004

007

Hijacking Large Language Models via Adversarial In-Context Learning

Anonymous ACL submission

Abstract

In-context learning (ICL) has emerged as a powerful paradigm leveraging LLMs for specific downstream tasks by utilizing labeled examples as demonstrations in the precondition prompts. Despite its promising performance, ICL suffers from instability with the choice and arrangement of examples. Additionally, crafted adversarial attacks pose a notable threat to the robustness of ICL. However, existing attacks are either easy to detect, rely on external models, or lack specificity towards ICL. This work introduces a novel transferable attack for ICL to address these issues, aiming to hijack LLMs to generate the targeted response. The proposed hijacking attack leverages a gradient-based prompt search method to learn and append imperceptible adversarial suffixes to the in-context demonstrations. Extensive experimental results on various tasks and datasets demonstrate the effectiveness of our hijacking attack, resulting in distracted attention towards adversarial tokens and consequently leading to unwanted target outputs. We also propose a defense strategy against hijacking attacks through the use of extra demonstrations, which enhances the robustness of LLMs during ICL. Broadly, this work reveals the significant security vulnerabilities of LLMs and emphasizes the necessity for in-depth studies on the robustness of LLMs related to ICL.

1 Introduction

In-context learning (ICL) is an emerging technique for rapidly adapting large language models (LLMs), i.e., GPT-4 (Achiam et al., 2023) and LLaMA2 (Touvron et al., 2023), to new tasks without finetuning the pre-trained parameters (Brown et al., 2020). The key idea behind ICL is to provide LLMs with labeled examples as in-context demonstrations (demos) within the prompt context before a test query. Through learning from the demos, LLMs are able to generate responses to queries based on ICL (Dong et al., 2022; Min et al., 2022).

Several existing works, however, have demonstrated the highly unstable nature of ICL (Zhao et al., 2021; Chen et al., 2022). Specifically, performance on target tasks using ICL can vary wildly based on the selection and order of demos, giving rise to highly volatile outcomes ranging from random to near state-of-the-art (SOTA) (Qiang et al., 2020; Lu et al., 2021; Min et al., 2022; Pezeshkpour and Hruschka, 2023; Qiang et al., 2024). Correspondingly, several approaches (Liu et al., 2021; Wu et al., 2022; Nguyen and Wong, 2023) have been proposed to address the unstable issue of ICL.

043

044

045

046

047

050

051

052

053

054

056

057

059

060

061

062

063

064

065

067

068

069

071

072

073

074

075

076

077

078

079

081

Further research has looked at how adversarial examples can undermine the performance of ICL (Zhu et al., 2023a; Wang et al., 2023c,b; Shayegani et al., 2023). These studies show that maliciously designed examples injected into the prompt instructions (Zhu et al., 2023a; Zou et al., 2023; Xu et al., 2023), demos (Wang et al., 2023c; Mo et al., 2023a), or queries (Wang et al., 2023b; Kandpal et al., 2023) can successfully attack LLMs to degrade their performance, revealing the significant vulnerabilities of ICL against adversarial inputs.

While existing adversarial attacks have been applied to evaluate LLM robustness, they have some limitations in practice. Most character-level attacks, e.g., TextAttack (Morris et al., 2020) and TextBugger (Li et al., 2018), can be easily detected and evaded through grammar checks, limiting realworld effectiveness (Qiang et al., 2022; Jain et al., 2023). Some other attacks like BERTAttack (Li et al., 2020) require an extra model to generate adversarial examples, which may not be feasible in real-world applications. Crucially, existing attacks are not specifically crafted to target techniques based on LLMs, i.e., ICL. As such, the inherent security risks of LLMs remain largely unexplored. There is an urgent need for red teaming tailored to ICL to expose the substantial risk of LLMs for further evaluating their adversarial robustness against potential real-world threats.

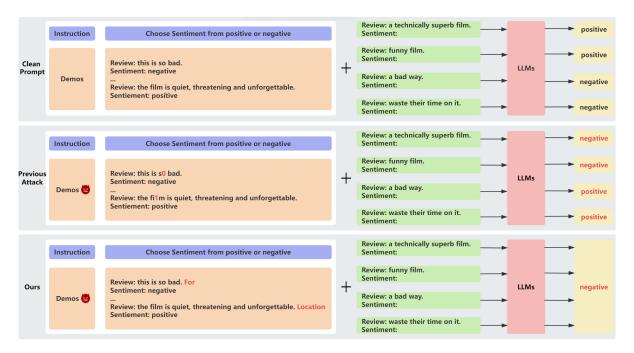


Figure 1: Illustrations of ICL using clean prompt and adversarial prompt. Given the clean in-context demos, LLMs can correctly generate the sentiment of the test queries. The previous attacks (Wang et al., 2023c) at the character level involve minor edits in some words, such as altering 'so' to 's0' and 'film' to 'film', of these in-context demos, leading to incorrect sentiment generated for the test queries. However, ours **learns** to **append** adversarial suffixes like 'For' and 'Location' to the in-context demos to efficiently and effectively **hijack** LLMs to generate the **unwanted target**, e.g., the 'negative' sentiment, **regardless** of the test query content.

This work proposes a novel adversarial attack specifically targeting ICL. We develop a gradientbased prompt search algorithm to learn adversarial suffixes in order to efficiently and effectively hijack LLMs via adversarial ICL. (Zhu et al., 2023a; Wang et al., 2023b) are the closest works to ours where they 'search' adversarial examples to simply manipulate model outputs. Yet, our attack method 'learns' adversarial tokens that directly hijack LLMs to generate the unwanted target that disrupts alignment with the desired output. In Figure 1, we illustrate the major difference between the previous attack and the proposed attack: ours hijacks LLMs to output the unwanted target response (e.g., 'negative') regardless of the query content. Furthermore, instead of manipulating the prompt instructions (Zhu et al., 2023a; Zou et al., 2023), demos (Wang et al., 2023c; Mo et al., 2023a), or queries (Wang et al., 2023b; Kandpal et al., 2023) leveraging standard adversarial examples, e.g., character-level attacks (Morris et al., 2020; Li et al., 2018), which are detectable easily, our hijacking attack is imperceptible in that it adds only 1-2 suffixes to the demos, as shown in Figure 1. Specifically, these suffixes are semantically incongruous but not easily identified as typos or gib-

086

087

100

101

102

104

105

106

108

109

berish compared to the existing ICL attack (Wang et al., 2023c). Finally, the backdoor attack during ICL (Kandpal et al., 2023) requires a trigger, which is impractical in real-world scenarios, whereas our attack hijacks the LLM to generate the unwanted target without triggering the queries.

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

Our extensive experiments validate the efficacy of our hijacking attacks. First, the attacks reliably induce LLMs to generate the targeted and misaligned output from the desired ones. Second, the adversarial suffixes learned via gradient optimization are transferable, remaining effective on different demo sets. Third, the suffix transferability holds even across different datasets for the same downstream task. Finally, our analysis shows that the adversarial suffixes distract LLMs' attention away from the task-relevant concepts. Our adversarial ICL attack poses a considerable threat to practical LLM applications due to its robust transferability and imperceptibility.

As this work represents one of the first efficient adversarial demonstration attacks during ICL, strategies for defending against such attacks have yet to be thoroughly investigated. Recently, (Mo et al., 2023b) introduced a method for defending against back-door attacks at test time, leveraging

few-shot demonstrations to correct the inference behavior of poisoned LLMs. Similarly, (Wei et al., 2023b) explored the power of in-context demos in manipulating the alignment ability of LLMs and proposed in-context attack and in-context defense methods for jailbreaking and guarding the aligned LLMs. Consequently, we explore the potential of using in-context demos exclusively to rectify the behavior of LLMs subjected to our hijacking attacks. Our defense strategy employs additional clean incontext demos at test time to safeguard LLMs from being hijacked by adversarial in-context demos. The experimental results demonstrate the efficacy of our proposed defense method against various adversarial demonstration attacks.

This work makes the following original contributions: (1) We propose a novel stealthy adversarial attack targeting in-context demos to hijack LLMs to generate unwanted target output during ICL. (2) We design a novel gradient-based prompt search algorithm to learn and append adversarial suffixes to the in-context demos. (3) Our extensive experiments demonstrate the effectiveness and transferability of the proposed adversarial ICL attack across various demo sets, LLMs, and tasks. (4) The proposed test time defense strategy effectively protects LLMs from being compromised by adversarial demonstration attacks.

2 Related Work

2.1 In-Context Learning

LLMs have shown impressive performance on numerous NLP tasks (Devlin et al., 2018; Lewis et al., 2019; Radford et al., 2019). Although fine-tuning has been a common method for adapting models to new tasks, it is often less feasible to fine-tune extremely large models with over 10 billion parameters. As an alternative, recent work has proposed ICL, where the model adapts to new tasks solely via inference conditioned on the provided in-context demos, without any gradient updates (Brown et al., 2020). By learning from the prompt context, ICL allows leveraging massive LLMs' knowledge without the costly fine-tuning process, showcasing an exemplar of the LLMs' emergent abilities (Schaeffer et al., 2023; Wei et al., 2022).

Intensive research has been dedicated to ICL. Initial works attempt to find better ways to select labeled examples for the demos (Liu et al., 2021; Rubin et al., 2021). For instance, (Liu et al., 2021) presents a simple yet effective retrieval-based method that selects the most semantically

similar examples as demos, leading to improved accuracy and higher stability. Follow-up works have been done to understand why ICL works (Xie et al., 2021; Razeghi et al., 2022; Min et al., 2022; Wei et al., 2023a; Kossen et al., 2023). (Xie et al., 2021) provides theoretical analysis that ICL can be formalized as Bayesian inference that uses the demos to recover latent concepts. Another line of research reveals the brittleness and instability of ICL approaches: small changes to the demo examples, labels, or order can significantly impact performance (Lu et al., 2021; Zhao et al., 2021; Min et al., 2022; Nguyen and Wong, 2023).

2.2 Adversarial Attacks on LLMs

Early adversarial attacks on LLMs apply simple character or token operations to trigger the LLMs to generate incorrect predictions, such as TextAttack (Morris et al., 2020) and BERT-Attack (Li et al., 2020). Since these attacks usually generate misspelled and/or gibberish prompts that can be detected using spell checker and perplexitybased filters, they are easy to block in real-world applications. Some other attacks struggled with optimizing over discrete text, leading to the manual or semi-automated discovery of vulnerabilities through trial-and-error (Li et al., 2021; Perez and Ribeiro, 2022; Li et al., 2023c; Qiang et al., 2023; Casper et al., 2023; Kang et al., 2023; Li et al., 2023a; Shen et al., 2023). For example, jailbreaking prompts are intentionally designed to bypass an LLM's built-in safeguard, eliciting it to generate harmful content that violates the usage policy set by the LLM vendor (Shen et al., 2023; Zhu et al., 2023b; Chao et al., 2023; Mehrotra et al., 2023; Jeong, 2023; Guo et al., 2024; Yu et al., 2024). These red teaming efforts craft malicious prompts in order to understand LLM's attack surface (Ganguli et al., 2022). However, the discrete nature of text has significantly impeded learning more effective adversarial attacks against LLMs.

Recent work has developed gradient-based optimizers for efficient text modality attacks. For example, (Wen et al., 2023) presented a gradient-based discrete optimizer that is suitable for attacking the text pipeline of CLIP, efficiently bypassing the safeguards in the commercial platform. (Zou et al., 2023), building on (Shin et al., 2020), described an optimizer that combines gradient guidance with random search to craft adversarial strings that induce LLMs to respond to the questions that would otherwise be banned. More recently, (Zhao et al.,

2024) proposed poisoning demonstration examples and prompts to make LLMs behave in alignment with pre-defined intentions.

238

239

240

241

243

247

248

249

250

251

253

256

265

268

269

270

272

273

275

276

277

281

286

288

Our hijacking attack algorithm falls into this stream of work, yet we target few-shot ICL instead of zero-shot queries. We use gradient-based prompt search to automatically learn effective adversarial suffixes rather than manually engineered prompts. Importantly, we show that LLMs can be hijacked to output the targeted unwanted output by appending optimized adversarial tokens to the ICL demos, which reveals a new lens of LLM vulnerabilities that may have been missed by prior approaches.

2.3 Defense Against Attacks on LLMs

The existing literature on the robustness of LLMs includes various strategies for defense (Goyal et al., 2023; Studnia et al., 2023; Liu et al., 2023; Xu et al., 2024; Wu et al., 2024). However, most of these defenses, such as those involving adversarial training (Liu et al., 2020; Li et al., 2023b; Formento et al., 2024; Wang et al., 2024) or data augmentation (Qiang et al., 2024; Yuan et al., 2024), need to re-train or fine-tune the models, which is computationally infeasible for LLM users. Moreover, the restriction of many closed-source LLMs to only permit query access for candidate defenses introduces new challenges.

Recent studies focus on developing defenses against attacks on LLMs that utilize adversarial prompting. (Jain et al., 2023) and (Alon and Kamfonas, 2023) have suggested the use of perplexity filters to detect adversarial prompts. While the filters are effective at catching the attack strings that contain gibberish words or character-level adversarial tokens with high perplexity scores, they fall short in detecting more subtle adversarial prompts, like the ones used in our adversarial demonstration attacks with as low perplexity as clean samples shown in Figure 2. Recently, (Mo et al., 2023b) introduced a method to mitigate backdoor attacks at test time by identifying the task and retrieving relevant defensive demonstrations. These demonstrations are combined with user queries to counteract the adverse effects of triggers present in backdoor attacks. This defense strategy eliminates the need for modifications or tuning of LLMs. Its objective is to re-calibrate and correct the behavior of LLMs during test-time evaluations. Similarly, (Wei et al., 2023b) investigated the role of in-context demonstrations in enhancing the robustness of LLMs and highlighted their effectiveness in defending against

jailbreaking attacks. The authors developed an in-context defense strategy that constructs a safe context to caution the model against generating any harmful content.

289

290

291

292

293

294

295

296

297

298

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

327

328

329

330

331

332

333

334

335

336

So far, defense mechanisms against adversarial demonstration attacks have not been extensively explored. Our approach introduces a test-time defense strategy that uses additional clean in-context demos to safeguard LLMs from adversarial incontext manipulations. In line with prior works (Mo et al., 2023b; Wei et al., 2023b; Wang et al., 2024), this defense strategy avoids the necessity for retraining or fine-tuning LLMs. Instead, it focuses on re-calibrating and correcting the behavior of LLMs during evaluations at test time.

3 Preliminaries

3.1 ICL Formulation

Formally, ICL is characterized as a problem involving the conditional generation of text (Liu et al., 2021), where an LLM \mathcal{M} is employed to generate response y_Q given an optimal task instruction I, a demo set C, and an input query x_Q . I specifies the downstream task that ${\mathcal M}$ should perform, e.g., "Choose sentiment from positive or negative" used in our sentiment generation task. C consists of N (e.g., 8) concatenated data-label pairs following a specific template S, formally: $C = [S(x_1, y_1); \dots; S(x_N, y_N)], ';'$ here denotes the concatenation operator. Thus, given the input prompt as $p = [I; C; S(x_Q, _)], \mathcal{M}$ generates the response as $\hat{y}_Q = \mathcal{M}(p)$. $S(x_Q, \underline{\ })$ here means using the same template as the demos but with the label empty.

3.2 Adversarial Attack on LLMs

In typical text-based adversarial attacks, the attackers manipulate the input x with the goal of misleading the model to produce inaccurate output or bypass safety guardrails (Zou et al., 2023; Maus et al., 2023). Specifically, given the input-output pair (x,y), the attackers aim to learn the adversarial perturbation δ adding to x by maximizing the model's objective function but without misleading humans by bounding the perturbation within the "perceptual" region Δ . The objective function of the attacking process thus can be formulated as:

$$\max_{\delta \in \Delta} \mathcal{L}(\mathcal{M}(x_Q + \delta), y_Q), \tag{1}$$

 \mathcal{L} here denotes the task-specific loss function, for instance, cross-entropy loss for classification tasks.

4 The Threat Model

4.1 ICL Hijacking Attack

ICL consists of an instruction I, a demo set C, and an input query x_Q , providing more potential attack vectors compared to the conventional text-based adversarial attacks. This work focuses on manipulating C without changing I and x_Q . Recently, (Wang et al., 2023c) applied character-level perturbation techniques (e.g., TextAttack (Morris et al., 2020) and TextBugger (Li et al., 2018)), such as character insertion, character deletion, neighboring character swap, and character substitution, to reverse the output.

Our hijacking attack learns the adversarial suffix tokens to the in-context demos to manipulate LLMs' output via a new greedy gradient-based prompt injection algorithm. Given a clean demo set $C = [S(x_1, y_1); \cdots; S(x_N, y_N)]$, our hijacking attack automatically produces an adversarial suffix for each demo in c, formally:

$$C' = [S(x_1 + \delta_1, y_1); \dots; S(x_N + \delta_N, y_N)], (2)$$

where C' denotes the perturbed demo set. To make it clear, the adversarial suffixes appended to each demo as perturbations are different, respectively. In this case, the attack or perturbation budget refers to the number of tokens in each adversarial suffix.

As a result, our hijacking attack induces \mathcal{M} to generate an unwanted target output y_T via appending adversarial suffix tokens on the in-context demos as $y_T = \mathcal{M}(p')$. In other words, \mathcal{M} generates the same or different responses for the clean and perturbed prompts depending on the True or False of $\mathcal{M}(p) = y_T$:

$$\begin{cases} \mathcal{M}(p) = \mathcal{M}(p'), & \text{True}, \\ \mathcal{M}(p) \neq \mathcal{M}(p'), & \text{False}, \end{cases}$$

where $p = [I; C; S(x_Q, _)]$ and $p' = [I; C'; S(x_Q, _)]$, respectively.

4.2 Hijacking Attack Objective

We express the goal of the hijacking attack as a formal objective function. Let us consider the LLM \mathcal{M} as a function that maps a sequence of tokens $x_{1:n}$, with $x \in \{1, \cdots, V\}$ where V denote the vocabulary size, namely, the number of tokens, to a probability distribution over the next token x_{n+1} . Specifically, $\mathcal{P}(x_{n+1}|x_{1:n})$ denotes the probability that x_{n+1} is the next token given the previous tokens $x_{1:n}$. In more detail, we formulate $x_{1:n}$ as

 $[I; C; S(x_Q, _)]$ and x_{n+1} as \hat{y}_Q in the context of an ICL task, respectively.

Using the notations defined earlier, the hijacking attack objective we want to optimize is simply the negative log probability of the target token x_{n+1} . The generated target output y_T is different from the ground truth label y_Q for the training query (x_Q, y_Q) . Formally:

$$\mathcal{L}(x_Q) = -\log \mathcal{P}(\mathcal{M}(y_T|p')), \tag{3}$$

where $y_T \neq y_Q$, demonstrating the attack hijacks \mathcal{M} to generate the target output. p' denotes the perturbed prompt as: $p' = [I; C'; S(x_Q, _)]$. In summary, the problem of optimizing the adversarial suffix tokens can be formulated as the following optimization objective:

$$\underset{\delta_i \in \{1, \dots, V\}^{|N|}}{\text{minimize}} \mathcal{L}(x_Q), \tag{4}$$

where i denotes the indices of the demos and N is the number of demos in the perturbed demos set C', respectively.

4.3 Greedy Gradient-guided Injection

A primary challenge in optimizing Eq. 4 is optimizing over a discrete set of possible token values. While there are some methods for discrete optimization, prior work (Carlini et al., 2023) has shown that those effective strategies often struggle to reliably attack the aligned LLMs.

Motivated by prior works (Shin et al., 2020; Zou et al., 2023; Wen et al., 2024), we propose a simple yet effective algorithm for LLMs hijacking attacks, called greedy gradient-guided injection (GGI) algorithm (Algorithm 1 in Appendix). The key idea comes from greedy coordinate descent: if we could evaluate all possible suffix token injections, we could substitute the tokens that maximize the adversarial loss reduction. Since exhaustively evaluating all tokens is infeasible due to the large candidate vocabulary size, we instead leverage gradients with respect to the suffix indicators to find promising candidate tokens for each position. We then evaluate all of these candidate injections with explicit forward passes to find the one that decreases the loss the most. This allows an efficient approximation of the true greedy selection. We can optimize the discrete adversarial suffixes by iteratively injecting the best tokens.

We compute the linearized approximation of replacing the demo x_i in C by evaluating the gradient $\nabla_{\mathbf{e}_{x_i^j}} \mathcal{L}(x_Q) \in \mathbb{R}^{|V|}$, where $\mathbf{e}_{x_i^j}$ denotes the

vector representing the current value of the j-th adversarial suffix token. Note that because LLMs typically form embeddings for each token, they can be written as functions of $\mathbf{e}_{x_i^j}$, and thus we can immediately take the gradient with respect to this quantity (Ebrahimi et al., 2017; Shin et al., 2020).

The key aspects of our GGI algorithm (Algorithm 1 in Appendix) are: firstly, it uses gradients of the selected token candidates to calculate the top candidates; secondly, it evaluates the top candidates explicitly to identify the most suitable one; and lastly, it iteratively injects the best token at each position to optimize the suffixes. This approximates an extensive greedy search in a computationally efficient manner.

5 The Defense Method

Having developed the ICL hijacking attack by incorporating adversarial tokens into the in-context demos, we now present a straightforward yet potent defense strategy to counter this attack. Initially, we assume that defenders treat LLMs as black-box, lacking any insight into their training processes or underlying parameters. The defenders apply defense on the input prompt p directly during test-time evaluation. Their goal is to rectify the behavior of LLMs and induce LLMs to generate desired responses to user queries.

Given an input prompt p' that includes adversarial tokens within the demos C', we assume that LLMs, when presented with demos containing clean data for the same tasks, can understand the genuine intent of the user's query through ICL, rather than being misled by the adversarial demos. In this context, 'clean data' refers to data without any adversarial tokens and is randomly selected from the training set. More precisely, the defenders modify the input prompt p' into \tilde{p} by appending or inserting more clean demos into the demo set C', as follows: $\tilde{p} = [I; C'; \tilde{C}; S(x_Q, _)]$. $\tilde{C} = [S(\tilde{x}_1, \tilde{y}_1); \cdots; S(\tilde{x}_N, \tilde{y}_N)]$ here denotes the clean demos selected from the training set. Through this approach, the defender guarantees that the in-context demos align with the user's query and possess resilience against adversarial attacks. In our experiments, we maintained an equal number of demos in C' and \tilde{C} and observed that this method resulted in effective defense across a range of datasets and tasks.

6 Experiment Setup

Datasets: We evaluate the performance of our LLM hijacking algorithm and other baseline al-

gorithms on three classification datasets covering sentiment analysis and topic generation tasks. SST-2 (Socher et al., 2013) and Rotten Tomatoes (RT) (Pang and Lee, 2005) are both binary sentiment analysis datasets of movie reviews. AG's News (Zhang et al., 2015) is a multi-class news topic classification dataset. These datasets allow us to evaluate the hijacking attacks on a diverse set of text classification benchmarks across both binary and multi-class settings. More details on the dataset statistics are provided in Table 4 of the Appendix. Large Language Models: The experiments are conducted using three different LLMs, GPT2-XL (Radford et al., 2019), LLaMA-7b/13b (Touvron et al., 2023), OPT-2.7b/6.7b (Zhang et al., 2022), and Vicuna-7b (Chiang et al., 2023) allowing us to evaluate attack effectiveness on both established and SOTA LLMs. The choice of LLMs covers a diverse set of architectures and model sizes, enabling a comprehensive evaluation of the vulnerability of LLMs via adversarial ICL.

ICL Settings: For ICL, we follow the setting in (Wang et al., 2023c) and use their template to incorporate the demos for prediction. The detailed template is provided in Figure 8 of the Appendix. We evaluate the 2-shot, 4-shot, and 8-shot settings for the number of demos. Specifically, for each test example, we randomly select the demos from the training set and repeat this process 5 times, reporting the average accuracy over the repetitions.

Evaluation Metrics: Several different metrics evaluate the performance of ICL and hijacking attacks. Clean accuracy evaluates the accuracy of ICL on downstream tasks using clean demos. Attack accuracy evaluates the accuracy of ICL given the perturbed demons. Defense accuracy demonstrates the accuracy of ICL with the defense method against the hijacking attack. We further evaluate the effectiveness of hijacking attacks using attack success rate (ASR). Given a test sample (x, y) from a test set D, the clean and perturbed prompts are denoted as p = [I; C; x] and p' = [I; C'; x], respectively. ASR is calculated as $ASR = \sum_{(x,y) \in D} \frac{\mathbb{I}(\mathcal{M}(p') = y_T)}{\mathbb{I}(\mathcal{M}(p) = y)}$,

where 1 denotes the indicator function.

7 Result and Discussion

7.1 ICL Performance

The rows identified as 'Clean' in Table 1 and Table 2 show the ICL performance on the respective tasks when using clean in-context demos. In particular, Table 1 presents the accuracies for the gen-

Table 1: The performance on sentiment analysis task with and without attacks on ICL. The row identified as 'Clean' in gray color represents the accuracy with clean in-context demos. Other rows illustrate the accuracies with adversarial in-context demos. The details of the baselines in green color are present in Section B of the Appendix. Specifically, we employ TextAttack (TA) (Morris et al., 2020) following the attack in (Wang et al., 2023c) as the most closely related baseline for our attack (GGI). The accuracies of positive (P) and negative (N) sentiments are reported separately to highlight the effectiveness of our hijacking attack.

				SS'	T-2					R	Т		
Model	Method	2-sl	nots	4-sl	nots	8-sl	nots	2-sl	hots	4-sl	nots	8-sl	nots
		P	N	P	N	P	N	P	N	P	N	P	N
	Clean	94.7	52.2	88.6	49.4	91.6	69.0	93.3	54.7	88.6	76.9	90.2	80.5
	Square	99.4	2.0	99.8	4.2	99.4	11.0	99.8	1.5	100	4.1	99.3	7.5
GPT2-XL	Greedy	100	10.8	100	6.2	100	0.2	100	5.3	100	2.8	100	0.0
	TA	95.0	2.2	99.8	17.8	99.6	21.6	95.9	8.1	96.3	41.3	96.4	47.3
	GGI	100	1.2	100	0.0	100	0.0	100	2.8	100	0.0	100	0.0
	Clean	69.4	87.8	70.2	93.8	77.8	93.0	84.4	91.4	84.4	93.1	88.6	92.8
	Square	99.2	31.4	93.8	72.2	99.6	29.0	98.1	42.2	97.0	68.7	99.4	33.2
OPT-6.7b	Greedy	100	25.0	97.8	39.0	100	2.0	99.4	31.7	99.8	4.7	100	0.8
	TA	94.8	80.8	54.8	98.6	91.6	89.4	92.5	86.1	77.6	96.4	94.0	86.3
	GGI	100	0.0	98.4	2.0	100	0.2	100	2.6	99.8	0.0	100	0.2
	Clean	91.4	81.2	88.2	81.4	94.6	82.6	84.8	78.4	85.9	80.5	90.4	85.4
	Square	89.2	84.4	86.6	85.8	94.0	83.8	85.9	85.4	84.6	88.6	91.6	88.4
Vicuna-7b	Greedy	93.0	83.4	88.4	87.0	94.6	80.0	91.2	82.8	86.9	88.7	91.9	85.9
	TA	87.0	85.2	76.2	88.2	94.2	80.6	83.3	84.2	79.6	88.6	92.1	84.4
	GGI	90.6	42.2	96.4	23.2	100	0.8	87.6	36.4	95.1	35.7	100	0.2
	Clean	81.4	86.3	74.4	91.9	82.7	92.4	86.0	83.6	81.9	91.6	89.3	97.8
	Square	86.8	80.0	96.8	58.6	98.0	56.4	86.9	57.4	97.4	50.1	97.8	57.4
LLaMA-7b	Greedy	95.0	47.6	100	0.0	100	0.0	88.9	2.8	99.8	0.0	100	0.0
	TA	87.2	77.8	93.8	69.0	99.8	8.8	83.1	57.4	94.2	68.9	99.6	3.80
	GGI	100	0.4	100	0.0	100	0.0	96.8	0.0	100	0.0	100	0.0
	Clean	97.8	76.4	95.6	88.0	95.8	90.0	94.2	84.8	92.7	92.1	91.4	91.9
	Square	98.4	72.8	98.2	78.4	97.8	85.4	93.6	87.4	94.4	84.1	94.2	87.6
LLaMA-13b	Greedy	98.0	41.4	100	3.0	100	0.0	55.9	11.3	92.9	0.0	100	0.4
	TA	98.2	72.2	92.8	92.8	97.5	87.6	94.8	81.8	88.0	94.0	92.5	89.3
	GGI	99.2	37.8	100	7.2	100	0.0	99.1	3.8	86.1	3.6	100	0.0

Table 2: The performance of AG's News topic generation task with and without attacks on ICL. The clean and attack accuracies are reported separately for the four topics. These results highlight the effectiveness of our hijacking attacks to induce LLMs to generate the target token, i.e., "tech", regardless of the query content.

Model	Method		4-8	shots	8-shots				
Model	Method	word	sports	business	tech	word	sports	business	tech
	Clean	48.5	87.0	64.9	71.9	48.2	50.6	71.0	83.6
	Square	2.0	66.0	26.8	96.0	19.6	65.6	28.0	97.2
GPT2-XL	Greedy	12.8	60.4	29.2	96.4	8.0	21.2	10.0	98.8
	TA	54.8	84.0	73.2	82.4	82.0	82.4	91.2	57.6
	GGI	0.0	2.0	0.4	100	0.0	0.0	0.0	100
	Clean	68.2	96.8	66.6	49.0	88.6	97.4	78.2	61.0
	Square	78.4	98.0	76.0	36.8	94.4	98.0	60.0	57.6
LLaMA-7b	Greedy	69.6	98.8	75.2	51.6	89.6	100	68.4	73.6
	TA	42.4	94.8	67.6	32.4	95.2	96.0	39.2	24.8
	GGI	0.0	20.0	0.00	98.0	29.6	56.0	0.0	100

eration of positive (P) and negative (N) sentiments in the SST-2 and RT datasets. All the tested LLMs perform well, achieving an average accuracy of 83.6% on SST-2 and 86.7% on RT across various in-context few-shot settings. Table 2 indicates that LLMs with ICL also perform well in the context of multi-class generation on AG's News dataset. The average accuracies stand at 69.1% for 4-shot settings and 72.3% for 8-shot settings across various LLMs. Additionally, LLMs with ICL exhibit improved performance with an increased number of in-context demos, particularly achieving best results with 8-shot settings.

532

533

534

536

537

538

541

543

544

545

547

7.2 ICL Performance with Hijacking Attack While LLMs utilizing ICL show strong performance with clean in-context demos, Tables 1 and 2

reveal that their effectiveness is significantly undermined by hijacking attacks. These attacks manipulate the models to produce target outputs by appending adversarial suffixes to the in-context demos. The baseline attacks successfully induce LLMs to generate the targeted positive sentiment through a few shots of adversarially perturbed demos, as shown in Table 1. As a result, the positive test samples achieve a predominantly higher accuracy than the negative ones. Furthermore, our GGI attack more effectively hijacks LLMs to generate the target output, i.e., exclusively positive sentiment token, regardless of the test query content. The positive test samples achieve almost 100% accuracy. On the contrary, the negative ones get nearly 0% accuracy on most settings, as shown in 1. For the

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

Table 3: The performance of the defenses using ASRs across various LLMs and datasets. Adv denotes our hijacking attack using the adversarial demos. Adv+Clean represents the proposed defense method, leveraging extra clean demos with adversarial demos. The numbers within the parenthesis indicate the reduction in the ASRs after defense.

SST-2					RT		AG's News			
Model	odel Adv+Clean			Adv	Adv+	Clean	Adv	Adv+Clean		
	Adv	Preceding Proceeding		Auv	Preceding	Proceeding	Auv	Preceding	Proceeding	
GPT2-XL	100	100 (-0)	99.6 (-0.4)	100	100 (-0)	97.4 (-2.6)	99.1	75.5 (-23.6)	80.5 (-18.6)	
OPT-6.7b	98.2	44.9 (-53.3)	52.5 (-45.7)	99.9	50.2 (-49.7)	57.8 (-42.1)	65.6	23.5 (-42.1)	22.5 (-43.1)	
LLaMA-7b	100	49.1 (-50.9)	98.3 (-1.7)	100	53.1 (-46.9)	99.8 (-0.2)	82.8	42.2 (-40.6)	88.2 (+5.4)	

more complex multi-class AG's News topic generation task, the effectiveness of those baseline attacks decreases significantly. Especially the TA method is not successful in hijacking LLMs to generate the target token, i.e., 'tech'. Only our GGI attack successfully hijacks the LLMs to generate "tech" by appending the adversarial suffixes to the in-context demos, as shown in Table 2.

7.3 Defense Method Performance

564

565

566

569

570

571

573

574

575

578

581

586

590

591

592

594

595

596

597

604

Table 3 presents ASRs of our hijacking attack when countered with the proposed defense mechanism that uses additional clean demos. The defense method is tested in two different settings. In the Preceding setting, clean demos are placed before the adversarial demos in the sequence $\tilde{p} = [I; \tilde{C}; C'; S(x_Q, _)]$. Conversely, in the Proceeding setting, clean demos are added after the adversarial demos, forming the sequence $\tilde{p} = [I; C'; \tilde{C}; S(x_Q, _)]$.

The results show a significant decrease in ASRs of our hijacking attack, affirming the effectiveness of the defense method. Notably, the Preceding setting results in considerably lower ASRs compared to the Proceeding setting. This relates to the mechanism through which our hijacking attack induces LLMs to generate target outputs. As depicted in Figure 6b, the adversarial suffixes divert the LLMs' attention away from the original query. Furthermore, Figure 7 of the Appendix illustrates that the LLM primarily focuses on the initial segments of the demos, which are indicated by a darker green color. Therefore, in the Preceding method, the model shifts its attention to these first few demos. which contain additional clean samples before the adversarial demos. These clean samples effectively re-calibrate and rectify the model's behavior, resulting in a larger reduction in ASRs, as shown in Table 3. In contrast, the first few demos remain adversarial in the Proceeding method, rendering it ineffective in defending against the adversarial demonstration attack.

Furthermore, the results indicate that our proposed defense methods are ineffective on small-

sized LLMs, such as the GPT2-XL used in our experiments. We hypothesize that this is due to their limited emergent abilities. In other words, employing additional demos during ICL cannot correct the behavior of small-sized LLMs under hijacking attacks.

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

7.4 Stealthiness of GGI

The perplexity scores shown in Figure 2 for both the baseline and our attacks exhibit minor (non-significant, P-value of $0.329 \gg 0.05$ cutoff is reported from one-way ANOVA test) increases compared to the score of the clean samples, highlighting the stealth of our proposed word-level adversarial attacks. Specifically, the adversarial triggers learned from our GGI algorithm are imperceptible and maintain the semantic integrity and coherence of the original content, as shown in the examples of Figures 9 and 10 in the Appendix.

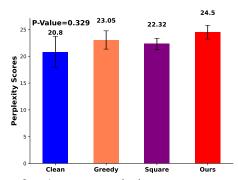


Figure 2: Average perplexity scores reported for LLaMA2-7b on 100 random samples under eight-shots setting from SST-2 derived from three separate runs under various attacks.

8 Conclusion

This work further reveals the vulnerability of ICL via crafted hijacking attacks. By appending imperceptible adversarial suffixes to the in-context demos using a greedy gradient-based algorithm, our attack effectively hijacks the LLMs to generate the unwanted targets by diverting LLMs' attention from the relevant context to the adversarial suffixes. We will continue on studying novel attack and defense techniques to enhance ICL security.

9 Limitations and Risks

635

637

638

639

641

643

644

647 648

650

652

655 656

658

662

This work uncovers a potential vulnerability of LLMs during in-context learning. By inserting adversarial tokens, which our algorithm has learned, into in-context demos, we can make the LLM produce undesired target outputs without the need for a trigger in the query. Our evaluation focuses on the attack success rate, particularly in the context of sentiment analysis and topic generation tasks. Our threat model, which is based on single token generation, has proven to be highly effective while maintaining content integrity. This efficiency eliminates the necessity for models that generate multiple tokens, which could compromise the content integrity. However, it is possible that our attack may be more effective for generation tasks across the LLMs that are similar in sizes (or smaller) and training approaches. Further studies are warranted to extend our approach to a wide range of downstream tasks and LLMs.

This work represents a purple teaming effort with the goal to discover the vulnerabilities of LLM during in-context learning and defend against the attacks. It offers a unified platform that enables both the red team and blue team to collaborate more effectively. Moreover, it facilitates a seamless knowledge transfer between the teams. As such, it will not pose risks for natural users nor LLM vendors. Rather, our findings can be utilized by these stakeholders to guard against malicious uses and enhance the robustness of LLMs to such threats.

667	Reference
-----	-----------

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. arXiv preprint arXiv:2303.08774.
- Gabriel Alon and Michael Kamfonas. 2023. Detecting language model attacks with perplexity. *arXiv* preprint arXiv:2308.14132.
- Maksym Andriushchenko, Francesco Croce, Nicolas Flammarion, and Matthias Hein. 2020. Square attack: a query-efficient black-box adversarial attack via random search. In *European conference on computer vision*, pages 484–501. Springer.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Nicholas Carlini, Milad Nasr, Christopher A Choquette-Choo, Matthew Jagielski, Irena Gao, Anas Awadalla, Pang Wei Koh, Daphne Ippolito, Katherine Lee, Florian Tramer, et al. 2023. Are aligned neural networks adversarially aligned? *arXiv preprint arXiv:2306.15447*.
- Stephen Casper, Jason Lin, Joe Kwon, Gatlen Culp, and Dylan Hadfield-Menell. 2023. Explore, establish, exploit: Red teaming language models from scratch. *arXiv preprint arXiv:2306.09442*.
- Patrick Chao, Alexander Robey, Edgar Dobriban, Hamed Hassani, George J Pappas, and Eric Wong. 2023. Jailbreaking black box large language models in twenty queries. *arXiv preprint arXiv:2310.08419*.
- Yanda Chen, Chen Zhao, Zhou Yu, Kathleen McKeown, and He He. 2022. On the relation between sensitivity and accuracy in in-context learning. *arXiv* preprint *arXiv*:2209.07661.
- Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E Gonzalez, et al. 2023. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality. *See https://vicuna.lmsys.org (accessed 14 April 2023)*, 2(3):6.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Zhiyong Wu, Baobao Chang, Xu Sun, Jingjing Xu, and Zhifang Sui. 2022. A survey for in-context learning. *arXiv preprint arXiv:2301.00234*.
- Javid Ebrahimi, Anyi Rao, Daniel Lowd, and Dejing Dou. 2017. Hotflip: White-box adversarial

examples for text classification. arXiv preprint arXiv:1712.06751.

- Brian Formento, Wenjie Feng, Chuan Sheng Foo, Luu Anh Tuan, and See-Kiong Ng. 2024. Semrode: Macro adversarial training to learn representations that are robust to word-level attacks. *arXiv preprint arXiv:2403.18423*.
- Deep Ganguli, Liane Lovitt, Jackson Kernion, Amanda Askell, Yuntao Bai, Saurav Kadavath, Ben Mann, Ethan Perez, Nicholas Schiefer, Kamal Ndousse, et al. 2022. Red teaming language models to reduce harms: Methods, scaling behaviors, and lessons learned. arXiv preprint arXiv:2209.07858.
- Shreya Goyal, Sumanth Doddapaneni, Mitesh M Khapra, and Balaraman Ravindran. 2023. A survey of adversarial defenses and robustness in nlp. *ACM Computing Surveys*, 55(14s):1–39.
- Xingang Guo, Fangxu Yu, Huan Zhang, Lianhui Qin, and Bin Hu. 2024. Cold-attack: Jailbreaking llms with stealthiness and controllability. *arXiv preprint arXiv:2402.08679*.
- Neel Jain, Avi Schwarzschild, Yuxin Wen, Gowthami Somepalli, John Kirchenbauer, Ping-yeh Chiang, Micah Goldblum, Aniruddha Saha, Jonas Geiping, and Tom Goldstein. 2023. Baseline defenses for adversarial attacks against aligned language models. arXiv preprint arXiv:2309.00614.
- Joonhyun Jeong. 2023. Hijacking context in large multimodal models. *arXiv preprint arXiv:2312.07553*.
- Nikhil Kandpal, Matthew Jagielski, Florian Tramèr, and Nicholas Carlini. 2023. Backdoor attacks for in-context learning with language models. *arXiv* preprint arXiv:2307.14692.
- Daniel Kang, Xuechen Li, Ion Stoica, Carlos Guestrin, Matei Zaharia, and Tatsunori Hashimoto. 2023. Exploiting programmatic behavior of llms: Dual-use through standard security attacks. *arXiv preprint arXiv:2302.05733*.
- Jannik Kossen, Yarin Gal, and Tom Rainforth. 2023. Incontext learning learns label relationships but is not conventional learning. In *The Twelfth International Conference on Learning Representations*.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Ves Stoyanov, and Luke Zettlemoyer. 2019. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. arXiv preprint arXiv:1910.13461.
- Haoran Li, Dadi Guo, Wei Fan, Mingshi Xu, and Yangqiu Song. 2023a. Multi-step jailbreaking privacy attacks on chatgpt. *arXiv preprint arXiv:2304.05197*.

Jiazhao Li, Zhuofeng Wu, Wei Ping, Chaowei Xiao, and VG Vydiswaran. 2023b. Defending against insertion-based textual backdoor attacks via attribution. *arXiv* preprint arXiv:2305.02394.

- Jinfeng Li, Shouling Ji, Tianyu Du, Bo Li, and Ting Wang. 2018. Textbugger: Generating adversarial text against real-world applications. *arXiv preprint arXiv:1812.05271*.
- Linyang Li, Ruotian Ma, Qipeng Guo, Xiangyang Xue, and Xipeng Qiu. 2020. Bert-attack: Adversarial attack against bert using bert. *arXiv* preprint *arXiv*:2004.09984.
- Xin Li, Xiangrui Li, Deng Pan, Yao Qiang, and Dongxiao Zhu. 2023c. Learning compact features via intraining representation alignment. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pages 8675–8683.
- Xin Li, Xiangrui Li, Deng Pan, and Dongxiao Zhu. 2021. Improving adversarial robustness via probabilistically compact loss with logit constraints. In *Proceedings of the AAAI conference on artificial intelligence*, volume 35, pages 8482–8490.
- Jiachang Liu, Dinghan Shen, Yizhe Zhang, Bill Dolan, Lawrence Carin, and Weizhu Chen. 2021. What makes good in-context examples for gpt-3? *arXiv* preprint arXiv:2101.06804.
- Qin Liu, Fei Wang, Chaowei Xiao, and Muhao Chen. 2023. From shortcuts to triggers: Backdoor defense with denoised poe. *arXiv preprint arXiv:2305.14910*.
- Xiaodong Liu, Hao Cheng, Pengcheng He, Weizhu Chen, Yu Wang, Hoifung Poon, and Jianfeng Gao. 2020. Adversarial training for large neural language models. *arXiv preprint arXiv:2004.08994*.
- Yao Lu, Max Bartolo, Alastair Moore, Sebastian Riedel, and Pontus Stenetorp. 2021. Fantastically ordered prompts and where to find them: Overcoming few-shot prompt order sensitivity. *arXiv* preprint *arXiv*:2104.08786.
- Natalie Maus, Patrick Chao, Eric Wong, and Jacob R Gardner. 2023. Black box adversarial prompting for foundation models. In *The Second Workshop on New Frontiers in Adversarial Machine Learning*.
- Anay Mehrotra, Manolis Zampetakis, Paul Kassianik, Blaine Nelson, Hyrum Anderson, Yaron Singer, and Amin Karbasi. 2023. Tree of attacks: Jailbreaking black-box llms automatically. *arXiv preprint arXiv:2312.02119*.
- Sewon Min, Xinxi Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2022. Rethinking the role of demonstrations: What makes in-context learning work? *arXiv* preprint arXiv:2202.12837.

Lingbo Mo, Boshi Wang, Muhao Chen, and Huan Sun. 2023a. How trustworthy are open-source llms? an assessment under malicious demonstrations shows their vulnerabilities. *arXiv preprint arXiv:2311.09447*.

- Wenjie Mo, Jiashu Xu, Qin Liu, Jiongxiao Wang, Jun Yan, Chaowei Xiao, and Muhao Chen. 2023b. Test-time backdoor mitigation for black-box large language models with defensive demonstrations. *arXiv* preprint arXiv:2311.09763.
- John X Morris, Eli Lifland, Jin Yong Yoo, Jake Grigsby, Di Jin, and Yanjun Qi. 2020. Textattack: A framework for adversarial attacks, data augmentation, and adversarial training in nlp. *arXiv preprint arXiv:2005.05909*.
- Tai Nguyen and Eric Wong. 2023. In-context example selection with influences. *arXiv preprint arXiv:2302.11042*.
- Bo Pang and Lillian Lee. 2005. Seeing stars: Exploiting class relationships for sentiment categorization with respect to rating scales. In *Proceedings of the ACL*.
- Fábio Perez and Ian Ribeiro. 2022. Ignore previous prompt: Attack techniques for language models. *arXiv preprint arXiv:2211.09527*.
- Pouya Pezeshkpour and Estevam Hruschka. 2023. Large language models sensitivity to the order of options in multiple-choice questions. *arXiv* preprint *arXiv*:2308.11483.
- Yao Qiang, Supriya Tumkur Suresh Kumar, Marco Brocanelli, and Dongxiao Zhu. 2022. Tiny rnn model with certified robustness for text classification. In 2022 International Joint Conference on Neural Networks (IJCNN), pages 1–8. IEEE.
- Yao Qiang, Chengyin Li, Prashant Khanduri, and Dongxiao Zhu. 2023. Interpretability-aware vision transformer. *arXiv preprint arXiv:2309.08035*.
- Yao Qiang, Xin Li, and Dongxiao Zhu. 2020. Toward tag-free aspect based sentiment analysis: A multiple attention network approach. In 2020 International Joint Conference on Neural Networks (IJCNN), pages 1–8. IEEE.
- Yao Qiang, Subhrangshu Nandi, Ninareh Mehrabi, Greg Ver Steeg, Anoop Kumar, Anna Rumshisky, and Aram Galstyan. 2024. Prompt perturbation consistency learning for robust language models. *arXiv* preprint arXiv:2402.15833.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.
- Yasaman Razeghi, Robert L Logan IV, Matt Gardner, and Sameer Singh. 2022. Impact of pretraining term frequencies on few-shot reasoning. *arXiv* preprint *arXiv*:2202.07206.

Ohad Rubin, Jonathan Herzig, and Jonathan Berant. 2021. Learning to retrieve prompts for in-context learning. *arXiv preprint arXiv:2112.08633*.

- Rylan Schaeffer, Brando Miranda, and Sanmi Koyejo. 2023. Are emergent abilities of large language models a mirage? *arXiv preprint arXiv:2304.15004*.
- Erfan Shayegani, Md Abdullah Al Mamun, Yu Fu, Pedram Zaree, Yue Dong, and Nael Abu-Ghazaleh. 2023. Survey of vulnerabilities in large language models revealed by adversarial attacks. *arXiv* preprint arXiv:2310.10844.
- Xinyue Shen, Zeyuan Chen, Michael Backes, Yun Shen, and Yang Zhang. 2023. "do anything now": Characterizing and evaluating in-the-wild jailbreak prompts on large language models. *arXiv preprint arXiv:2308.03825*.
- Taylor Shin, Yasaman Razeghi, Robert L Logan IV, Eric Wallace, and Sameer Singh. 2020. Autoprompt: Eliciting knowledge from language models with automatically generated prompts. *arXiv preprint arXiv:2010.15980*.
- Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D Manning, Andrew Y Ng, and Christopher Potts. 2013. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of the 2013 conference on empirical methods in natural language processing*, pages 1631–1642.
- Joachim Studnia, Simiao Zuo, Xiaodong Liu, Qiang Lou, Jian Jiao, and Denis Charles. 2023. Evaluating adversarial defense in the era of large language models. In R0-FoMo: Robustness of Few-shot and Zero-shot Learning in Large Foundation Models.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.
- Haoyu Wang, Guozheng Ma, Cong Yu, Ning Gui, Linrui Zhang, Zhiqi Huang, Suwei Ma, Yongzhe Chang, Sen Zhang, Li Shen, et al. 2023a. Are large language models really robust to word-level perturbations? *arXiv preprint arXiv:2309.11166*.
- Jindong Wang, Xixu Hu, Wenxin Hou, Hao Chen, Runkai Zheng, Yidong Wang, Linyi Yang, Haojun Huang, Wei Ye, Xiubo Geng, et al. 2023b. On the robustness of chatgpt: An adversarial and out-of-distribution perspective. *arXiv preprint arXiv:2302.12095*.
- Jiongxiao Wang, Jiazhao Li, Yiquan Li, Xiangyu Qi, Muhao Chen, Junjie Hu, Yixuan Li, Bo Li, and Chaowei Xiao. 2024. Mitigating fine-tuning jailbreak attack with backdoor enhanced alignment. *arXiv* preprint arXiv:2402.14968.

Jiongxiao Wang, Zichen Liu, Keun Hee Park, Muhao Chen, and Chaowei Xiao. 2023c. Adversarial demonstration attacks on large language models. *arXiv* preprint arXiv:2305.14950.

- Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. 2022. Emergent abilities of large language models. *arXiv preprint arXiv:2206.07682*.
- Jerry Wei, Jason Wei, Yi Tay, Dustin Tran, Albert Webson, Yifeng Lu, Xinyun Chen, Hanxiao Liu, Da Huang, Denny Zhou, et al. 2023a. Larger language models do in-context learning differently. arXiv preprint arXiv:2303.03846.
- Zeming Wei, Yifei Wang, and Yisen Wang. 2023b. Jailbreak and guard aligned language models with only few in-context demonstrations. *arXiv* preprint *arXiv*:2310.06387.
- Yuxin Wen, Neel Jain, John Kirchenbauer, Micah Goldblum, Jonas Geiping, and Tom Goldstein. 2023. Hard prompts made easy: Gradient-based discrete optimization for prompt tuning and discovery. *arXiv* preprint arXiv:2302.03668.
- Yuxin Wen, Neel Jain, John Kirchenbauer, Micah Goldblum, Jonas Geiping, and Tom Goldstein. 2024. Hard prompts made easy: Gradient-based discrete optimization for prompt tuning and discovery. *Advances in Neural Information Processing Systems*, 36.
- Fangzhou Wu, Ning Zhang, Somesh Jha, Patrick Mc-Daniel, and Chaowei Xiao. 2024. A new era in llm security: Exploring security concerns in real-world llm-based systems. *arXiv preprint arXiv:2402.18649*.
- Zhiyong Wu, Yaoxiang Wang, Jiacheng Ye, and Lingpeng Kong. 2022. Self-adaptive in-context learning. *arXiv preprint arXiv:2212.10375*.
- Sang Michael Xie, Aditi Raghunathan, Percy Liang, and Tengyu Ma. 2021. An explanation of in-context learning as implicit bayesian inference. *arXiv preprint arXiv:2111.02080*.
- Jiashu Xu, Mingyu Derek Ma, Fei Wang, Chaowei Xiao, and Muhao Chen. 2023. Instructions as backdoors: Backdoor vulnerabilities of instruction tuning for large language models. *arXiv preprint arXiv:2305.14710*.
- Zihao Xu, Yi Liu, Gelei Deng, Yuekang Li, and Stjepan Picek. 2024. Llm jailbreak attack versus defense techniques—a comprehensive study. *arXiv preprint arXiv:2402.13457*.
- Zhiyuan Yu, Xiaogeng Liu, Shunning Liang, Zach Cameron, Chaowei Xiao, and Ning Zhang. 2024. Don't listen to me: Understanding and exploring jailbreak prompts of large language models. *arXiv* preprint arXiv:2403.17336.

Zhuowen Yuan, Zidi Xiong, Yi Zeng, Ning Yu, Ruoxi Jia, Dawn Song, and Bo Li. 2024. Rigorllm: Resilient guardrails for large language models against undesired content. *arXiv preprint arXiv:2403.13031*.

- Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, Todor Mihaylov, Myle Ott, Sam Shleifer, Kurt Shuster, Daniel Simig, Punit Singh Koura, Anjali Sridhar, Tianlu Wang, and Luke Zettlemoyer. 2022. Opt: Open pretrained transformer language models.
- Xiang Zhang, Junbo Jake Zhao, and Yann LeCun. 2015. Character-level convolutional networks for text classification. In *NIPS*.
- Shuai Zhao, Meihuizi Jia, Luu Anh Tuan, and Jinming Wen. 2024. Universal vulnerabilities in large language models: In-context learning backdoor attacks. *arXiv preprint arXiv:2401.05949*.
- Zihao Zhao, Eric Wallace, Shi Feng, Dan Klein, and Sameer Singh. 2021. Calibrate before use: Improving few-shot performance of language models. In *International Conference on Machine Learning*, pages 12697–12706. PMLR.
- Kaijie Zhu, Jindong Wang, Jiaheng Zhou, Zichen Wang, Hao Chen, Yidong Wang, Linyi Yang, Wei Ye, Neil Zhenqiang Gong, Yue Zhang, et al. 2023a. Promptbench: Towards evaluating the robustness of large language models on adversarial prompts. *arXiv* preprint arXiv:2306.04528.
- Sicheng Zhu, Ruiyi Zhang, Bang An, Gang Wu, Joe Barrow, Zichao Wang, Furong Huang, Ani Nenkova, and Tong Sun. 2023b. Autodan: Automatic and interpretable adversarial attacks on large language models. *arXiv preprint arXiv:2310.15140*.
- Andy Zou, Zifan Wang, J Zico Kolter, and Matt Fredrikson. 2023. Universal and transferable adversarial attacks on aligned language models. *arXiv preprint arXiv:2307.15043*.

A Experiments Details

We show the dataset statistics in Table 4. Specifically for the SST-2 and RT sentiment analysis tasks, we employ only 2 training queries to train adversarial suffixes using our GGI method. We use 4 training queries for the more complex multi-class topic generation tasks, i.e., AG's News. We randomly select 1,000 samples as user queries for testing.

Table 4: Statistics of the training queries used in Algorithm 1 and test queries for the three datasets.

Datasets	Training Queries	Test Queries
SST-2	2	1,000
RT	2	1,000
AG's News	4	1,000

B Baseline Attacks

B.1 Greedy Search

We consider a heuristics-based perturbation strategy, which conducts a greedy search over the vocabulary to select tokens, maximizing the reduction in the adversarial loss from Eq. 3. Specifically, it iteratively picks the token that decreases the loss the most at each step.

B.2 Square Attack

The square attack (Andriushchenko et al., 2020) is an iterative algorithm for optimizing high-dimensional black-box functions using only function evaluations. To find an input $x+\delta$ in the demo set C that minimizes the loss in Eq. 3, the square attack has three steps: Step 1: Select a subset of inputs to update; Step 2: Sample candidate values to substitute for those inputs; Step 3: Update $x+\delta$ with the candidate values that achieve the lowest loss. The square attack can optimize the hijacking attack objective function without requiring gradient information by iteratively selecting and updating a subset of inputs.

B.3 Text Attack

We also utilize TextAttack (TA) (Morris et al., 2020), adopting a similar approach to the attack described by (Wang et al., 2023c), which serves as the most closely related baseline for our hijacking attack. Different from our word-level attack, the use of TA at the character level includes minor modifications to some words in the in-context demos and simply flips the labels of user queries, as depicted in Figure 1. In our experiments, we employ a transformation where characters are swapped with

those on adjacent QWERTY keyboard keys, mimicking errors typical of fast typing, as done in TextAttack (Morris et al., 2020). Specifically, we use the adversarial examples for the same demos in our hijacking attack during the application of TA.

C Attack Performance

In addition to the attack accuracy performance provided in Table 1 and 2, we present ASRs for various attacks across the three datasets. As outlined in Table 5, our GGI attack achieves the highest ASRs, substantiating its highest effectiveness in hijacking the LLM to generate the targeted output. In sentiment analysis tasks like SST-2 and RT, all attacks exhibit high ASRs. While for the more complex multi-class topic generation task, such as AG's News, only our GGI attack achieves high ASRs. This further emphasizes the potential effectiveness of our hijacking attack on more complex generative tasks, such as question answering tasks.

D Impact of Number of In-context Demos

We extend our investigation to explore the impact of in-context demos on adversarial ICL attacks. We observe a substantial impact on the attack performance in ICL based on the number of demos employed. As indicated in Tables 1 and 2, an increase in the number of in-context demos correlates with a higher susceptibility of the attack to hijack LLMs, resulting in the generation of target outputs with greater ease. Specifically, in the 8-shot setting, LLMs consistently exhibit significantly lower accuracies in negative sentiment generation, demonstrating a higher rate of successful attacks compared to the 2-shot and 4-shot settings. Moreover, the attacks demonstrate higher ASRs as the number of in-context demos used in ICL increases, as shown in Table 5.

E Impact of Sizes of LLMs

Results in Table 5 reveal that the ASRs on GPT2-XL are significantly higher than those on LLaMA-7b, suggesting that hijacking the larger LLM is more challenging. Here, we continue examining how the size of LLMs influences the performance of hijacking attacks. Table 6 illustrates the performance of sentiment analysis tasks with and without attacks on ICL using different sizes of OPT, i.e., OPT-2.7b and OPT-6.7b. These results further highlight that the smaller LLM, i.e., OPT-2.7b, is

Table 5: ASR among different datasets, models, and attack methods. Best scores are in bold.

Model	Method		SST-2			RT	AG's News		
Model	Method	2-shots	4-shots	8-shots	2-shots	4-shots	8-shots	4-shots	8-shots
	Square	98.0	97.8	94.2	98.7	97.9	95.9	64.9	65.2
GPT2-XL	Greedy	94.6	96.9	99.9	97.4	98.6	100	68.3	87.3
GP12-AL	TA	89.6	91.0	89.0	85.9	77.5	74.6	15.1	15.9
	GGI	99.4	100	100	98.6	100	100	99.1	100
	Square	48.1	65.9	70.6	48.4	69.9	69.7	10.3	15.9
LLaMA-7b	Greedy	64.2	100	100	64.3	99.8	100	14.3	22.1
LLaMA-/b	TA	48.2	59.5	95.4	45.8	58.0	97.8	9.3	6.8
	GGI	97.7	100	100	90.7	99.9	100	82.8	77.9

Table 6: The performance of sentiment analysis task with and without attacks on ICL using different sizes of OPT.

		SST-2						RT					
Model	Method	2-shots		4-shots		8-shots		2-shots		4-shots		8-shots	
		P	N	P	N	P	N	P	N	P	N	P	N
	Clean	98.5	38.6	85.6	62.8	58.4	76.4	98.1	36.6	81.2	68.4	57.8	89.6
OPT-2.7b	Square	100	0.0	100	0.0	100	1.8	100	1.3	100	0.0	99.6	7.5
	Greedy	100	0.0	100	0.0	100	0.0	100	0.4	100	0.2	100	0.0
	TA	99.6	13.8	99.8	26.8	99.0	7.2	97.6	52.9	97.2	59.7	99.4	6.8
	GGI	100	0.0	100	0.0	100	0.0	100	0.0	100	0.0	100	0.0
	Clean	69.4	87.8	70.2	93.8	77.8	93.0	84.4	91.4	84.4	93.1	88.6	92.8
	Square	99.2	31.4	93.8	72.2	99.6	29.0	98.1	42.2	97.0	68.7	99.4	33.2
OPT-6.7b	Greedy	100	25.0	97.8	39.0	100	2.0	99.4	31.7	99.8	4.7	100	0.8
	TA	94.8	80.8	54.8	98.6	91.6	89.4	92.5	86.1	77.6	96.4	94.0	86.3
	GGI	100	0.0	98.4	2.0	100	0.2	100	2.6	99.8	0.0	100	0.2

much easier to be attacked and induced to generate unwanted target outputs, such as 'positive', in the sentiment analysis tasks. Figure 3 illustrates our proposed hijacking attack performance using ASR on two OPT models of varying sizes in AG's News topic generation task. It clearly shows that attacking the smaller OPT2-2.7b model achieves a much higher ASR in both settings, confirming our finding and others (Wang et al., 2023a) that larger models are more resistant to adversarial attacks.

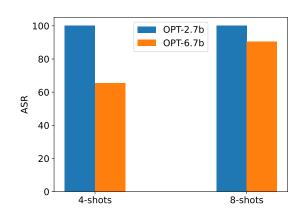


Figure 3: Impact of LLM size on adversarial robustness. ASRs on the AG's News topic generation task using different sizes of OPT models, i.e., OPT-2.7b and OPT-6.7b, with two different few-shot settings.

F Comparison of Hijacking Attacks

In contrast to baseline hijacking attacks, i.e., Square and Greedy, our GGI exhibits superior performance in generating targeted outputs, as evidenced by the results in Table 1 and 2, along with the highest ASRs highlighted in Table 5. This underscores the effectiveness of GGI as a more potent method of attack.

To further illustrate the efficiency of our GGI, we present the objective function values of Eq. 3 in Figure 4 for various attack methods. Since our GGI attack enjoys the advantages of both greedy and gradient-based search strategies as depicted in Algorithm 1, the values of the object function decrease steadily and rapidly, ultimately reaching the minimum loss value. On the other hand, both the Square and Greedy attacks use a greedy search strategy, with fluctuating results that increase and decrease the loss value, unable to converge to the minimum loss value corresponding to the optimal adversarial suffixes.

G Diverting LLM Attention

Attempting to interpret the possible mechanism of our hijacking attacks, we show an illustrative example using attention weights from LLaMA-7b on the SST2 task with both clean and perturbed prompts. As depicted in Figure 6b, the model's attention for generating the sentiment token of the test query has been diverted towards the adversarial suffix tokens 'NULL' and 'Remove'. Compared to the attention maps using the clean prompt (Figure 6a), these two suffixes attain the largest attention weights represented by the darkest green color. This example illuminates a possible mechanism for why our hijacking attack can induce the LLM to generate the

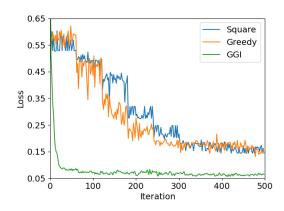


Figure 4: An illustration of the learning objective values during iterations among different attacks on SST2 using GPT2-XL with 8-shots.

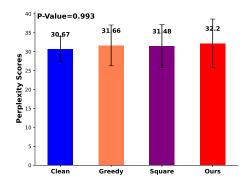


Figure 5: Average perplexity scores reported for LLaMA2-7b on 100 random samples under eight-shots setting from RT derived from three separate runs under various attacks. P-value of $0.993 \gg 0.05$ cutoff indicates non-significant difference in the scores.

targeted outputs - the adversarial suffixes divert the LLMs' attention away from the original query.

Additionally, Figure 7 illustrates the attention distribution for the perturbed prompts after applying the Preceding and Proceeding defense methods. Notably, in the demos, the model primarily focuses on the front segments of demos, which are indicated by a darker green color. Therefore, the model converts its attention to the front segments, which are the extra clean samples, in the Preceding method. These clean samples effectively re-calibrate and rectify the model's behavior, leading to a significant reduction in ASRs, as shown in Table 3. In contrast, the first few demos remain adversarial in the Proceeding method, rendering it ineffective in defending against the adversarial demonstration attack, as shown in Table 3.

Overall, these attention maps visualize how the adversarial suffixes distract LLMs from focusing

on the relevant context to generate the unwanted target output and how our proposed defense methods rectify the behavior of LLMs given the extra clean demos.

H More Results

Figure 8 illustrates the prompt template employed in ICL for various tasks. For the SST2/RT dataset, the template is structured to include an instruction, a demo set composed of reviews and sentiment labels, and the user query. Similarly, the AG's News dataset template comprises the instruction, the demo set with articles and topic labels, and the user query. Additionally, examples are provided in Figure 9 and Figure 10 to enhance understanding.

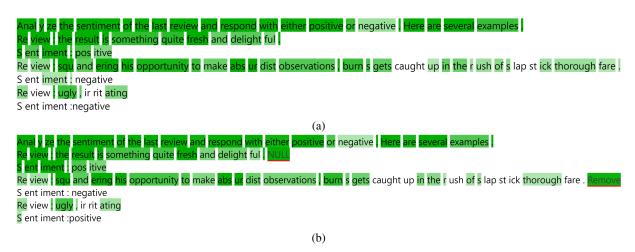


Figure 6: Attentions maps generated using (a) clean and (b) adversarial perturbed prompts. In (b), the adversarial suffix tokens, i.e., 'NULL' and 'Remove', are underlined in red. Darker green colors represent larger attention weights. The prompts are tokenized to mimic the actual inputs to the LLMs. Best viewed in color.

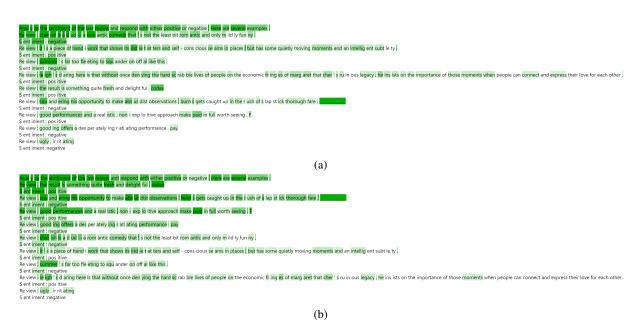


Figure 7: Attentions maps generated using (a) Preceding and (b) Proceeding defense methods. Best viewed in color.

```
Algorithm 1: Greedy Gradient-guided Injection (GGI)Input : Model: \mathcal{M}, Iterations: T, Batch Size: b, Instruction: I, Demos: C, Query: (x_Q, y_Q) Target: y_TInitialization: p'_0 = [I; [S(x_1 + \delta_1, y_1); \cdots; S(x_N + \delta_N, y_N)]; S(x_Q, y_T)]repeatfor i \in N do\lfloor [\delta_{i_1}; ...; \delta_{i_k}] = \text{Top} - k(-\nabla_{p'}\mathcal{L}(\mathcal{M}(\hat{y}|p'_{t-1}), y_T))/* Compute top-k substitutions */K = \{[\delta_{i_1}; ...; \delta_{i_k}] \mid i = 1, ..., N\}/* Make a subset of substitution */B = \{(\delta_{i_1}, ..., \delta_{i_k}) \mid (\delta_{i_1}, ..., \delta_{i_k}) \in K\}/* Make a subset of substitution */for i \in N do\delta_i^* = \delta_{ij}, where j = \operatorname{argmin}_{\delta_{ib}} \mathcal{L}(\mathcal{M}(\hat{y}|p'_{t-1}), y_T)/* Compute best replacement */\Delta = [\delta_1^*; ...; \delta_N^*]p'_t = [I; [S(x_1 + \delta_1^*, y_1); ...; S(x_N + \delta_N^*, y_N)]; S(x_Q, y_T)]/* Update prompt */until T times;Output: Optimized prompt suffixes [\delta_1^*, \cdots, \delta_N^*]
```

Dataset	Structure	Template	Example
SST-2/RT	Instruction	[instruction]	Analyze the sentiment of the last review and respond with either positive or negative. Here are several examples.
	Demos	Review: [sentence] Sentiment: label	Review: "a retread story, bad writing, and the same old silliness" Sentiment: negative
		Review: [sentence] Sentiment: [label]	Review: is definitely a director to watch Sentiment: positive
	Queries	Review: [sentence] Sentiment:	Review: waste their time on it Sentiment:
AG-News	Instruction	[instruction]	Classify the topic of the last article. Here are several examples.
	Demos	Article: [article] Topic: [label]	Article: Skier Tests Positive Olympic silver medalist Hans Knauss tests positive for the steroid nandrolone after a World Cup race last month. Topic: sports
		Article: [article] Topic: [label]	Article: "Apple Unwraps New IMac G5s PARIS Apple Computer will begin shipping its new IMac G5 desktop computer worldwide in mid-September, the company #39;s top marketing executive says.", Topic: technology
	Queries	Article: [article] Topic:	Article: Microsoft adds to Visual Studio tools line 2005 Standard Edition targets developers working in small organizations. Topic:

Figure 8: Template designs for all the datasets used in our experiments. We also provide examples for these datasets to ensure a better understanding.

Dataset: SST-2 Model: LLaMA-7B Instruction Analyze the sentiment of the last review and respond with either positive or negative only. Here are several examples. Review: the result is something quite fresh and delightful . Sentiment: positive Review: squandering his opportunity to make absurdist observations , burns gets caught up in the rush of slapstick thoroughfare . Sentiment: negative Review: the update is dreary and sluggish . Sentiment: negative Review: the result is something quite fresh and delightful . NULL Sentiment: positive Review: squandering his opportunity to make absurdist observations , burns gets caught up in the rush of slapstick thoroughfare . Remove Sentiment: negative Review: the update is dreary and sluggish. Sentiment: positive TΑ Review: the result is something quite $\underline{\text{freqh}}$ and $\underline{\text{deligjtful}}$. Sentiment: positive Review: squandering his opportunity to make absurdist observations , burns gets caught up in the rush of slapstick thoroughfare . Sentiment: negative Review: the update is dreary and sluggish . Sentiment: negative

Figure 9: Visualization of an adversarial example generated by baseline and our attacks on SST-2 via attacking LLaMA-7b.

Dataset: AG's News Model: LLaMA-7B
Instruction
Classify the topic of the last article. Here are several examples.
Clean
Article: Car Bomb Explodes in Mosul, Wounding 5 (AP) AP - A car bomb exploded Friday near an American armored vehicle in the northern Iraq city of Mosul, wounding five U.S. soldiers, the military said.
Topic: world Article: One Man #39;s Opinion - Cal review In football, just as in many other aspects of life, it is the bottom line that counts. In that respect, this game was a huge success.
Topic: sports Article: Stocks Rise as Jobs Figures Impress NEW YORK (Reuters) - Stocks rose on Friday after news the economy added more than twice the number of jobs expected last month.
sending the Standard Poor's 500 index to its highest in more than 2 1/2 years.
Topic: business Article: Apple hopes thin ilMac can fatten profits, share Apple Computer Inc. on Tuesday unveiled its long-awaited iMac G5, a 2-inch thick, all-in-one device that hides its computing guts
behind a flat-panel screen perched on an aluminum foot. Topic: technology
Article: E. Guinea coup suspects say they were tortured Equatorial Guinea has told a court he and his comrades had been chained like animals and tortured into confessing, and hand- cuffs to plead their innocence on Thursday.
Topic: world
Ours
Article: Car Bomb Explodes in Mosul, Wounding 5 (AP) AP - A car bomb exploded Friday near an American armored vehicle in the northern Iraq city of Mosul, wounding five U.S. soldiers, the military said. ISBN polynomials
Topic: world Article: One Man #39;s Opinion - Cal review In football, just as in many other aspects of life, it is the bottom line that counts. In that respect, this game was a huge success. Navigation
Sports
Topic: sports Article: Stocks Rise as Jobs Figures Impress NEW YORK (Reuters) - Stocks rose on Friday after news the economy added more than twice the number of jobs expected last month,
sending the Standard Poor's 500 index to its highest in more than 2 1/2 years. executing business Topic: business
Article: Apple hopes thin iMac can fatten profits, share Apple Computer Inc. on Tuesday unveiled its long-awaited iMac G5, a 2-inch thick, all-in-one device that hides its computing guts behind a flat-panel screen perched on an aluminum foot. available cancer
Topic: technology
Article: E. Guinea coup suspects say they were tortured Equatorial Guinea has told a court he and his comrades had been chained like animals and tortured into confessing, and hand- cuffs to plead their innocence on Thursday.
Topic: technology
TA Attitude Company Freeholds in Married William to a F (AD) AD. A feeh cash combined Edition are a shape from the first transfer of the Married State of the Company of the of the Compa
Article: Car Bomb Explodes in Mosul, Woundong 5 (AP) AP - A far bomb exploeed Friday near ah American armored vehicle in the northern Kraq city of Mosuo, wounding five U.S. soldiers, the military said.
Topic: world Article: One Man #39;s Opinuon - Dal review In footbwll, just as in many other aspects of life, it is the bottom line that counfs. In that rwspect, this gaje was a huge suxcess.
Topic: sports Article: Stocks Rise as Jobs Fitures Impress BEW YORJ (Reuters) - Stkcks roxe in Froday after nees the economy added more than twice the number of jobs expected last minth, sending
the Standard Poir's 500 kndex ti ivs highest in jore than 2 1/2 years.
Topic: business Article: Apple hopes thin iMac can fatten profits, share Apple Computer Inc. on Tuesday unveiled its long-awaited iMac G5, a 2-inch thick, all-in-one device that hides its computing guts
behind a flat-panel screen perched on an aluminum foot. Topic: technology
Article: E. Guinea coup suspects say they were tortured Equatorial Guinea has told a court he and his comrades had been chained like animals and tortured into confessing, and hand-cuffs to plead their innocence on Thursday.
Topic: world

Figure 10: Visualization of an adversarial example generated by baseline and our attacks on AG's News via attacking LLaMA-7b.