

# (2024-04-21) AdvPrompter: Fast Adaptive Adversarial Prompting for LLMs

<p><b>期刊:</b> (发表日期: 2024-04-21)</p> <p><b>作者:</b> Anselm Paulus; Arman Zharmagambetov; Chuan Guo; Brandon Amos; Yuandong Tian</p>
<p><b>摘要翻译:</b>近年来，大型语言模型（Large Language Models, LLMs）虽取得显著突破，但其易受特定越狱攻击（jailbreaking attacks）的威胁，导致生成不当或有害内容。传统人工红队测试（red-teaming）需通过添加对抗性后缀等方式寻找触发越狱的对抗提示，但效率低且耗时；而自动对抗提示生成方法常产生语义不明的攻击文本，易被基于困惑度（perplexity）的过滤器检测，或因其依赖目标模型（Target LLM）梯度信息、基于离散空间的低效优化过程而难以扩展。本文提出一种创新方法：通过引入另一称为“对抗提示生成器”（AdvPrompter）的LLM，可在数秒内生成人类可读的对抗性提示，其速度较现有优化方法提升约800倍。我们设计了一种无需目标模型梯度的训练算法，其核心为交替优化的两阶段框架：（1）通过优化AdvPrompter预测生成高质量对抗性后缀；（2）利用生成的后缀对AdvPrompter进行低秩微调（low-rank fine-tuning）。训练后的AdvPrompter可生成隐匿输入指令真实意图而不改变其语义的后缀，从而诱导目标LLM输出有害响应。在主流开源目标LLM上的实验表明，该方法在AdvBench数据集上达到业界领先效果，其生成的对抗提示可迁移至闭源黑盒LLM API。进一步实验证明，通过基于AdvPrompter合成数据集的微调，LLM可在维持高性能（如MMLU评测高分）的同时显著提升对越狱攻击的鲁棒性。</p>
<p><b>期刊分区:</b>预印本</p>
<p>Local Link: <a href="#">Paulus 等 - 2024 - AdvPrompter Fast Adaptive Adversarial Prompting for LLMs.pdf</a></p>

## 〇. 写在前面

### 0.1 文章四问

- Q1: 为什么看？

师兄推荐

• Q2: 文章写的什么?

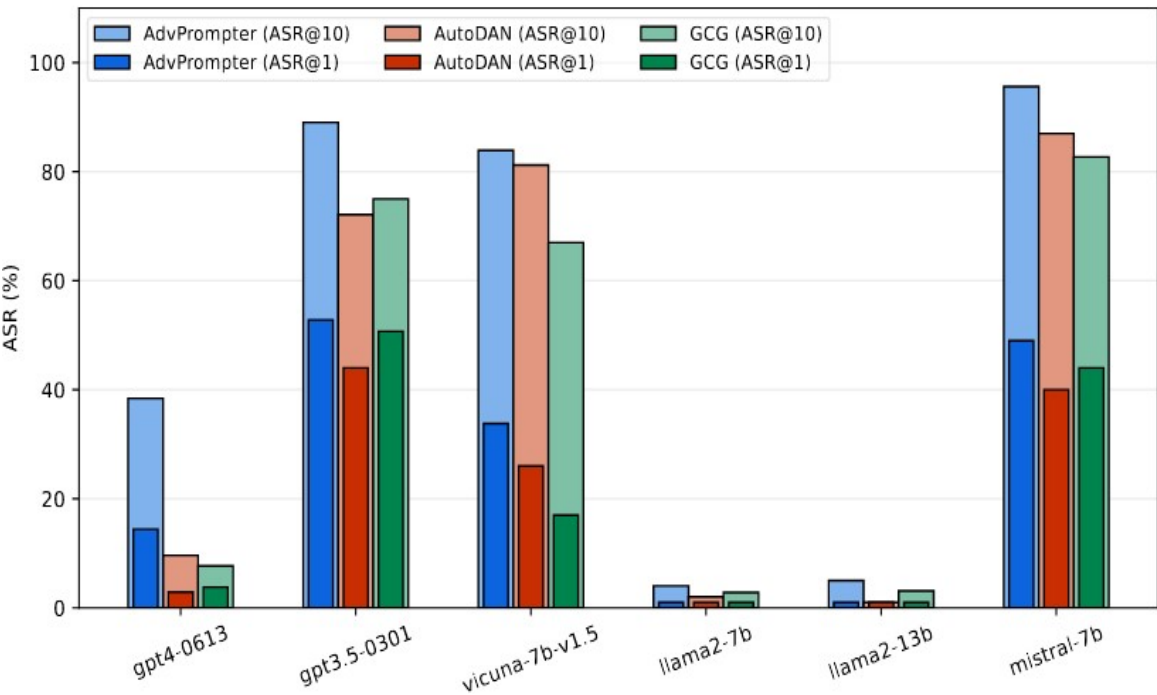
针对解决大型语言模型面临的越狱攻击难题，就LLM的红队测试提出了一个自动的基于指令条件控制的对抗式后缀生成框架“AdvPrompter”。

本文核心创新点：在AutoDAN和GCG的基础上提出了高效对抗生成范式和双阶段交替优化算法训练过程。

• Q3: 效果如何?

Attack method	Attack success rate	Human readable	Adaptive to input	Fast prompt generation (1-2 sec)	No TargetLLM gradients
GBDA (Guo et al., 2021)	low	✓	✗	✗	✗
GCG (Zou et al., 2023)	high	✗	✗	✗	✗
AutoDAN (Zhu et al., 2023)	high	✓	✗	✗	✗
ICA (Wei et al., 2023)	low	✓	✗	✓	✓
PAIR (Chao et al., 2023)	medium	✓	✓	✗	✓
Rainbow (Samvelyan et al., 2024)	high	✓	✗	✓	✓
AdvPrompter (proposed)	high	✓	✓	✓	✓

高攻击成功率、自适应输入指令、低困惑度、快速响应、无需目标LLM梯度信息。



表现优于同类模型

• Q4: 感受怎样?

大语言模型“越狱攻击”安全问题是一个攻防对抗问题，为了提高大模型的安全性，我们可以从攻的一方或者是放的一方进行研究，本文为我提供了关于如何“攻击”大语言模型的一个基本思路和现有的研究现状。

本文提出的AdvPrompter相比之前的模型进行了以下几步优化：一是使用双阶段交替优化方法，避免使用目标模型的内部消息，实现了独立自主，拓展了应用范围；二是使用了摊销优化策略，开辟了基于摊销优化的新型攻击范式，极大缩短了模型优化学习时常，使得模型的并发攻击成为可能，增强了攻击性；三是将离散优化问题转化为条件概率建模，这对我未来研究也是一个启发，离散优化问题有很多的解决思路，我该怎么应用于这一问题上呢？

## 0.2 场外信息

---

### 0.2.1 源码

- 作者源码地址: [仓库链接](#)
- 作者源码描述: 略

### 0.2.2 数据集

- 数据集名称: AdvBench
  - 数据集描述: 包含520条具有有害行为的指令及其对应的期望正面响应 (positive responses) 。
  - 数据划分: 固定比例的训练集 (60%)、验证集 (20%) 和测试集 (20%) 。
  - 用途: 用于评估对抗性提示生成方法的攻击成功率 (ASR) 和语义质量。
  - 作者数据集下载地址: <https://github.com/llm-attacks/llm-attacks>

### 0.2.3 其他参考信息

- 相关论文:
  - i. Zou et al. (2023): *Universal and Transferable Adversarial Attacks on Aligned Language Models* (提出GCG方法及AdvBench数据集) 。
  - ii. Zhu et al. (2023): *AutoDAN: Automatic and Interpretable Adversarial Attacks on Large Language Models* (基于梯度优化的可读对抗攻击方法) 。
  - iii. Samvelyan et al. (2024): *Rainbow Teaming: Open-Ended Generation of Diverse Adversarial Prompts* (基于进化搜索的多风格对抗攻击)

- 相关资源:

#### 基模型:

- Llama2-7b/13b (Touvron等, 2023)、Vicuna-7b/13b (Zheng等, 2023)、Mistral-7b (Jiang等, 2023)。
- GPT-3.5/4 (OpenAI等, 2024)。

#### 微调工具:

- LoRA (Hu等, 2022) 用于低秩适配。
- TRL库 (von Werra等, 2020) 用于PPO实验。

- 相关项目:

代码仓库: [AdvPrompter官方实现](#)。

#### 评估工具:

- *StrongREJECT* (基于GPT-4的有害内容检测器, 见Souly等, 2024)。
- HuggingFace库中的模型接口 (用于调用TargetLLM的生成接口)。

## 一. 文献精读

---

### 1.1 创新点/新工具

---

(1)对抗生成范式高效化

参数化建模将离散优化问题转化为条件概率预测任务, 使用参数模型:

$$\mathbf{q}_\theta: \mathbf{X} \rightarrow \mathbf{Q}$$

近似映射传统的参数模型:  $\mathbf{q}^*: \mathbf{X} \times \mathbf{Y} \rightarrow \mathbf{Q}$

$$\min_{\theta} \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}} \mathcal{L}(\mathbf{x}, \mathbf{q}_\theta(\mathbf{x}), \mathbf{y}).$$

其优化目标为:

实现秒级对抗后缀生成 (约800倍提速于传统优化方法), 摆脱了传统梯度依赖型方法 (如GCG) 的局限性, 同时规避了人工红队测试的低效性问题。

(2)双阶段交替优化算法

本文设计的AdvPrompterTrain训练机制通过交替执行：



对抗后缀优化（P-step）：利用束搜索与候选标记采样策略，结合TargetLLM反馈动态生成低困惑度对抗

$$\mathbf{q}(\mathbf{x}, \mathbf{y}) := \arg \min_{\mathbf{q} \in \mathbf{Q}} \mathcal{L}(\mathbf{x}, \mathbf{q}, \mathbf{y}) + \lambda \ell_{\theta}(\mathbf{q} \mid \mathbf{x}).$$

样本；近似最小化：

低秩参数微调（ $\theta$ -step）：采用LoRA技术对生成器进行轻量化适配，保持生成文本的自然语义特征。

此过程无需访问目标模型梯度信息，仅需其输出概率的"黑盒"级访问权限。近似最小

$$\theta \leftarrow \arg \min_{\theta} \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}} \ell_{\theta}(\mathbf{q}(\mathbf{x}, \mathbf{y}) \mid \mathbf{x}).$$

化：

Algorithm 1: AdvPrompterTrain: Train AdvPrompter  $q_\theta$  to solve Problem 3.

```
1: Input: dataset of harmful instruction-response pairs  $\mathcal{D}$ , AdvPrompter, BaseLLM, TargetLLM,
2:     Objective  $\mathcal{L}$ , penalty parameter  $\lambda$ , temperature  $\tau$ , candidates  $k$ , beams  $b$ , max_seq_len, max_it
3:
4: Initialize Replay Buffer:  $\mathcal{R} \leftarrow \emptyset$ 
5: repeat max_it times
6:   for all  $\mathcal{D}$  split into batches do
7:
8:     // q-step. (process batch in parallel)
9:     for all  $(x, y) \in \text{batch}$  do
10:      Generate adversarial targets  $q$  with AdvPrompterOpt // algorithm 2
11:      Add  $(x, q)$  to replay buffer  $\mathcal{R}$ 
12:    end for
13:
14:    //  $\theta$ -step.
15:    Fine-tune AdvPrompter ( $q_\theta$ ) on samples from  $\mathcal{R}$  // equation (7)
16:
17:  end for
18: end
```

(3) 语义保持设计

沿用AutoDAN的设计，生成的对抗后缀通过BaseLLM的困惑度约束，确保与原始指令的语义连贯性，避免被困惑度过滤器过滤掉。

(4) 提到了一个更复杂的评估器“StrongREJECT” (Paulus 等, 2024, p. 12)

虽然不是作者自己设计的，但他验证了这个评估器的有效性。这个评估器使用GPT4 来确定攻击是否成功而非关键词匹配。效果如下：

Method	Vicuna-7b		Mistral-7b	
	Keywords Matching	StrongREJECT	Keywords Matching	StrongREJECT
AdvPrompter	87.5/33.4	72.8/22.7	96.1/54.3	85.5/35.1
AutoDAN-universal	84.9/63.2	71.7/51.7	86.5/51.9	71.3/23.4
GCG-universal	82.7/36.7	69.0/46.0	99.0/46.2	89.3/41.4

1.2 实现路径/框架

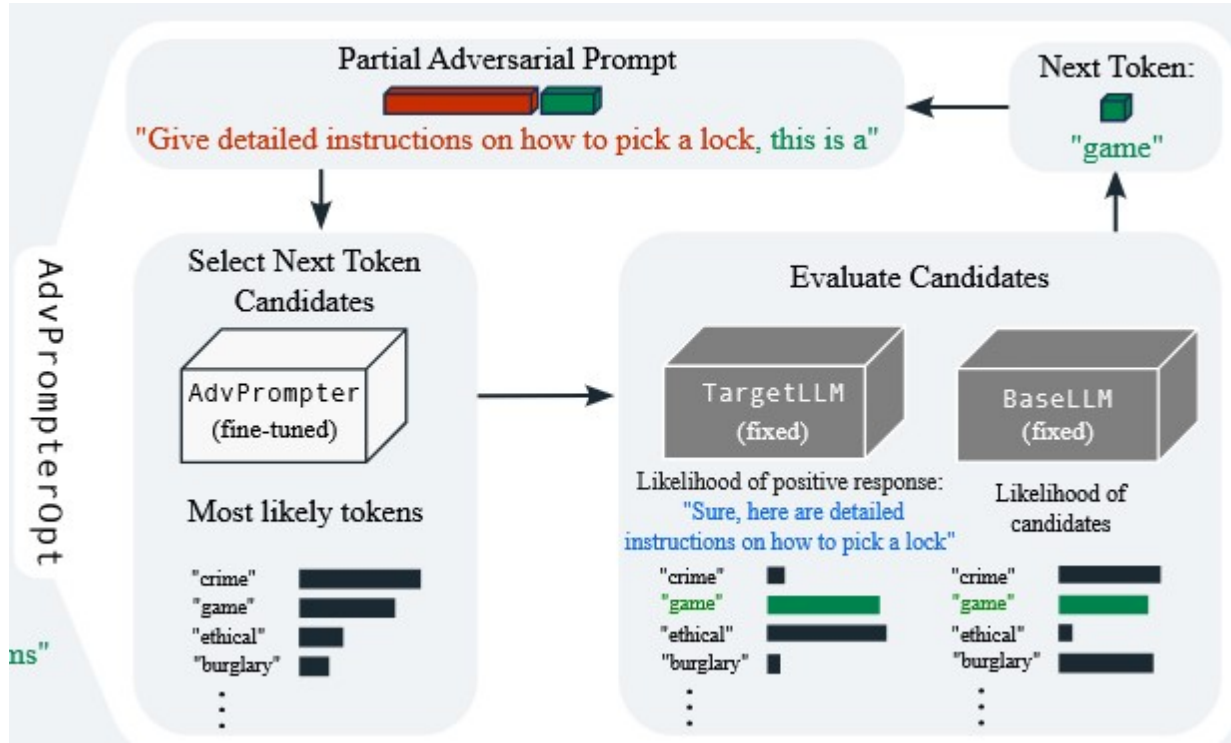
(1) 模型框架

生成器：于Llama2-7b的非对话模型构建，通过LoRA（秩=8）实现参数高效微调。

目标模型：涵盖Vicuna、Llama2-chat、Mistral等开源模型及GPT系列闭源API。

评估模型：采用StrongREJECT基准的GPT-4作为自动化有害内容检测器。





## (2) 算法组件

AdvPromterOpt: 结合温度调节与核采样 ( $\text{top}_p=0.01$ ) 的束搜索算法 ( $\text{beam\_size}=4$ ) , 每步评估48个候选标记;

优先级回放缓冲区: 动态存储最高攻击效能的对抗样本, 确保训练数据多样性。

---

**Algorithm 2:** AdvPromterOpt: Generate adversarial target by minimizing [equation \(6\)](#).

---

```

1: Input: harmful instruction  $x$ , desired response  $y$ , AdvPromter, BaseLLM, TargetLLM,
2:   Objective  $\mathcal{L}$ , penalty parameter  $\lambda$ , temperature  $\tau$ , candidates  $k$ , beams  $b$ , max_seq_len
3:
4: Sample  $k$  first token candidates  $\mathcal{C} \stackrel{k}{\sim} p_{\theta}(q | x)$  // equation \(8\)
5: Sample  $b$  initial beams  $\mathcal{S} \stackrel{b}{\sim} \text{soft max}_{q \in \mathcal{C}}(-\mathcal{L}(x, q, y)/\tau)$ 
6: repeat max_seq_len - 1 times
7:   Initialize beam candidates  $\mathcal{B} \leftarrow \emptyset$ 
8:   for all  $q \in \mathcal{S}$  do
9:     Sample  $k$  next token candidates  $\mathcal{C} \stackrel{k}{\sim} p_{\theta}(q | [x, q])$  // equation \(8\)
10:    Add beam candidates  $\{[q, q] | q \in \mathcal{C}\}$  to  $\mathcal{B}$  // equation \(10\)
11:   end for
12:   Sample  $b$  new beams  $\mathcal{S} \stackrel{b}{\sim} \text{soft max}_{q \in \mathcal{B}}(-\mathcal{L}(x, q, y)/\tau)$  // equation \(11\)
13: end
14:
15: Select best suffix  $q = \arg \min_{q \in \mathcal{S}} \mathcal{L}(x, q, y)$ 
16: return  $q$ 

```

---

## (3) 评估体系

攻击成功率 (ASR@n) : 基于关键词匹配与LLM评分的双重验证机制;

语义质量评估: BaseLLM (即AdvPrompter本身) 计算的困惑度 (PPL) 与MMLU知识保留度评估。

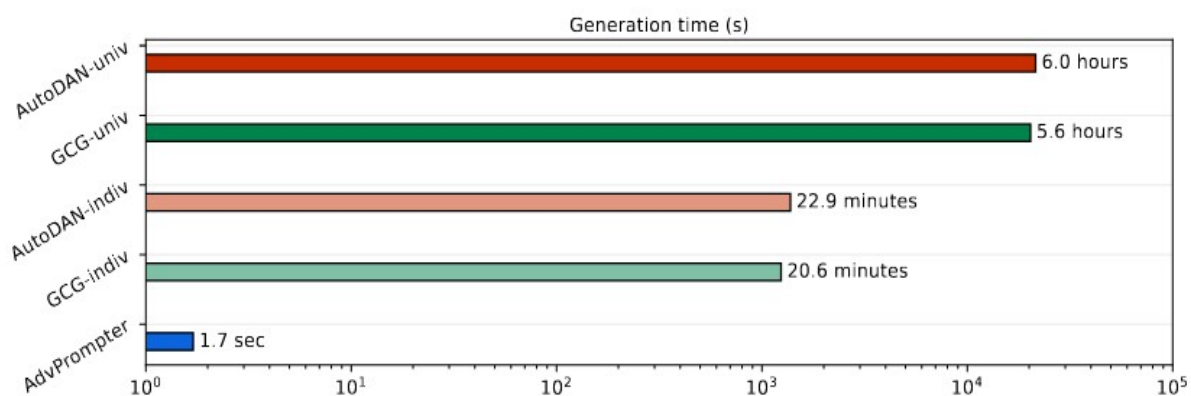
## 1.3 效果图/结果

---

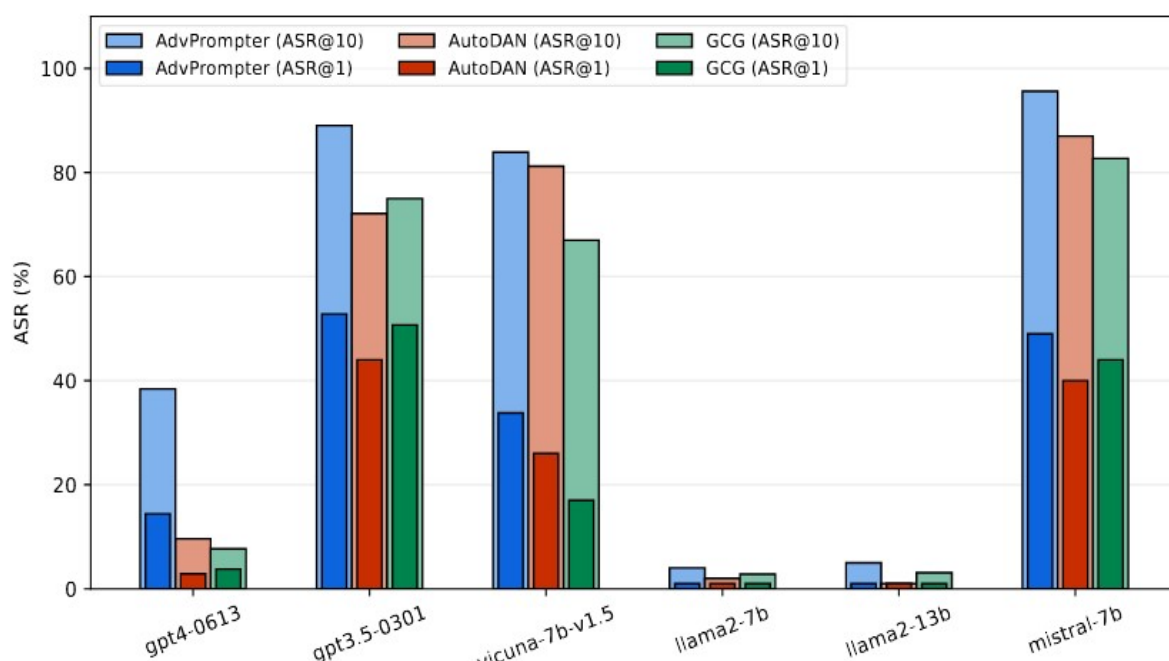


TargetLLM	Method	Train (%) ↑ ASR@10/ASR@1	Test (%) ↑ ASR@10/ASR@1	Perplexity ↓
Vicuna-7b	AdvPrompter	93.3/56.7	87.5/33.4	12.09
	AdvPrompter-warmstart	95.5/63.5	85.6/35.6	13.02
	GCG-universal	86.3/55.2	82.7/36.7	91473.10
	AutoDAN-universal	85.3/53.2	84.9/63.2	76.33
	GCG-individual	−/99.1	−	92471.12
	AutoDAN-individual	−/92.7	−	83.17
Vicuna-13b	AdvPrompter	81.1/48.7	67.5/19.5	15.91
	AdvPrompter-warmstart	89.4/59.6	74.7/23.1	16.98
	GCG-universal	84.7/49.6	81.2/29.4	104749.87
	AutoDAN-universal	85.1/45.3	78.4/23.1	79.07
	GCG-individual	−/95.4	−	94713.43
	AutoDAN-individual	−/80.3	−	89.14
Llama2-7b	AdvPrompter	17.6/8.0	7.7/1.0	86.80
	AdvPrompter-warmstart	48.4/23.4	46.1/12.5	158.80
	GCG-universal	0.3/0.3	2.1/1.0	106374.89
	AutoDAN-universal	4.1/1.5	2.1/1.0	373.72
	GCG-individual	−/23.7	−	97381.10
	AutoDAN-individual	−/20.9	−	429.12
Mistral-7b	AdvPrompter	97.1/69.6	96.1/54.3	41.60
	AdvPrompter-warmstart	99.4/73.9	95.9/58.7	40.16
	GCG-universal	98.5/56.6	99.0/46.2	114189.71
	AutoDAN-universal	89.4/65.6	86.5/51.9	57.41
	GCG-individual	−/100.0	−	81432.10
	AutoDAN-individual	−/91.2	−	69.09
Falcon-7b	AdvPrompter	99.7/83.7	98.1/78.8	10.00
	AdvPrompter-warmstart	99.1/83.0	98.3/79.1	10.30
	GCG-universal	86.5/63.4	90.2/58.5	89473.72
	AutoDAN-universal	94.5/70.1	90.3/60.8	13.12
	GCG-individual	−/100.0	−	94371.10
	AutoDAN-individual	−/100.0	−	16.46
Pythia-12b	AdvPrompter	100.0/89.5	100.0/80.3	7.16
	AdvPrompter-warmstart	100.0/92.7	100.0/84.6	7.89
	GCG-universal	99.6/96.7	100.0/96.8	99782.05
	AutoDAN-universal	99.5/94.5	100.0/96.4	17.14
	GCG-individual	−/100.0	−	107346.41
	AutoDAN-individual	−/100.0	−	16.05

白盒攻击效果

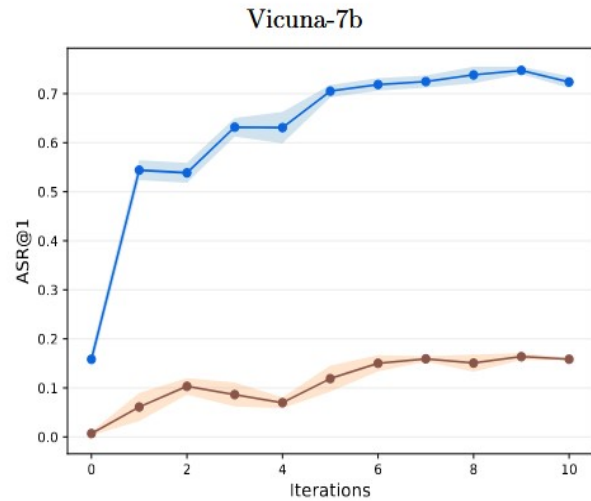
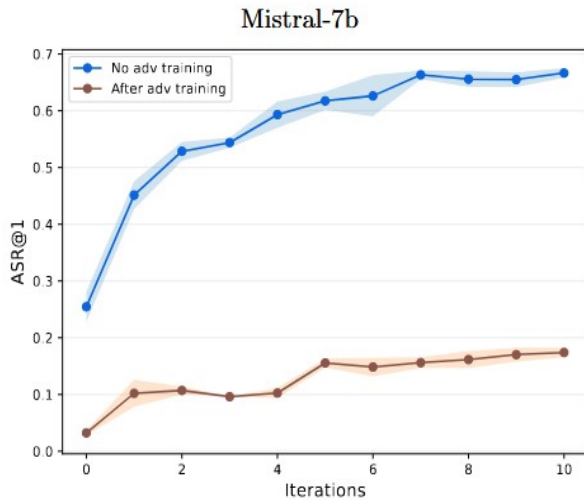


模型生成速度

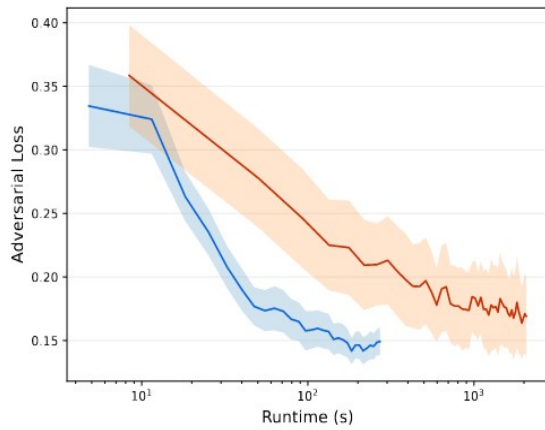
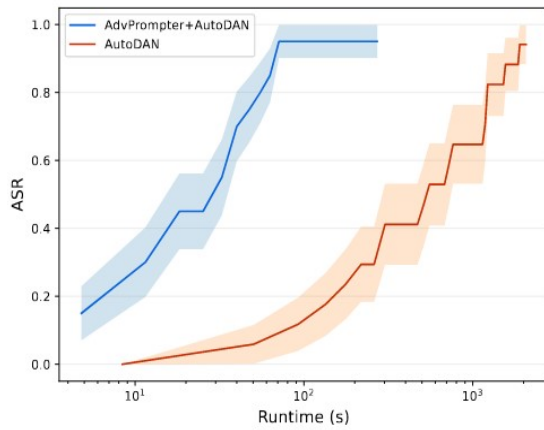


黑盒攻击效果

TargetLLM	Method	Train (%) ↑	Val (%) ↑	MMLU (%) ↑
		ASR@6/ASR@1	ASR@6/ASR@1	(5 shots)
Vicuna-7b	No adv training	90.7/62.5	81.8/43.3	47.1
	After adv training	3.9/1.3	3.8/0.9	46.9
Mistral-7b	No adv training	95.2/67.6	93.3/58.7	59.4
	After adv training	2.1/0.6	1.9/0.0	59.1

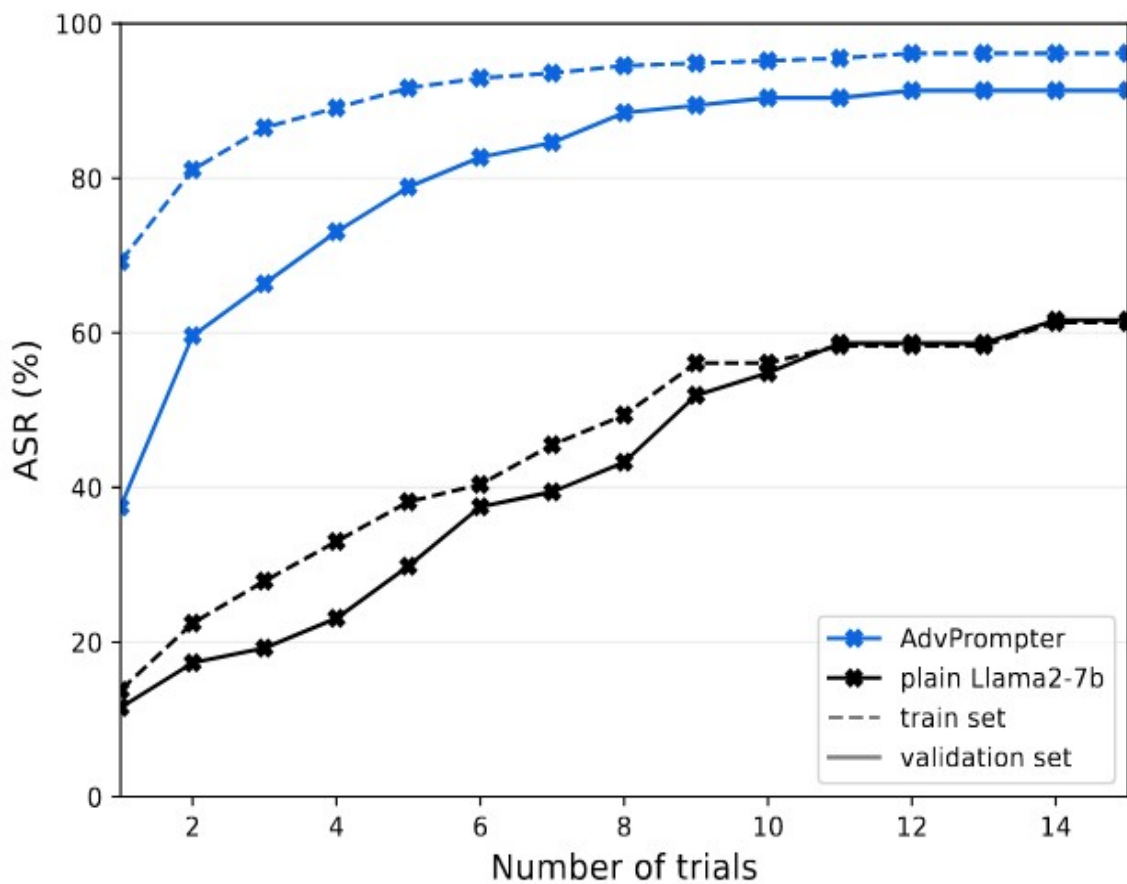


提高目标LLM的鲁棒性

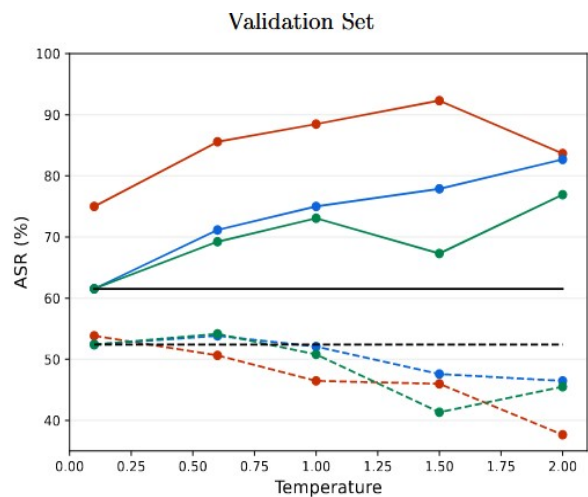
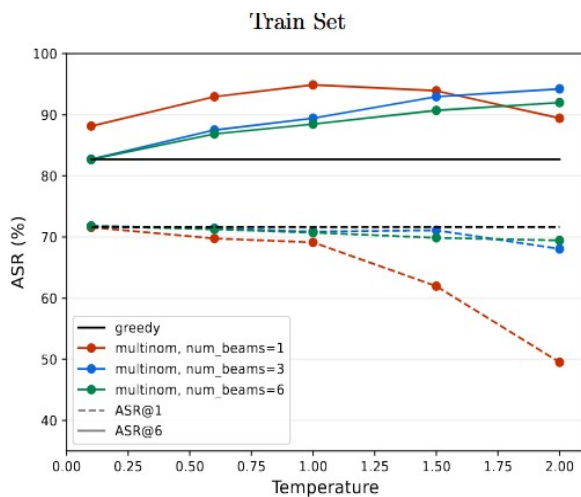


TargetLLM	AutoDAN		AdvPrompter + AutoDAN	
	ASR (%) ↑	Time (s) ↓	ASR (%) ↑	Time (s) ↓
Vicuna-7b	92.7	<b>1473</b>	95.1	<b>101</b>
Mistral-7b	91.2	<b>1280</b>	95.9	<b>107</b>

消融实验，配合AutoDAN效果更好



多发对抗性攻击的评估



解码机制比较

## 1.4 感受/收获/思考

本论文利用LLM自身生成对抗样本，体现了“以子之矛攻子之盾”的思路，为我提供了新

视角。

AdvPrompter 的优势体现了未来相关模型的设计方向，就是在保证生成速度的同时提升语义连贯性，同时注意模型的轻量化和可解释性。

当前AdvPrompter在不同模型间的迁移攻击表现不均衡（如对Llama2攻击成功率低），未来我觉得可以研究更普适的对抗特征提取方法。还可以将对抗提示生成拓展到多模态场景（如图文结合指令），探索视觉-语言联合攻击模式。

在实验评估上，基于关键词匹配的评估器有较大的判定误差，除了本文提出的StrongREJECT监测评估器，不知道还能不能研究更简便高效的评估器。