

# AdvPrompter: Fast Adaptive Adversarial Prompting for LLMs

## 读书笔记

### 一、传统方法有什么缺陷

- 生成内容无意义**：传统方法（如GCG）生成的对抗提示通常语义不连贯，容易被基于困惑度的过滤器检测到。
- 目标模型梯度依赖**：需要获取目标LLM的梯度信息（白盒攻击），难以应用于黑盒模型。
- 计算效率低**：生成单个对抗提示需数小时，且每次攻击需重新优化。
- 适应性缺乏**：生成的对抗提示是固定后缀，无法根据不同输入动态调整。

### 二、本文针对传统方法的缺陷要解决的是什么问题

- 提升提示生成的速度和可读性**：如何快速生成人类可读的对抗性提示，同时保持高效的攻击成功率。
- 提高模型自适应能力**：如何使对抗性提示能够适应不同的指令，即使在未见过的测试指令上也能保持有效。
- 灰盒消除模型梯度依赖**：如何在不需要目标模型梯度信息的情况下，生成有效的对抗性提示，实现白盒访问目标模型。

### 三、本文核心思想和步骤（怎么解决问题的）

#### 3.1 核心思想

训练一个专用对抗生成器（**AdvPrompter**），通过自回归生成人类可读的对抗提示，动态适配不同输入指令，实现高效、隐蔽的越狱攻击，同时利用对抗数据反向增强目标模型的防御能力。

#### 3.2 具体步骤

##### 3.2.1 对抗生成器训练（AdvPrompterTrain）

提出一种新的训练方法AdvPrompterTrain，通过交替执行两个步骤来训练：

- q-step（对抗样本生成）**：基于当前AdvPrompter参数，使用优化算法AdvPrompterOpt（从AdvPrompter的预测分布中采样候选词元，结合贪心策略或随机束搜索选择最优词元，迭代生成完整对抗后缀）生成高攻击成功率的对抗提示。
- θ-step（模型微调）**：使用q-step生成的对抗性后缀作为目标，通过低秩适配（LoRA）微调AdvPrompter，提升其生成能力。

##### 3.3 实时对抗提示生成

- 动态适配**：AdvPrompter接收用户指令后，自回归生成语义连贯的后缀，绕过基于困惑度的过滤器。
- 多轮攻击加速**：训练后的AdvPrompter单次生成快，支持批量生成候选提示，攻击效率大幅提升。

##### 3.4 防御-攻击闭环增强

- 对抗性微调防御**：将AdvPrompter生成的越狱数据用于TargetLLM微调，使其对同类攻击的鲁棒性提升。
- 迁移攻击验证**：在黑盒模型上直接迁移攻击，成功率高。

### 四、最终效果

- 攻击成功率（ASR）更高**：在白盒和黑盒设置下，AdvPrompter的攻击成功率均高于现有方法（GCG和AutoDAN等）。
- 生成速度快**：单次生成仅1-2秒，比传统方法快800倍。
- 可读性优**：困惑度（PPL）显著低于传统方法，更接近人类语言。
- 迁移攻击强**：AdvPrompter能够根据具体指令生成适应性强的对抗性后缀，即使在未见过的指令上也能保持有效。

## 五、创新点

- 无需目标模型梯度信息**：AdvPrompter的训练和应用不依赖于目标模型的梯度信息，使其能够应用于黑盒模型，扩大了对抗性攻击方法的适用范围。
- 动态生成可读对抗提示，适应性强**：AdvPrompter可以生成自然语言对抗提示，突破传统无意义后缀的限制，显著提升隐蔽性（绕过基于困惑度的检测）。
- 自动化交替训练框架**：通过对抗后缀生成（q-step）与模型微调（ $\theta$ -step）交替迭代，实现高效优化，支持训练阶段和推理阶段的灵活适配。

## 六、本文不足和可改进方向

### 6.1 本文不足

- 依赖预训练模型**：AdvPrompter需预先生成对抗样本，可能受限于数据集多样性，性能可能受到预训练模型质量的影响。
- 防御验证局限**：仅测试了基础防御机制（如安全系统消息），未覆盖高级防御（如输入/输出端的多模型联合过滤）。

### 6.2 可改进方向

- 提高对抗性提示的隐蔽性和适应性**：研究如何使对抗性提示更难被检测到，同时保持对不同模型和未知指令的高适应性。
- 机制可解释性分析**：深入解析AdvPrompter如何提升TargetLLM抗攻击能力，量化对抗提示与模型鲁棒性提升的因果关联（如关键语义模式、扰动阈值）。