attacks, going beyond 25 for the marginal im-431 provements seems less appealing.

Table 1: Latency (Averaged across 100 samples).The table shows the latency times for various defense methods.

derense methods.		
	Method	Time (s)
	Standard Inference	2.80
	Repeat	0.45
	Interjection	0.89
	Reversal (25 steps)	12.62
	Full Defense	16.31



Figure 7: Ablation Attack Success Rate for UAT and Jailbreak. On the left we show performance across multiple combinations of the defense system. In our defense, each of the layers complements the others, so that when placed sequentially it catches the majority of attacks, and is robust against adaptive attackers. On the right we show ASR for only the reversal layer with increasing iterations, both Vicuna and Llama2 level off in terms of gains from increasing optimization steps.

### 5 CONCLUSION

In this paper, we show that self-supervised metrics are sufficient in detecting adversarial attacks, and that the attacks are also repairable by targeting perplexity. We introduce an inference time defense system to detect and repair input generation through prompt injection. Our defense works with existing models and alignment, while not requiring any additional training or fine-tuning. This ensures that our defense system can react to new types of attacks that were not present in the alignment training set. The creation of these self-supervised defenses allow for multiple variations that can detect attacks. The layered nature of the method also shows resilience against adaptive attackers.

### 6 ETHICS STATEMENT

This paper presents work on Large-Language Model safety, in hopes of presenting results that will make future conversational models safer to use. We acknowledge that our discussion does contain mentions of attacks and harmful content, but it is done to illustrate how our techniques specifically address those problems. These attacks have already been published, and we believe that there is value in bringing up discussions of new defense methods, even if it might bring increased promi-nence to these attacks. Whether researchers should publicly share defense methods because it might help attackers is a point up for discussion, but we do present the scenario of an adaptive attacker to help alleviate those concerns. There is already substantial discussion about the safety and prevalence of conversational AI, and the models we tested on are already pub-licly available for use. Our meth-ods are designed with practicality in mind, but they cannot account for every scenario in real-life use cases. As with any other safety system, the methods we present must be constantly monitored, and updated if necessary as attackers discover new avenues. This is to ensure that the work this paper produces does not lead to eventual harm. 

# 486 REFERENCES

502

515

- Bochuan Cao, Yuanpu Cao, Lu Lin, and Jinghui Chen. Defending against alignment-breaking at tacks via robustly aligned llm, 2023.
- Micah Carroll, Alan Chan, Henry Ashton, and David Krueger. Characterizing manipulation from ai systems, 2023.
- Stephen Casper, Xander Davies, Claudia Shi, Thomas Krendl Gilbert, Jérémy Scheurer, Javier Rando, Rachel Freedman, Tomasz Korbak, David Lindner, Pedro Freire, Tony Wang, Samuel Marks, Charbel-Raphaël Segerie, Micah Carroll, Andi Peng, Phillip Christoffersen, Mehul Damani, Stewart Slocum, Usman Anwar, Anand Siththaranjan, Max Nadeau, Eric J. Michaud, Jacob Pfau, Dmitrii Krasheninnikov, Xin Chen, Lauro Langosco, Peter Hase, Erdem B1y1k, Anca Dragan, David Krueger, Dorsa Sadigh, and Dylan Hadfield-Menell. Open problems and fundamental limitations of reinforcement learning from human feedback, 2023.
- Sizhe Chen, Julien Piet, Chawin Sitawarin, and David Wagner. Struq: Defending against prompt injection with structured queries, 2024. URL https://arxiv.org/abs/2402.06363.
- Yanai Elazar, Akshita Bhagia, Ian Magnusson, Abhilasha Ravichander, Dustin Schwenk, Alane
  Suhr, Pete Walsh, Dirk Groeneveld, Luca Soldaini, Sameer Singh, Hanna Hajishirzi, Noah A.
  Smith, and Jesse Dodge. What's in my big data?, 2023.
- Samuel Gehman, Suchin Gururangan, Maarten Sap, Yejin Choi, and Noah A. Smith. Realtoxici typrompts: Evaluating neural toxic degeneration in language models, 2020.
- 509 Chuan Guo, Alexandre Sablayrolles, Hervé Jégou, and Douwe Kiela. Gradient-based adversarial
   510 attacks against text transformers, 2021.
- <sup>511</sup> Julian Hazell. Spear phishing with large language models, 2023.
- Dan Hendrycks, Nicholas Carlini, John Schulman, and Jacob Steinhardt. Unsolved problems in ml safety, 2022.
- Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hennigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Jack W. Rae, Oriol Vinyals, and Laurent Sifre. Training compute-optimal large language models, 2022.
- Hakan Inan, Kartikeya Upasani, Jianfeng Chi, Rashi Rungta, Krithika Iyer, Yuning Mao, Michael
  Tontchev, Qing Hu, Brian Fuller, Davide Testuggine, and Madian Khabsa. Llama guard: Llmbased input-output safeguard for human-ai conversations, 2023. URL https://arxiv.org/
  abs/2312.06674.
- Neel Jain, Avi Schwarzschild, Yuxin Wen, Gowthami Somepalli, John Kirchenbauer, Ping yeh Chiang, Micah Goldblum, Aniruddha Saha, Jonas Geiping, and Tom Goldstein. Baseline defenses for adversarial attacks against aligned language models, 2023.
- Jiaming Ji, Tianyi Qiu, Boyuan Chen, Borong Zhang, Hantao Lou, Kaile Wang, Yawen Duan,
  Zhonghao He, Jiayi Zhou, Zhaowei Zhang, Fanzhi Zeng, Kwan Yee Ng, Juntao Dai, Xuehai
  Pan, Aidan O'Gara, Yingshan Lei, Hua Xu, Brian Tse, Jie Fu, Stephen McAleer, Yaodong Yang,
  Yizhou Wang, Song-Chun Zhu, Yike Guo, and Wen Gao. Ai alignment: A comprehensive survey,
  2024.
- Mandar Joshi, Eunsol Choi, Daniel S. Weld, and Luke Zettlemoyer. Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Vancouver, Canada, July 2017. Association for Computational Linguistics.
- 539 Daniel Kang, Xuechen Li, Ion Stoica, Carlos Guestrin, Matei Zaharia, and Tatsunori Hashimoto. Exploiting programmatic behavior of llms: Dual-use through standard security attacks, 2023.

549

550

551

554

566

567

568

569

- Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child,
   Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language
   models, 2020.
- Aounon Kumar, Chirag Agarwal, Suraj Srinivas, Aaron Jiaxun Li, Soheil Feizi, and Himabindu Lakkaraju. Certifying llm safety against adversarial prompting, 2023.
- Harrison Lee, Samrat Phatale, Hassan Mansoor, Thomas Mesnard, Johan Ferret, Kellie Lu, Colton
  Bishop, Ethan Hall, Victor Carbune, Abhinav Rastogi, and Sushant Prakash. Rlaif: Scaling reinforcement learning from human feedback with ai feedback, 2023.
  - Jan Leike, David Krueger, Tom Everitt, Miljan Martic, Vishal Maini, and Shane Legg. Scalable agent alignment via reward modeling: a research direction, 2018.
- Yuhui Li, Fangyun Wei, Jinjing Zhao, Chao Zhang, and Hongyang Zhang. Rain: Your language
   models can align themselves without finetuning, 2023.
- Xiaogeng Liu, Nan Xu, Muhao Chen, and Chaowei Xiao. Autodan: Generating stealthy jailbreak
   prompts on aligned large language models, 2024. URL https://arxiv.org/abs/2310.
   04451.
- Yi Liu, Gelei Deng, Zhengzi Xu, Yuekang Li, Yaowen Zheng, Ying Zhang, Lida Zhao, Tianwei
   Zhang, and Yang Liu. Jailbreaking chatgpt via prompt engineering: An empirical study, 2023a.
- Yiheng Liu, Tianle Han, Siyuan Ma, Jiayue Zhang, Yuanyuan Yang, Jiaming Tian, Hao He, Antong Li, Mengshen He, Zhengliang Liu, Zihao Wu, Lin Zhao, Dajiang Zhu, Xiang Li, Ning Qiang, Dingang Shen, Tianming Liu, and Bao Ge. Summary of chatgpt-related research and perspective towards the future of large language models. *Meta-Radiology*, 1(2):100017, September 2023b. ISSN 2950-1628. doi: 10.1016/j.metrad.2023.100017. URL http://dx.doi.org/10.1016/j.metrad.2023.100017.
  - Huijie Lv, Xiao Wang, Yuansen Zhang, Caishuang Huang, Shihan Dou, Junjie Ye, Tao Gui, Qi Zhang, and Xuanjing Huang. Codechameleon: Personalized encryption framework for jailbreaking large language models, 2024. URL https://arxiv.org/abs/2402.16717.
- 570 Chengzhi Mao, Mia Chiquier, Hao Wang, Junfeng Yang, and Carl Vondrick. Adversarial attacks are
   571 reversible with natural supervision. In *Proceedings of the IEEE/CVF International Conference* 572 on Computer Vision (ICCV), pp. 661–671, October 2021.
- Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback, 2022.
- Peter S. Park, Simon Goldstein, Aidan O'Gara, Michael Chen, and Dan Hendrycks. Ai deception:
   A survey of examples, risks, and potential solutions, 2023.
- 580 Ethan Perez, Sam Ringer, Kamilé Lukovsiūtė, Karina Nguyen, Edwin Chen, Scott Heiner, Craig Pet-581 tit, Catherine Olsson, Sandipan Kundu, Saurav Kadavath, Andy Jones, Anna Chen, Ben Mann, 582 Brian Israel, Bryan Seethor, Cameron McKinnon, Christopher Olah, Da Yan, Daniela Amodei, Dario Amodei, Dawn Drain, Dustin Li, Eli Tran-Johnson, Guro Khundadze, Jackson Kernion, 583 James Landis, Jamie Kerr, Jared Mueller, Jeeyoon Hyun, Joshua Landau, Kamal Ndousse, Lan-584 don Goldberg, Liane Lovitt, Martin Lucas, Michael Sellitto, Miranda Zhang, Neerav Kingsland, 585 Nelson Elhage, Nicholas Joseph, Noemí Mercado, Nova DasSarma, Oliver Rausch, Robin Lar-586 son, Sam McCandlish, Scott Johnston, Shauna Kravec, Sheer El Showk, Tamera Lanham, Timothy Telleen-Lawton, Tom Brown, Tom Henighan, Tristan Hume, Yuntao Bai, Zac Hatfield-Dodds, 588 Jack Clark, Samuel R. Bowman, Amanda Askell, Roger Grosse, Danny Hernandez, Deep Gan-589 guli, Evan Hubinger, Nicholas Schiefer, and Jared Kaplan. Discovering language model behaviors 590 with model-written evaluations, 2022. 591
- Mansi Phute, Alec Helbling, Matthew Hull, ShengYun Peng, Sebastian Szyller, Cory Cornelius, and
   Duen Horng Chau. Llm self defense: By self examination, llms know they are being tricked, 2023.

594 595 596	Xiangyu Qi, Yi Zeng, Tinghao Xie, Pin-Yu Chen, Ruoxi Jia, Prateek Mittal, and Peter Henderson. Fine-tuning aligned language models compromises safety, even when users do not intend to!, 2023.
597 598 599 600	Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D Manning, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. <i>arXiv preprint arXiv:2305.18290</i> , 2023.
601 602 603	Alexander Robey, Eric Wong, Hamed Hassani, and George J. Pappas. Smoothllm: Defending large language models against jailbreaking attacks, 2024. URL https://arxiv.org/abs/2310.03684.
605 606 607 608 609	Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, Artyom Kozhevnikov, Ivan Evtimov, Joanna Bitton, Manish Bhatt, Cristian Canton Ferrer, Aaron Grattafiori, Wenhan Xiong, Alexandre Défossez, Jade Copet, Faisal Azhar, Hugo Touvron, Louis Martin, Nicolas Usunier, Thomas Scialom, and Gabriel Synnaeve. Code Ilama: Open foundation models for code, 2023.
610 611 612	Xinyue Shen, Zeyuan Chen, Michael Backes, Yun Shen, and Yang Zhang. "do anything now": Characterizing and evaluating in-the-wild jailbreak prompts on large language models, 2024.
613 614	Nisan Stiennon, Long Ouyang, Jeff Wu, Daniel M. Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul Christiano. Learning to summarize from human feedback, 2022.
615 616 617	Dídac Surís, Sachit Menon, and Carl Vondrick. Vipergpt: Visual inference via python execution for reasoning. <i>arXiv preprint arXiv:2303.08128</i> , 2023.
618 619	Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks, 2014.
621 622 623	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko- lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open founda- tion and fine-tuned chat models. <i>arXiv preprint arXiv:2307.09288</i> , 2023.
624 625	Eric Wallace, Shi Feng, Nikhil Kandpal, Matt Gardner, and Sameer Singh. Universal adversarial triggers for attacking and analyzing nlp, 2021.
626 627 628	Alexander Wei, Nika Haghtalab, and Jacob Steinhardt. Jailbroken: How does llm safety training fail? <i>arXiv preprint arXiv:2307.02483</i> , 2023.
629 630 631 632	Zeming Wei, Yifei Wang, Ang Li, Yichuan Mo, and Yisen Wang. Jailbreak and guard aligned language models with only few in-context demonstrations, 2024. URL https://arxiv.org/abs/2310.06387.
633 634	Yotam Wolf, Noam Wies, Oshri Avnery, Yoav Levine, and Amnon Shashua. Fundamental limita- tions of alignment in large language models, 2023.
635 636 637 638	Youliang Yuan, Wenxiang Jiao, Wenxuan Wang, Jen-tse Huang, Pinjia He, Shuming Shi, and Zhaopeng Tu. Gpt-4 is too smart to be safe: Stealthy chat with llms via cipher. <i>arXiv preprint arXiv:2308.06463</i> , 2023.
639 640 641	Yi Zeng, Hongpeng Lin, Jingwen Zhang, Diyi Yang, Ruoxi Jia, and Weiyan Shi. How johnny can persuade llms to jailbreak them: Rethinking persuasion to challenge ai safety by humanizing llms. <i>arXiv preprint arXiv:2401.06373</i> , 2024.
642 643 644	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. Judging llm-as-a-judge with mt-bench and chatbot arena, 2023.
646 647	Chunting Zhou, Pengfei Liu, Puxin Xu, Srini Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, Susan Zhang, Gargi Ghosh, Mike Lewis, Luke Zettlemoyer, and Omer Levy. Lima: Less is more for alignment, 2023.

648 649	Daniel M. Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B. Brown, Alec Radford, Dario Amodei, Paul Christiano, and Geoffrey Irving. Fine-tuning language models from human preferences, 2020.	
651	Andy Zou, Zifan Wang, J. Zico Kolter, and Matt Fredrikson. Universal and transferable adversarial attacks on aligned language models, 2023.	
052		
000		
004		
000		
657		
658		
659		
660		
661		
662		
663		
664		
665		
666		
667		
668		
669		
670		
671		
672		
673		
674		
675		
676		
677		
678		
679		
680		
681		
682		
683		
684		
685		
686		
600		
680		
600		
691		
692		
693		
694		
695		
696		
697		
698		
699		
700		
701		