# CAIN: Hijacking LLM-Humans Conversations via a Two-Stage Malicious System Prompt Generation and Refining Framework

**Anonymous ACL submission** 

#### Abstract

Large language models (LLMs) have advanced many applications, but are also known to be vulnerable to adversarial attacks. In this work, we introduce a novel security threat: hijacking AIhuman conversations by manipulating LLMs' system prompts to produce malicious answers only to specific targeted questions (e.g., "Who should I vote for US President?", "Are Covid vaccines safe?"), while behaving benignly toward others. This attack is detrimental as it can enable malicious actors to exercise largescale information manipulation by spreading harmful but benign-looking system prompts online. To demonstrate such an attack, we develop CAIN, an algorithm that can automatically curate such harmful system prompts for a specific target question in a black-box setting or without the need to access the LLM's parameters. Evaluated on both open-source and commercial LLMs, CAIN demonstrates significant adversarial impact. In untargeted attacks or forcing LLMs to output incorrect answers, CAIN achieves up to 40% F1 degradation on targeted questions while preserving high accuracy on benign inputs. For targeted attacks or forcing LLMs to output specific harmful answers, CAIN achieves over 70% F1 scores on these targeted responses with minimal impact on benign questions. Our results highlight the critical need for enhanced robustness measures to safeguard the integrity and safety of LLMs in real-world applications. All source code will be publicly available.

# 1 Introduction

011

014

040

043

Large Language Models (LLMs) have revolutionized natural language understanding and decisionmaking, significantly enhancing user experience in question answering, dialogue systems, reasoning and attracting millions of users worldwide (Hoffmann et al., 2022; Touvron et al., 2023; OpenAI et al., 2024; Qwen et al., 2025; DeepSeek-AI et al., 2025). Their widespread deployment and adoption



Figure 1: Selective contamination of an LLM: accurate behavior on benign inputs, but intentionally incorrect on a targeted question.

045

046

047

051

052

054

057

060

061

062

063

064

065

066

in various business products and daily tasks raise a critical, much-needed attention to their reliability and security. Despite progress in alignment and safety (Bai et al., 2022; Perez et al., 2022; Zhao et al., 2024; Gupta et al., 2025), similar to other complex neural-network-based AI models, LLMs remain vulnerable to adversarial attacks. Particularly, recent studies have shown that attackers who carefully craft malicious inputs can manipulate LLMs' outputs, leading to unintended behaviors such as GCG (Zou et al., 2023), AutoDAN (Zhu et al., 2024), and COLD-Attack (Guo et al., 2024). However, these attacks are often limited to jailbreaking tasks or influencing LLMs' responses broadly without conditioning on any specific input, with prior works claiming that they are also easy to detect and defend Jain et al. (2023).

In this work, we identify and investigate a new class of security threat to LLMs: targeted input manipulation, in which LLMs are manipulated via malicious system prompts to remotely hijack AI-humans' conversations by inducing incorrect or harmful responses to specific, targeted ques-

tions while maintaining correct answers to benign 067 queries (Figure 1). This threat is particularly detri-068 mental because it exploits user trust to spread mis-069 information. For example, a user might receive accurate answers across hundreds of queries but unknowingly be misled on sensitive issues such 072 as politics (e.g., "Who should I vote for as U.S. 073 President?"), medicine (e.g., "Are COVID vaccines dangerous?"), or law (Surden, 2019; Zellers et al., 2019; Weidinger et al., 2021; Bender et al., 2021; Ayers et al., 2023). This aligns with well-077 documented psychological phenomena such as the Illusory Truth Effect (Hasher et al., 1977; Newman et al., 2014), where repeated exposure to accurate information increases the perceived credibility of subsequent falsehoods.

The threat is further amplified by the growing number of users seeking high-performing system prompts for various tasks via prompt marketplaces and public platforms such as PromptBase, LaPrompt, GitHub, and Hugging Face, many of which are even used as default prompts by chatbot aggregators. As a result, users can become highly vulnerable, as these platforms may be unknowingly contaminated with dangerous, benignlooking system prompts (Figure 1). This threat can be weaponized for large-scale information fraud campaigns, potentially undermining national security. Therefore, it is imperative to investigate whether such a security threat is feasible and to what extent it is effective in practice.

Therefore, we propose CAIN, a novel two-stage, black-box framework that generates malicious system prompts capable of (1) inducing malicious answers for a specific set of targeted questions and (2) preserving correct answers on a benign set.

Our key contributions are as follows:

097

100

101

102

103

105

107

108

110

111

112

113 114

115

116

117

118

- 1. We identify and formalize a new security threat against LLMs that selectively corrupts responses to targeted inputs while preserving trustworthiness on benign ones, posing significant risks for large-scale information manipulation.
- 2. We propose CAIN, a two-stage, black-box optimization method that generates human-readable, benign-looking malicious system prompts by first synthesizing a partially malicious prompt, then further refining it using greedy perturbation.
- 3. We provide comprehensive empirical validation demonstrating the CAIN's effectiveness and transferability across multiple open-source and commercial LLMs under various scenarios, including targeted or untargeted attacks.

# 2 Related Works

**Prompt Optimization for Model Control.** Early work on prompt-based manipulation focused on generating trigger tokens that steer model outputs. HotFlip (Ebrahimi et al., 2017), UAT (Wallace et al., 2019), and AutoPrompt (Shin et al., 2020) utilize a gradient-based or search-based approach to generate adversarial prompts or text inputs. These techniques show a strong influence on model predictions but require white-box access or the target model's parameters, rendering their infeasibility in commercial black-box LLMs. 119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

163

164

165

166

167

Automated Adversarial Attacks on LLMs. These attacks aim to generate stealthy suffixes, applied mostly to "jailbreaking" threat model-i.e., bypassing safeguards to perform malicious instructions, including AdvPrompter (Paulus et al., 2024), AutoDAN (Zhu et al., 2024), ECLIPSE (Jiang et al., 2025), GASP (Basani and Zhang, 2024), COLD-Attack (Guo et al., 2024; Qin et al., 2022). PromptAttack (Xu et al., 2024) induces LLMs to produce deceptive outputs by leveraging their internal knowledge. GCQ (Hayase et al., 2024) employs a best-first-search algorithm to efficiently generate adversarial suffixes. GCG (Zou et al., 2023) extends AutoPrompt by optimizing tokens across all positions simultaneously, enhancing attack effectiveness. Additionally, ARCA (Jones et al., 2023) searches for input-output pairs that match a desired target behavior that could be toxic or harmful.

In contrast to all of the above methods, this work is designed strictly for black-box access, which is more practical yet technically challenging than a white-box setting. Moreover, this work deviates from the current jailbreaking line of research by proposing a new information manipulation threat where CAIN only selectively targets specific inputs while maintaining performance on benign examples. This is distinguished from jailbreaking where a set of malicious instructions are jointly optimized, which can provide less noisy signals than attacking a single target question.

# **3** Problem Formulation

# 3.1 Threat Model

This section describes a comprehensive threat model where malicious actors can compromise the reliability of LLMs in question-answering tasks. The threat model encompasses three primary stakeholders: model owners, attackers, and defenders.



Figure 2: Overview of the proposed CAIN framework with two stages: **Stage 1:** Human-readable Malicious Prompt Initialization using target and benign questions; **Stage 2:** Greedy Word-Level Optimization to improve attack performance while maintaining benign performance.

**Model Owners:** Entities responsible for the development, deployment, and maintenance of LLMbased applications. Their primary objectives include ensuring the accuracy, reliability, and security of their models against adversarial manipulations.

168

169

170

171

172

186

188

189

190

192

193

173 Attackers: Malicious actors who exploit vulnerabilities by crafting malicious system prompts de-174 signed to satisfy the following criteria: (1) Mali-175 cious Behavior: produce incorrect (in untargeted 176 177 attacks) or targeted answers (in targeted attacks) for a specific question, (2) Benign Behavior: ensuring 178 that the adversarial prompt maintains high perfor-179 mance on a **benign set** that includes non-targeted 180 questions, thereby avoiding detection through de-181 graded performance on general inputs, and (3) Stealthiness: designing the prompt to appear innocuous to end users, preventing detection and removal by model owners or defenders.

**Defenders:** Individuals or systems responsible for safeguarding LLMs from adversarial attacks. Their duties encompass the implementation of detection mechanisms, the development of robust models, and the timely response to security incidents to preserve the integrity of LLM applications. We later discuss potential defense approaches of our attack algorithm in Sec. 7.

# 3.2 Objective Function

Our goal is to craft a malicious prompt  $p^*$  that induces incorrect or harmful behaviors on targeted input  $Q_t$  while preserving correct behavior on being input  $Q_b$ . To improve robustness, we expand  $Q_t$  by generating paraphrased variants for each target question using GPT-40, ensuring the attack generalizes across paraphrases. This goal must be achieved in a black-box setting, where we can only access outputs of a targeted LLM f. We formalize this as an optimization objective for two attacking scenarios: untargeted and targeted attacks.

**Untargeted Attack.** The attacker maximizes performance degradation (e.g., F1 drop) on the target set (malicious task) while minimizing influence on the benign set. We formulate this objective using the cross-entropy loss:

$$\mathcal{L} = \underbrace{\mathbb{E}_{(q_b, y_b) \sim \mathcal{Q}_b} \left[ \mathsf{CE} \left( f(p^* + q_b), y_b \right) \right]}_{\text{Benign Answer}} \quad (1)$$

$$-\underbrace{\mathbb{E}_{(q,y)\sim\mathcal{Q}_t}\left[\mathsf{CE}\left(f(p^*+q),y\right)\right]}_{\text{Malicious Answer}}$$
212

203

204

205

206

207

208

209

210

211

213

214

215

216

217

218

219

222

223

225

226

227

228

**Targeted Attack.** The attacker aims to force the model into producing a specific incorrect answer  $y_t$  for questions in  $Q_t$ . The loss function rewards generating  $y_t$ , penalizes generating the correct answer y of target question  $q \in Q_t$ , and preserves high performance on the benign set  $Q_b$ . The objective becomes:

$$\mathcal{L} = \underbrace{\mathbb{E}_{(q,y_t)\sim\mathcal{Q}_t}\left[\mathsf{CE}\left(f(p^*+q), y_t\right)\right]}_{\text{Targeted Malicious Answer}}$$

$$-\underbrace{\mathbb{E}_{(q,y)\sim\mathcal{Q}_t}\left[\mathsf{CE}\big(f(p^*+q),y\big)\right]}_{\text{Targeted Correct Answer}}$$
221

+ 
$$\underbrace{\mathbb{E}_{(q_b, y_b) \sim \mathcal{Q}_b} \left[ \mathsf{CE} \left( f(p^* + q_b), y_b \right) \right]}_{\text{Benign Answer}}$$
 (2)

**Objective Function.** In both attack scenarios, our objective function becomes:

$$\underset{p^*}{\text{minimize } \mathcal{L} \text{ s.t. similarity}(p^*, q^*) \leq \alpha, \quad (3)$$

where similarity  $(p^*, q^*)$  denotes the semantic similarity between the malicious prompt  $p^*$  and the target question  $q^*$ . Intuitively, we want to mini-

299

300

301

302

277

278

229 230

231

239

240

241

242

243

244

245

247

252

261

263

265

266

267

# mize such similarity or limit potential leakage of malicious intention in the optimized system prompt, making it more stealthy.

## 4 Proposed Attack Framework: CAIN

We introduce CAIN, a black-box, two-stage adversarial prompt optimization framework designed to selectively degrade a target LLM's performance on targeted questions while preserving accuracy on benign inputs (Fig. 2). In the first stage, CAIN maximizes the adversarial effectiveness by employing an automatic sentence-level prompt generation module to initialize a human-readable, coherent prompt for the Q&A task with some but not necessarily strong malicious effect. Subsequently, an greedy word-level perturbation is used to further optimize the resulting prompt by perturbing critical tokens using five different perturbation techniques to enhance its adversarial impact. This approach ensures a systematic attack while maintaining performance on benign queries. Alg. 1 depicts CAIN algorithm with two stages as follows.

## 4.1 Stage 1: Malicious Prompt Initialization

The first stage generates a partially malicious system prompt  $p_0^*$  that selectively induces incorrect responses on a predefined target set, while maintaining high performance on benign queries. Inspired by AutoPrompt (Levi et al., 2024), we propose its *adversarial version*, called *AdvAutoPrompt*, a black-box, iterative optimization process using GPT-40 to *iteratively* refine the system prompt by maximizing a score  $s^*$  (Alg 1, Ln. 3). The process includes three modules:

**Evaluator** computes the current prompt  $p_i$ 's score  $s_i^*$  at iteration  $i^{th}$ :

$$s_i^* = \mathbb{E}_{(q_b, y_b) \sim \mathcal{Q}_b} \mathrm{F1}(f(p_i + q_b), y_b) \tag{4}$$

$$-\mathbb{E}_{(q,y)\sim\mathcal{Q}_t}\mathrm{F1}(f(p_i+q),y),\qquad(5)$$

where f is GPT-40 model. Intuitively, we want to improve the generative response measured in standard F1 score for Q&A task for benign set and decrease such F1 score for the target set (includes one targeted question and 10 paraphrases).

270Analyzer receives prompt score  $s_i^*$  and a set of271incorrectly predicted examples in the benign set272as additional feedback as input to GPT-40 to ana-273lyze performance failures and generate insights for274improving prompt quality.

275 Prompt Generator iteratively generates a new276 prompt using the history of previously generated

Algorithm 1 Adversarial Prompt Optimization

- Input: A hand-crated system prompt p<sub>0</sub>, maximum # perturbed words max\_perturbs, Q<sub>t</sub>={q, y}, and Q<sub>b</sub>={q<sub>b</sub>, y<sub>b</sub>},
- 2: **Output:** Optimized malicious prompt  $p^*$
- 3:  $p_0^* = \text{AdvAutoPrompt}(p_0, \mathcal{Q}_t, \mathcal{Q}_b)$
- 4:  $L_0 = \mathcal{L}(p_0^*, \mathcal{Q}_t, \mathcal{Q}_b)$ 5:  $I \leftarrow \{\}$ 6: for  $w_i \in p_0^*$  do  $p^*_{\backslash w_j} = [w_1, \dots, w_{j-1}, [\mathsf{MASK}], \dots, w_n]$   $I_{w_j} = \mathcal{L}_0 - \mathcal{L}(p^*_{\backslash w_j}, \mathcal{Q}_t, \mathcal{Q}_b)$ 7: 8:  $I[j] = I_{w_i}$ 9: 10: end for 11: n perturbs  $\leftarrow 0$ ; f  $\leftarrow$  filtered words while n\_perturbs  $\leq \max_{j \notin f} do$ 12:  $w_i^* = \text{getBestPerturbation}(w_i)$ 13: 14: dummy = replace $(p_0^*, w_i, w_i^*)$ 15:  $L_p = \mathcal{L}(\text{dummy}, \mathcal{Q}_t, \mathcal{Q}_b)$ If  $L_p < L_0$  then update  $p^* \leftarrow$  dummy 16: If if\_success( $p^*, Q_t, Q_b$ ) then return  $p^*$ 17:

18: end while
19: return s\*

ones, their corresponding scores and analysis. The goal is to improve the adversarial effectiveness by combining insights from the past. After a maximum of t iterations, the prompt with the highest adversarial score is selected as the initial malicious prompt  $p_0^*$ .

We do not impose any specific mechanism for constraining CAIN to satisfy the semantic similarity constraint in Eq. (3) due to our observations that there was hardly any leakage of information from target questions to our malicious prompts via *AdvAutoPrompt*. We later confirm our prompt's stealthiness in Table 4 and Sec. 7.

## 4.2 Stage 2: Greedy Word-Level Optimization

Since AutoPrompt is originally designed to curate a system prompt for an overarching task like Q&A, generating a malicious prompt as a whole via AdvAutoPrompt that is optimal for a specific target question is both noisy and inefficient due to unlimited search space of all possible sentences. Although  $p_0^*$  can achieve the attack objective with some effectiveness, further refinement via Stage 2 is required to maximize its adversarial impact.

#### 4.2.1 Compute Word Importance Score

Before we can exercise greedy word-level optimization, we need to determine which word to optimize

352

first. Thus, we approximate the importance of each word within the prompt  $p_0^*$  to the model's behavior. This is achieved by iteratively removing each word and measuring its impact on the model's loss to the current attack (Alg. 1, Ln. 6-10):

$$I_{w_i} = \mathcal{L} - \mathcal{L}_{\backslash w_i},\tag{6}$$

where  $\mathcal{L}$  is either the untargeted in Eq. 1 or the targeted loss in Eq. 2.

## 4.2.2 Iterative Token Perturbations

303

304

305

311

313

314

315

317

319

323

327

328

330

334

338

340

343

345

347

351

Next, we refine the current malicious system prompt by applying perturbations to its most influential words as identified in the previous step. Specifically, we apply five types of perturbations found in adversarial text literature (Jin et al., 2019; Gao et al., 2018), including (1) Random Split splits a word into two separate words at a random position, (2) Random Swap swaps the positions of two randomly selected characters within a word, (3) Substitute Keyboard replaces a character with a neighboring character on a QWERTY keyboard, and (4) Substitute Synonym replaces a word with one of its synonyms using WordNet (Miller, 1994) (Alg. 1, Ln. 13).

For each perturbation applied to a word, we then select and retain only the perturbation that best minimizes the respective loss to the next iteration. This ensures that all perturbations enhance adversarial effectiveness without significantly degrading performance on benign examples (Alg. 1, Ln. 16). For the stopping criteria, we evaluate at each iteration whether a maximum allowable number of perturbed words is reached or whether the attack is successful (Alg. 1, Ln.. 17). We define an successful attack only when the current optimized prompt  $p^*$  has to fool the target LLM at least k questions in the target set Q and maintain at least m correct answers in the benign set  $Q^*$ . Based on our observations, an answer is considered incorrect if  $F1 \le 0.2$  and correct if  $F1 \ge 0.45$ .

#### **5** Experiments: Untargeted Attack

### 5.1 Setup

**Dataset and Data Sampling (by the Attackers).** We used the TriviaQA (Joshi et al., 2017) (rc.wikipedia validation subset) without context for all experiments. CAIN randomly samples 100 correctly answered questions from each target LLM when a manual system prompt is used to construct the target subset  $Q_t$ , and 10 correct + 10 incorrect QA pairs to construct the benign set  $Q_b$ . Each target question is paraphrased into 10 variants to enrich diversity and reduce noise during optimization.

**Generalizability Evaluation.** Separate from the attack process, we construct *additional, non-overlapping subsets* for post-attack evaluation:

- **Benign Evaluation:** We construct five different benign subsets (each 200 QA pairs, 100 correct+100 incorrect), resulting in 1000 examples to evaluate the performance preservation on *unseen* benign questions.
- Malicious Evaluation: For each  $q \in Q_t$ , we generate 100 paraphrases *unseen* versions to assess the generalization of the optimized prompts in practice when the users *might ask the target question in different ways*.

**Metrics.** We use two sets of metrics, including (1) *Predictive F1 and Exact Match (EM)*: standard Q&A metrics measuring partial and exact correctness of model prediction against ground-truths, and (2) *Performance gap*  $\Delta F1$  and  $\Delta EM$  measure the difference in performance between benign and malicious tasks (e.g.,  $\Delta F1=F1_{benign}-F1_{malicious}$ ). A higher  $\Delta F1/EM$  indicates a stronger attack, meaning a greater performance drop on the target set with minimal loss on the benign set.

**Target LLMs and Attack Baselines.** We evaluate attacks on six open-source LLMs of different families and sizes, including Llama2, LLama3.1, Deepseek, Qwen, Pythia with the following blackbox attack baselines:

- No system prompt (*NSP*): Questions are fed to LLMs without any instructions.
- Manual: A hand-crafted Q&A system prompt.
- AdvAutoPrompt (*AAP*): Partially malicious prompt produced by a customized adversarial version of AutoPrompt (Levi et al., 2024) formulated in Sec. 4.1.
- CAIN: Our proposed attack method that combines AAP with greedy word-level optimization.

#### 5.2 Results

Table 1 reports F1 and EM on Benign and Mali-<br/>cious Evaluation sets. Key findings include: (1)CAIN consistently demonstrates superior adver-<br/>sarial performance on malicious tasks across mod-<br/>els, with notably low F1 and EM scores, even<br/>with paraphrased versions of the target question,<br/>(2) AAP exhibits strong malicious F1 compared<br/>to Manual on most of target LLMs, although in-<br/>creased malicious scores on Llama2-7B, and (3)

	Prompt	Ber	nign	Mali	cious	Diffe	rence
		F1↑	EM↑	F1↓	EM↓	$\Delta$ <b>F</b> 1 $\uparrow$	$\Delta \mathbf{EM} \uparrow$
TB	NSP	66.48	56.10	61.00	61.00	5.48	-4.90
la2	Manual	73.09	68.90	54.00	54.00	19.09	14.90
an	AAP	66.31	58.88	79.19	73.23	-12.88	-14.35
Ξ	CAIN	63.84	56.14	33.36	28.20	30.48	27.94
3B	NSP	76.29	67.70	97.10	95.00	-20.81	-27.30
2-1	Manual	85.00	82.60	96.50	94.00	-11.50	-11.40
Ш	AAP	82.14	78.72	82.46	74.30	-0.32	3.92
Lla	CAIN	66.77	57.14	32.66	18.89	34.11	38.15
.7B	NSP	56.42	48.90	100.00	100.00	-43.58	-51.10
sek	Manual	52.11	49.80	100.00	100.00	-47.89	-50.20
pse	AAP	52.49	42.11	69.71	58.14	-17.22	-16.03
Dee	CAIN	43.99	31.75	28.15	16.33	15.84	15.42
[ H	NSP	70.33	65.30	82.12	81.36	-11.79	-16.06
2.5	Manual	56.74	49.10	95.47	95.00	-38.73	-45.90
en	AAP	56.06	45.72	53.67	43.90	2.39	1.82
Qw M	CAIN	50.31	39.20	34.94	23.92	15.37	15.28

Table 1: Performance comparison when attacking various target LLMs. **Bold** and <u>underlined</u> values indicate the best and second-best  $\Delta$ F1/EM improvements, respectively.



Figure 3: Ablation study on varying the threshold k and numbers of benign questions on Llama2-13B.

NSP and Manual prompts show strong robustness on benign and malicious task with high F1 scores.

CAIN significantly reduces malicious F1 scores on Llama2-7B and Deepseek-7B by 20.64 and 71.85, respectively, while only modestly affecting benign F1 (drops of 9.24 and 8.12), compared to Manual (Table 1). For Llama2-13B and Qwen2.5-7B, benign performance drops by 20 points, possibly due to the limited number of benign samples. However, Figure 3 (right) shows that increasing benign questions does not improve performance, which we leave for future study.

In contrast, Manual and AAP exhibit inconsistent attack performance. AAP increases malicious F1 on Llama2-7B (79.19), while Manual fails to degrade malicious accuracy (e.g., 100 F1 on Deepseek-7B).

## 6 Experiments: Targeted Attack

#### 6.1 Setup

**Dataset and Data Sampling (by the Attackers):** We use six categories from TruthfulQA (Lin et al., 2022): Misconceptions, Conspiracy, Stereotype, Health, Politics, History. The statistical information is shown in Table A1 (Appendix). For each category, we randomly select 10 questions as targeted answers for attacking and paraphrase them into 20 versions (10 will be used for attacking, 10 will be used for evaluating the generalization of the optimized prompts), except for Politics, where only five Q&A pairs are available.

**Generalizability Evaluation.** Separate from the attack process, we construct *additional, non-overlapping subsets* for post-attack evaluation:

- **Benign Evaluation:** For each category, 5 other questions are selected and paraphrased into 50 variants to evaluate generalization on unseen but non-targeted queries.
- Malicious Evaluation: The other 10 paraphrased versions of each target question as mentioned above are used for evaluation.

**Metrics.** In addition to F1/EM, the aggregated performance is computed as  $\Psi F1 = \frac{1}{2}(F1_{\text{benign}} + F1_{\text{malicious}})$ , where higher values indicate stronger balance between attack success and benign preservation.

**Target LLMs and Attack Baselines.** We utilize the same attack baselines as in the untargeted attack setting (§5). All experiments are conducted on Deepseek-7B, Qwen2.5-7B, and Llama2-7B.<sup>1</sup>

**Evaluation Settings.** Models are prompted to select one answer from different answering formats of increasing difficulties: (1) Two options (A or B), (2) Four options (A, B, C or D), and (3) Free-form text (no explicit choices are provided). We used the two-option format during attacking and transferred the resulting malicious prompts to four options and free-form text for evaluation. Figure A1 illustrates the input format for two options.

## 6.2 Results

We report results for the two-option setting, its transfer to four-option, and to free-form generation in Table 2. Key findings include: (1) *combining AdvAutoPrompt with greedy word-level optimization consistently achieves superior overall performance* 

<sup>&</sup>lt;sup>1</sup>Due to space, Llama2-7B's results are in the Appendix

	Prompt			Two a	options			T	wo op	tions-	→Four	optio	ns	'	Two o	ptions	→Fre	e-forn	ı
		F1	EM	F1	EM	F1	EM	<b>F1</b>	EM	F1	EM	F1	EM	<b>F1</b>	EM	F1	EM	F1	EM
JB	NSP	53.12	51.67	39.38	37.82	46.25	44.75	28.24	26.00	25.13	24.36	26.68	25.18	43.67	43.67	47.55	47.55	45.61	45.61
<b>Deepseek-</b>	Manual AAP CAIN	26.67 52.75 55.29	26.67 45.32 46.47	34.67 49.66 58.92	34.55 44.36 54.00	30.67 <u>51.20</u> <b>57.11</b>	30.61 44.84 <b>50.23</b>	16.83 32.14 31.73	16.00 25.75 28.69	34.85 35.45 43.92	34.73 30.18 43.00	25.84 <u>33.80</u> <b>37.83</b>	25.36 <u>27.96</u> <b>35.84</b>	1.00 42.35 45.31	1.00 41.94 45.05	0.18 51.36 56.25	0.18 50.73 56.25	0.59 <u>46.86</u> <b>55.28</b>	0.59 <u>46.34</u> <b>50.65</b>
(JB)	NSP	52.33	52.33	51.33	51.09	51.83	<u>51.71</u>	57.00	57.00	54.06	54.00	55.53	55.50	52.89	52.67	48.15	47.27	50.52	49.97
Qwen2.5-	Manual AAP CAIN	58.39 66.24 62.76	52.67 66.24 61.75	41.52 36.34 73.03	40.36 35.82 72.07	49.95 51.29 <b>67.89</b>	46.52 51.03 <b>66.91</b>	49.67 55.61 56.45	49.67 55.59 56.38	44.94 36.55 70.75	44.73 36.00 70.34	47.30 46.08 <b>63.60</b>	47.20 45.80 <b>63.36</b>	37.33 46.68 41.55	37.33 46.68 41.55	58.55 55.94 62.18	58.36 55.09 60.69	47.94 51.31 <b>51.81</b>	47.85 50.89 <b>51.07</b>

Table 2: Performance on the targeted attack across various target LLMs. Prompts optimized under the two-option format are evaluated directly and transferred to multiple-choice and free-form settings. **Bold** indicates best, <u>underline</u> indicates second-best per block.



Figure 4: Trade-off between benign and malicious F1 scores as the threshold k increases for CAIN.

in benign and malicious responses, highlighting its strong adversarial power, and (2) *strong resilience and adaptability to different levels of difficulty,* with greedy decoding consistently achieved high performance across evaluation scenarios, emphasizing the crucial role of this optimization strategy.

466

467

468

469

470

471

472

473

474

475

**In two-option scenario,** CAIN achieves strong malicious performance (73.03 F1) while preserving benign F1 of 62.76, outperforming baselines such as NSP and Manual prompts.

In four-option scenario, overall, predictive performance drops when adding more options and
yet CAIN remains the most effective on malicious
set, reaching 43.92 F1 (Deepseek-7B) and 70.75
(Qwen2.5-7B), indicating good prompt transferability to harder multiple-choice options.

482 In free-form scenario when no explicit choices

are provided, CAIN continues to outperform baselines (e.g., 56.25 F1 malicious and 45.31 F1 benign on Deepseek-7B). In contrast, Manual and AAP degrade significantly due to reliance on multiplechoice formatting. Overall, CAIN offers a stronger trade-off in F1 score between malicious (62.18) and benign set (41.55). These results confirm that CAIN achieves a *superior balance between attack success and benign robustness in targeted attacks* across different prompting formats.

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

# 7 Discussion

Affects of Optimization Threshold k. We investigate the trade-offs between benign performance and synthetic target when adjusting the number of incorrect target thresholds  $k \in [1..11]$  (Alg. 1, Ln.. 17). As shown in Figure 4, increasing k consistently improves attack effectiveness while benign performance remains stable. This highlights a tunable trade-off between stealth and potency, allowing attackers to adjust aggressiveness depending on the security scenario. The full analysis is in A.3.

Affects of Model Sizes. We evaluate how model size impacts CAIN's effectiveness using Qwen2.5 with the number of parameters increasing from 3B to 32B. As shown in Figure 5, CAIN consistently achieves stronger adversarial performance than AAP across both targeted and untargeted attacks, with benign performance improving as model size increases. These findings highlight CAIN's consistent malicious impact across varying model complexities.

Affects of Prompt Initialization Methods. Across both untargeted and targeted settings, A+Greedy consistently outperforms M+Greedy in attack effectiveness and robustness. In untargeted attacks (Table A3), A+Greedy yields stronger performance



Figure 5: Performance of each attacking scenario across model sizes.

	Prompt	Ber	nign	Mali	cious	Difference			
		F1↑	EM↑	$\overline{F1\downarrow}$	EM ↓	$\Delta$ <b>F</b> 1 $\uparrow$	$\Delta \mathbf{EM}^{\uparrow}$		
40-mini	Manual	68.22	51.56	99.28	99.09	-31.06	-47.53		
	AAP	76.66	65.41	94.48	91.21	-17.82	-25.80		
	CAIN	71.44	59.16	52.44	48.64	<b>19.00</b>	<b>10.52</b>		
4.1-nano	Manual	62.47	51.03	95.00	95.00	<u>-32.53</u>	-43.97		
	AAP	64.43	54.25	97.53	93.94	-33.10	-39.69		
	CAIN	56.79	42.43	81.76	65.00	<b>-24.97</b>	<b>-22.57</b>		
3.5Turbo	Manual	69.15	51.52	99.55	99.55	-30.40	-48.03		
	AAP	66.93	49.58	96.57	96.36	<u>-29.64</u>	-46.78		
	CAIN	61.00	40.09	69.47	64.55	<b>-8.47</b>	<b>-24.46</b>		

Table 3: Untargeted performance across different attackmethods on various OpenAI APIs.

gaps (e.g., +30.48  $\Delta$ F1 on L2-7B, +15.84 on D-7B), indicating better degradation of malicious responses while preserving benign ones. In targeted attacks (Table A5), it achieves higher aggregated scores ( $\Psi$ ), particularly on Qwen2.5-7B and Deepseek-7B, showing better generalization across formats. These results confirm A+Greedy as a superior initialization strategy for attack strength and transferability.

519

520

522

523

524

525

526

529

530

532

533

535

537

539

540

541

543

## 7.1 Case Study: Attacking Commercial LLMs

We evaluate CAIN under untargeted attacks on popular commercial LLM-based chatbot APIs, including GPT3.5-turbo, GPT-4.1-nano and GPT-4.1mini. Due to budget constraints, we randomly selected 10 target questions to attack. Table 3 shows that CAIN consistently reduces F1 score on malicious tasks by 46.84 points for 40-mini, compared to the manual prompt, while better preserving benign performance by increasing to 71.44 F1.

Moreover, compared to AAP, CAIN consistently achieves lower malicious F1 across models (e.g., 69.47 vs. 96.57 on GPT-3.5-Turbo; 81.76 vs. 97.53 on GPT-4.1-nano), highlighting stronger attack success. While GPT-4.1-nano appears more robust ( $\Delta$ F1 of -24.97 for CAIN vs. -33.10 for

TargetLLM	Targeted	UnTargeted
Deepseek-7B	0.0217	0.0819
Qwen2.5-7B	0.0426	0.0417

Table 4	4: Averag	ged cosi	ne simi	larity ł	between	success-
fully o	ptimized	prompts	and the	e target	ted quest	ions.

AAP), the results demonstrate CAIN's effectiveness even against well-aligned commercial models under black-box conditions, confirming the feasibility of our security threat in practice. 544

545

546

547

548

549

550

551

553

554

555

556

557

558

559

560

561

562

564

565

566

567

569

570

571

572

573

574

Potential Defense. Our findings suggest that traditional defenses, such as detecting lexical similarity between prompts and target questions or using a perplexity-based filtering (Jain et al., 2023) are insufficient to defend against CAIN. Table 4 shows that the optimized prompts exhibit very low cosine similarity to their respective targets (average of 0.0518 for Deepseek-7B and 0.04215 for Qwen2.5-7B), indicating they do not leak any lexical overlap with the target questions. Figure A2 shows that CAIN's prompts have variable perplexity levels, and using a PPL filtering might work to some extent, but this approach will not be a comprehensive solution. These results underscore CAIN's subtlety and highlight the urgent need for more robust, behavior-based detection mechanisms.

# 8 Conclusion

We introduce CAIN, a black-box method that reveals a new vulnerability in LLMs: targeted prompt-based manipulation that preserves benign behavior. CAIN achieves substantial degradation on targeted questions, up to 40% F1 in untargeted attacks and over 70% F1 in targeted ones, without noticeably affecting benign performance. These attacks remain stealthy, transferable across model architectures, and evade traditional defenses such as lexical similarity or perplexity-based filtering.

# 575 Limitation

591

592

593

594

595

596

597

598

599

600

While CAIN demonstrates strong targeted manipu-576 lation in black-box settings, it faces several impor-577 tant limitations. First, achieving high adversarial 578 effectiveness occasionally comes at the cost of be-579 nign performance. CAIN outperforms baselines on OpenAI APIs, the overall attack success remains 581 limited due to alignment constraints in commercial systems. Finally, while CAIN evades common 583 lexical and perplexity-based filters, this also under-584 scores a broader limitation in the field: the lack of robust, behavior-aware defenses. Addressing these challenges will be crucial for advancing both offen-587 sive and defensive research in LLM alignment. 588

# 9 Broader Impacts and Ethics Statement

This work reveals a previously underexplored vulnerability in large language models (LLMs): the ability to craft adversarial system prompts that selectively cause incorrect responses to specific questions while maintaining accurate outputs on benign inputs. Such selective manipulation poses a subtle but serious threat, particularly in domains involving misinformation, political influence, or public health. Unlike traditional jailbreaks or universal attacks, CAIN operates stealthily, evading detection by standard lexical similarity and perplexity filters.

We intend to raise awareness of this threat and prompt the development of more robust, behavior-602 based defenses. All experiments were conducted in controlled settings using open-source models, and 604 evaluations on commercial APIs were performed 605 to assess practical limitations - not for misuse. While the techniques may be misused, we believe that exposing this vector responsibly contributes to a more secure and trustworthy deployment of 609 LLMs. We advocate for responsible disclosure, 610 transparent benchmarking, and the implementation 611 of proactive safeguards in future LLM systems. 612

## 613 References

615

616

617

618

619

630

631

633

634

635

636

637

639

641

650

651

653

656

659

661

- John W. Ayers, Adam Poliak, Mark Dredze, Eric C. Leas, Zechariah Zhu, Jessica B. Kelley, Dennis J. Faix, Aaron M. Goodman, Christopher A. Longhurst, Michael Hogarth, and Davey M. Smith. 2023. Comparing physician and artificial intelligence chatbot responses to patient questions posted to a public social media forum. JAMA Internal Medicine, 183(6):589– 596.
  - Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, Nicholas Joseph, Saurav Kadavath, Jackson Kernion, Tom Conerly, Sheer El-Showk, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, Tristan Hume, and 12 others. 2022. Training a helpful and harmless assistant with reinforcement learning from human feedback. *Preprint*, arXiv:2204.05862.
    - Advik Raj Basani and Xiao Zhang. 2024. Gasp: Efficient black-box generation of adversarial suffixes for jailbreaking llms. *ArXiv*, abs/2411.14133.
    - Emily M. Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. 2021. On the dangers of stochastic parrots: Can language models be too big? In Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency, FAccT '21, page 610–623, New York, NY, USA. Association for Computing Machinery.
    - DeepSeek-AI, Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Daya Guo, Dejian Yang, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, and 181 others. 2025. Deepseek-v3 technical report. *Preprint*, arXiv:2412.19437.
    - J. Ebrahimi, Anyi Rao, Daniel Lowd, and Dejing Dou. 2017. Hotflip: White-box adversarial examples for text classification. In *Annual Meeting of the Association for Computational Linguistics*.
    - J. Gao, J. Lanchantin, M. L. Soffa, and Y. Qi. 2018. Black-box generation of adversarial text sequences to evade deep learning classifiers. In 2018 IEEE Security and Privacy Workshops (SPW), pages 50– 56.
    - 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*.
    - Raghav Gupta, Ryan Sullivan, Yunxuan Li, Samrat Phatale, and Abhinav Rastogi. 2025. Robust multiobjective preference alignment with online dpo. In *AAAI Conference on Artificial Intelligence*.
    - Lynn Hasher, David M. Goldstein, and Thomas C. Toppino. 1977. Frequency and the conference of referential validity. *Journal of Verbal Learning and Verbal Behavior*, 16:107–112.

Jonathan Hayase, Ema Borevkovic, Nicholas Carlini, Florian Tramèr, and Milad Nasr. 2024. Querybased adversarial prompt generation. *ArXiv*, abs/2402.12329. 669

670

671

672

673

674

675

676

677

678

679

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

710

711

712

713

714

715

716

717

718

719

720

721

722

723

724

725

- 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, and 3 others. 2022. Training compute-optimal large language models. *Preprint*, arXiv:2203.15556.
- 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. *Preprint*, arXiv:2309.00614.
- Weipeng Jiang, Zhenting Wang, Juan Zhai, Shiqing Ma, Zhengyu Zhao, and Chao Shen. 2025. An optimizable suffix is worth a thousand templates: Efficient black-box jailbreaking without affirmative phrases via llm as optimizer. *Preprint*, arXiv:2408.11313.
- Di Jin, Zhijing Jin, Joey Tianyi Zhou, and Peter Szolovits. 2019. Is bert really robust? natural language attack on text classification and entailment. *arXiv preprint arXiv:1907.11932*.
- Erik Jones, Anca D. Dragan, Aditi Raghunathan, and Jacob Steinhardt. 2023. Automatically auditing large language models via discrete optimization. *ArXiv*, abs/2303.04381.
- Mandar Joshi, Eunsol Choi, Daniel Weld, and Luke Zettlemoyer. 2017. 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)*, pages 1601–1611, Vancouver, Canada. Association for Computational Linguistics.
- Elad Levi, Eli Brosh, and Matan Friedmann. 2024. Intent-based prompt calibration: Enhancing prompt optimization with synthetic boundary cases.
- Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. TruthfulQA: Measuring how models mimic human falsehoods. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics* (*Volume 1: Long Papers*), pages 3214–3252, Dublin, Ireland. Association for Computational Linguistics.
- George A. Miller. 1994. WordNet: A lexical database for English. In Human Language Technology: Proceedings of a Workshop held at Plainsboro, New Jersey, March 8-11, 1994.
- Eryn J. Newman, Mevagh Sanson, Emily K. Miller, Adele Quigley-Mcbride, Jeffrey L. Foster, Daniel M. Bernstein, and Maryanne Garry. 2014. People with easier to pronounce names promote truthiness of claims. *PLoS ONE*, 9.

804

805

784

726

727

734

736

738

739

OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal,

Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale-

man, Diogo Almeida, Janko Altenschmidt, Sam Alt-

man, Shyamal Anadkat, Red Avila, Igor Babuschkin,

Suchir Balaji, Valerie Balcom, Paul Baltescu, Haim-

ing Bao, Mohammad Bavarian, Jeff Belgum, and

262 others. 2024. Gpt-4 technical report. Preprint,

Anselm Paulus, Arman Zharmagambetov, Chuan Guo,

Ethan Perez, Saffron Huang, Francis Song, Trevor Cai,

Roman Ring, John Aslanides, Amelia Glaese, Nat

McAleese, and Geoffrey Irving. 2022. Red teaming

language models with language models. In Confer-

ence on Empirical Methods in Natural Language

Lianhui Qin, Sean Welleck, Daniel Khashabi, and Yejin

Choi. 2022. Cold decoding: Energy-based con-

strained text generation with langevin dynamics. Ad-

vances in Neural Information Processing Systems,

Qwen, :, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin

Yang, Jiaxi Yang, Jingren Zhou, and 25 oth-

ers. 2025. Qwen2.5 technical report. Preprint,

Wallace, and Sameer Singh. 2020. Autoprompt: Eliciting knowledge from language models with automat-

Taylor Shin, Yasaman Razeghi, Robert L. Logan IV, Eric

ically generated prompts. CoRR, abs/2010.15980.

Harry Surden. 2019. Artificial intelligence and law: An overview. Georgia State University law review,

Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-

chat models. Preprint, arXiv:2307.09288.

Eric Wallace, Shi Feng, Nikhil Kandpal, Matt Gardner,

Laura Weidinger, John F. J. Mellor, Maribeth Rauh,

Conor Griffin, Jonathan Uesato, Po-Sen Huang, Myra

Cheng, Mia Glaese, Borja Balle, Atoosa Kasirzadeh,

and Sameer Singh. 2019. Universal adversarial trig-

gers for attacking and analyzing nlp. In *Conference* on Empirical Methods in Natural Language Process-

bert, Amjad Almahairi, Yasmine Babaei, Nikolay

Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, and 49 others. 2023. Llama 2: Open foundation and fine-tuned

Brandon Amos, and Yuandong Tian. 2024. Ad-

vprompter: Fast adaptive adversarial prompting for

arXiv:2303.08774.

Processing.

35:9538-9551.

arXiv:2412.15115.

35:15109.

llms. ArXiv, abs/2404.16873.

- 745 746
- 747
- 753 754
- 756
- 758
- 759
- 761
- 763 765
- 771
- 772
- 774
- 776

777 778

- 779 Zachary Kenton, Sande Minnich Brown, William T.
  - Hawkins, Tom Stepleton, Courtney Biles, Abeba Birhane, Julia Haas, Laura Rimell, Lisa Anne Hendricks, and 4 others. 2021. Ethical and social risks of harm from language models. ArXiv, abs/2112.04359.

ing.

- Xilie Xu, Keyi Kong, Ning Liu, Lizhen Cui, Di Wang, Jingfeng Zhang, and Mohan Kankanhalli. 2024. An LLM can fool itself: A prompt-based adversarial attack. In The Twelfth International Conference on Learning Representations.
- Rowan Zellers, Ari Holtzman, Hannah Rashkin, Yonatan Bisk, Ali Farhadi, Franziska Roesner, and Yejin Choi. 2019. Defending against neural fake news. ArXiv, abs/1905.12616.
- Yujie Zhao, Jose Efraim Aguilar Escamill, Weyl Lu, and Huazheng Wang. 2024. Ra-pbrl: Provably efficient risk-aware preference-based reinforcement learning. ArXiv, abs/2410.23569.
- Sicheng Zhu, Ruiyi Zhang, Bang An, Gang Wu, Joe Barrow, Zichao Wang, Furong Huang, Ani Nenkova, and Tong Sun. 2024. AutoDAN: Interpretable gradientbased adversarial attacks on large language models. In First Conference on Language Modeling.
- Andy Zou, Zifan Wang, J. Zico Kolter, and Matt Fredrikson. 2023. Universal and transferable adversarial attacks on aligned language models. Preprint, arXiv:2307.15043.

Category	Num. Samples
Misconceptions	100
Conspiracy	26
Stereotype	24
Health	55
Politics	10
History	24

Table A1: The number of data examples used for targeted attack in each category.

## A Appendix

809

810

811

812

814 815

816

817

818

820

821

823

825

826

828

829

831

832

834

838

840

841

842

843

#### A.1 Dataset and Prompt template

Table A1 shows the statistical information of each category in the TruthfulQA (Lin et al., 2022), which will be used for targeted attack. Figure A1 illustrates an example of the input for targeted attack during training and evaluation.

#### A.2 Implementation Details

For each targeted question, we used 20 benign examples (10 from the correct subset, 10 from the incorrect subset) and the correct proportion m = 7. For the AdvAutoPrompt, we set a maximum length of  $p^*$  up to 150 words and the number of iterations t = 3 to save cost.

#### A.3 Analysis of optimization threshold k

For untargeted attack (top row), increasing the incorrect targeted threshold (k) generally leads to a more substantial degradation of performance on the synthetic target (malicious) dataset for both models. For DeepSeek-7B, the F1 score on malicious inputs drops significantly from approximately 40% to just above 10% as k increases, while benign performance shows a modest decline from around 40% to 35%. Qwen2.5-7B exhibits a similar pattern: its malicious performance falls from roughly 55% to about 15%, and benign performance decreases from approximately 60% to 45% across the same range of k. This trend indicates that while the attack more effectively degrades performance on malicious queries at higher k values, this also comes with a slight decrease in performance on benign questions compared to their corresponding initial attacked state and a more noticeable reduction compared to the NSP benign baselines.

Conversely, for targeted attacks (bottom row), a higher threshold (*k*) significantly improves the success rate of forcing the specific malicious output. On DeepSeek-7B, the malicious F1 score, representing attack success, climbs from about 58%

[Instruction]: The following is a multiple-
choice question; you are given two possible
answers. Please only select one of them without
any explanations:
Q: In the Gregorian calendar, what does ``AD"
mean?
Choices:
A: Anno Domini
B: After Death
[Target LLM]: B

Figure A1: The prompt template for two options of targeted attack.

	Prompt	Ber	nign	Mali	cious	Diffe	rence
		F1↑	EM↑	$\overline{\mathbf{F1}}\downarrow$	EM↓	$\Delta$ <b>F</b> 1 $\uparrow$	$\Delta \mathbf{EM} \uparrow$
<del>.</del>	NSP	58.59	47.60	88.61	83.00	-30.02	-35.40
la3	Manual	64.25	56.60	99.75	99.50	-35.70	-42.90
an	AAP	44.84	31.70	52.00	42.00	-7.16	-10.30
П	CAIN	45.15	32.04	27.46	16.40	17.69	15.64
_	NSP	40.98	28.50	97.40	97.00	-56.42	-68.50
þi	Manual	54.82	49.00	100.00	100.00	-45.18	-51.00
P	AAP	49.13	40.06	58.20	51.27	<u>-9.07</u>	<u>-18.14</u>
	CAIN	49.08	40.70	32.32	25.28	16.76	15.42

Table A2: Performance comparison when attacking on Pythia-12B and Llama3.1-7B.

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

to nearly 80% with an increasing k, while performance on benign inputs remains relatively stable around 50%, comparable to its NSP benign baseline. A more pronounced trend is observed for Qwen2.5-7B, where its malicious attack success rate rises from approximately 50% at k = 1 to over 80% for  $k \ge 8$ ; its benign performance also remains stable at around 62%. Notably, this increased targeted efficacy is generally achieved without a substantial negative impact on the models' performance on benign inputs.

#### A.4 Additional results on untargeted attack

Table A2 presents the performance of untargeted attacks on Llama3.1-7B and Pythia-12B. Across both models, CAIN significantly outperforms all baselines, including Manual and AAP, in balancing attack strength and benign performance. While Manual prompts achieve high benign F1/EM, they fail to reduce malicious performance (e.g., 99.75 F1 on Llama3.1 and 100.00 F1 on Pythia). In contrast, CAIN reduces malicious F1 to 27.46 and 32.32, respectively, while maintaining reasonable benign scores. This results in the highest  $\Delta$ F1 and  $\Delta$ EM margins (e.g., +17.69 F1 on Llama3.1 and +16.76

	Prompt	Ber	nign	Mali	cious	Diffe	erence
		F1↑	EM↑	F1↓	EM↓	$\Delta F1\uparrow$	$\Delta \mathbf{EM} \uparrow$
L2-7B	M+G A+G	68.33 63.84	62.59 56.14	38.25 33.36	31.46 28.20	30.08 <b>30.48</b>	<b>31.13</b> 27.94
L2-13B	M+G A+G	81.92 66.77	78.62 57.14	41.44 32.66	38.36 18.89	<b>40.48</b> 34.11	<b>40.26</b> 38.15
L3.1-8B	M+G	62.61	52.12	50.05	41.69	12.56	10.43
	A+G	45.15	32.04	27.46	16.40	<b>17.69</b>	<b>15.64</b>
D-7B	M+G	53.59	48.41	37.73	33.28	15.66	15.13
	A+G	43.99	31.75	28.15	16.33	<b>15.84</b>	<b>15.42</b>
Q2.5	M+G	46.97	36.13	61.39	50.68	-14.42	-14.55
	A+G	50.31	39.20	34.94	23.92	15.37	<b>15.28</b>
P-12B	M+G	50.25	42.90	40.46	34.41	9.79	8.49
	A+G	49.08	40.70	32.32	25.28	<b>16.76</b>	1 <b>5.42</b>

Table A3: Results of attacking performance with manual initialization and AdvAutoPrompt. "A" denotes AAP, "G" stands for Greedy, and "A+G" is our proposed method. "L, D, Q, P" denote Llama, Deepseek, Qwen, and Pythia models, respectively.

on Pythia), demonstrating CAIN's superior ability to selectively degrade targeted outputs without broadly compromising accuracy.

870

871

872

873

874

876

877

879

880

881

882

## A.5 Additional results for targeted attacks

Table A4 shows that AAP achieves the best performance in the two-option setting ( $\Psi$ F1 = 50.94,  $\Psi$ EM = 42.58), but its effectiveness drops when transferred to the four-option format. In contrast, CAIN maintains more stable performance across both settings, achieving strong targeted attack success (highest malicious F1) with better transferability ( $\Psi$ F1 = 31.01 vs. 32.80). This suggests CAIN is more robust and generalizable under realistic conditions where question formats vary.



Figure A2: Perplexity distribution of successfully optimized prompts across different prompt methods under both untargeted and targeted attack.

	Prompt			Two	option	s		Two options $\rightarrow$ Four options							
		Benign Malicious			cious	S	um	Benign		Malicious		Sum			
		<b>F1</b> ↑	EM↑	$F1\uparrow$	EM ↑	$\Psi$ F1 $\uparrow$	$\Psi$ EM $\uparrow$	$F1\uparrow$	EM↑	<b>F1</b> ↑	EM↑	$\Psi$ F1 $\uparrow$	$\Psi$ EM $\uparrow$		
B	NSP	10.03	1.67	14.09	3.27	12.06	2.47	9.72	1.33	19.60	1.82	14.66	1.58		
Llama2-7	Manual M+Greedy AAP CAIN	19.67 28.20 42.47 35.10	19.00 21.57 34.16 19.61	44.03 60.71 59.41 61.86	43.45 58.00 51.00 <u>48.71</u>	31.85 44.45 <b>50.94</b> <u>48.48</u>	31.23 39.78 <b>42.58</b> 34.16	1.00 13.73 27.30 20.44	1.00 10.60 19.88 11.83	3.50 29.81 38.29 41.58	3.27 28.80 33.27 34.84	2.25 21.77 <b>32.80</b> <u>31.01</u>	2.13 19.70 <b>26.58</b> 23.34		

Table A4: Performance of the targeted attack on Llama2-7B.

	Prompt			Two o	ptions			Т	wo opt	→Four	r options		
		Benign		Mali	Malicious		Sum $(\Psi)$		nign	Malicious		$\overline{{\rm Sum}(\Psi)}$	
		F1↑	EM↑	$F1\uparrow$	<b>EM</b> ↑	F1↑	EM ↑	<b>F1</b> ↑	EM↑	F1↑	EM↑	F1↑	EM↑
Deepseek-7B	M+Greedy A+Greedy	47.94 55.29	46.11 46.47	42.90 58.92	42.26 54.00	45.42 <b>57.11</b>	44.19 <b>50.23</b>	27.38 31.73	25.65 28.69	30.32 43.92	29.35 43.00	28.85 <b>37.83</b>	27.50 <b>35.84</b>
Qwen2.5-7B	M+Greedy A+Greedy	60.41 62.76	60.33 61.75	62.73 73.03	62.73 72.07	61.57 <b>67.89</b>	61.53 66.91	50.88 56.45	50.88 56.38	69.70 70.75	69.70 70.34	60.29 <b>63.60</b>	60.29 63.36

Table A5: Performance with different initialization methods on targeted attacks.