# EasyJailbreak: A Unified Framework for Jailbreaking Large Language Models

**Authored by:**

Weikang Zhou  Xiao Wang  Limao Xiong  Han Xia  Yingshuang Gu  Mingxu Chai  Fukang Zhu  Caishuang Huang  Tao Gui  Qi Zhang  Xuanjing Huang

**Presented by:**

Himadri Gobinda Biswas-2105047
Shams Hossain Simanto-2105048
Nawriz Ahmed Turjo-2105032

**1** Introduction

**2** Related Work

**3** Framework

**4** Usage

**5** LLM Benchmarking

**6** Conclusion

# 1 Introduction

2 Related Work

3 Framework

4 Usage

5 LLM Benchmarking

6 Conclusion

# Introduction to Jailbreaking

## What is jailbreaking?

# Introduction to Jailbreaking

## What is jailbreaking?

*A method to bypass safeguards in Large Language Models (LLMs) to elicit prohibited or unintended outputs.*

## Introduction to Jailbreaking

# Jailbreaking as a research topic



Objectives

# Jailbreaking as a research topic

# Jailbreaking as a research topic



Identify vulnerabilities in LLMs

Improve model safety and robustness

Objectives

Introduction
○○○●

Related Work
○○○○

Framework
○○○○○○○○○○○○○

Usage
○○○○

LLM Benchmarking
○○○○○○○○○○○○○○

Conclusion
○○○○○

# Jailbreaking as a research topic

Introduction
○○○●

Related Work
○○○○

Framework
○○○○○○○○○○○○

Usage
○○○○

LLM Benchmarking
○○○○○○○○○○○○○○○

Conclusion
○○○○○

# Jailbreaking as a research topic

## Related Work

Current jailbreaking methodologies fall into 3 categories.

Current
Methods

## Related Work

Current jailbreaking methodologies fall into 3 categories.

Introduction
0000

**Related Work**
0●00

Framework
0000000000000

Usage
0000

LLM Benchmarking
000000000000000
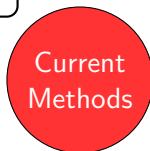
Conclusion
00000

## Related Work

Current jailbreaking methodologies fall into 3 categories.

Introduction
oooo
**Related Work**
o●oo
Framework
ooooooooooooo
Usage
oooo
LLM Benchmarking
ooooooooooooo
Conclusion
ooooo

## Related Work

Current jailbreaking methodologies fall into 3 categories.

# Limitations of Existing Approaches

## Limitations of Existing Approaches

- Fair comparison is hard due to varying datasets and models.

Introduction
0000

**Related Work**
0000

Framework
0000000000000

Usage
0000

LLM Benchmarking
00000000000000

Conclusion
00000

# Limitations of Existing Approaches

- Fair comparison is hard due to varying datasets and models.
- Lack of source code makes reproducing prior work slow and error-prone.
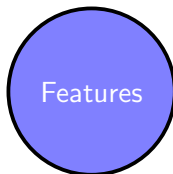
## Limitations of Existing Approaches

- Fair comparison is hard due to varying datasets and models.
- Lack of source code makes reproducing prior work slow and error-prone.
- These barriers complicate identifying and addressing LLM vulnerabilities.

# Features of EasyJailbreak

EasyJailbreak addresses these limitations by offering the following features:

Features

## Features of EasyJailbreak

EasyJailbreak addresses these limitations by offering the following features:

Introduction
0000

**Related Work**
000●

Framework
0000000000000

Usage
0000

LLM Benchmarking
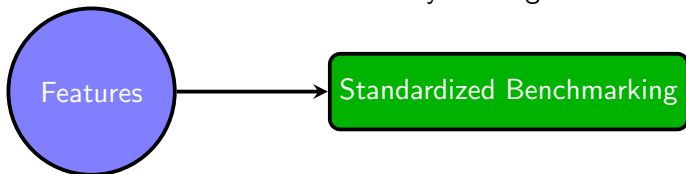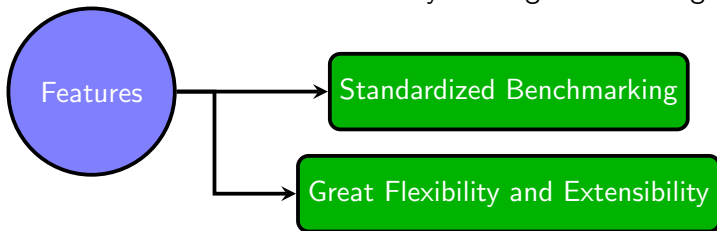00000000000000

Conclusion
00000

## Features of EasyJailbreak

EasyJailbreak addresses these limitations by offering the following features:

## Features of EasyJailbreak

EasyJailbreak addresses these limitations by offering the following features:

Introduction
oooo

Related Work
oooo

**Framework**
○●○○○○○○○○○○○○

Usage
oooo

LLM Benchmarking
○○○○○○○○○○○○○○○

Conclusion
ooooo

# Framework diagram

Introduction
oooo

Related Work
oooo

Framework
oo●ooooooooooooo

Usage
oooo

LLM Benchmarking
oooooooooooooooo

Conclusion
ooooo

# Steps To Conduct A Jailbreak Attack

# Preparation

**What do we even do in the preparation phase?**

Introduction
oooo

Related Work
oooo

**Framework**
oooo●oooooooooo

Usage
oooo

LLM Benchmarking
oooooooooooooooo

Conclusion
ooooo

# Preparation

**What do we even do in the preparation phase?**

**Define the "Queries"**

Introduction
oooo

Related Work
oooo

**Framework**
ooo●ooooooooooo

Usage
oooo

LLM Benchmarking
oooooooooooooooo

Conclusion
ooooo

# Preparation

> What do we even do in the preparation phase?

**Define the "Queries"**

**Choose the "Seeds"**

Introduction
oooo

Related Work
oooo

**Framework**
ooo●ooooooooooo

Usage
oooo

LLM Benchmarking
ooooooooooooooooo

Conclusion
ooooo

# Preparation

**What do we even do in the preparation phase?**

**Define the "Queries"**

**Select the "Models"**

**Choose the "Seeds"**

Introduction
0000

Related Work
0000

**Framework**
0000●0000000000

Usage
0000

LLM Benchmarking
000000000000000

Conclusion
00000

## Preparation

**What do we even do in the preparation phase?**

**Define the "Queries"**

**Select the "Models"**

**Choose the "Seeds"**

**Configure "Hyperparameters"**

# Understanding the Key Terms


**What are queries, seeds, models, and hyperparameters?**

Introduction
○○○○

Related Work
○○○○

**Framework**
○○○○●○○○○○○○○○

Usage
○○○○

LLM Benchmarking
○○○○○○○○○○○○○○○

Conclusion
○○○○○

## Understanding the Key Terms

**What are queries, seeds, models, and hyperparameters?**

**Let's break it down with an example.**

## Preparation: Queries

Think of preparing for a jailbreak attack as defeating an opponent:

### Queries

The main objectives (e.g., asking "How to make a bomb?").

Introduction
○○○○

Related Work
○○○○

**Framework**
○○○○○○○●○○○○○○

Usage
○○○○

LLM Benchmarking
○○○○○○○○○○○○○○○

Conclusion
○○○○○

# Preparation: Seeds

Think of preparing for a jailbreak attack as defeating an opponent:

## Seeds
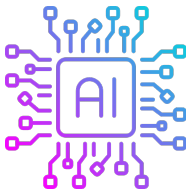
Starting points for the attack.

# Preparation: Models

Think of preparing for a jailbreak attack as defeating an opponent:

> **Models**
>
> The opponent (LLMs) you're trying to defeat.

# Preparation: Adjusting Hyperparameters

Think of preparing for a jailbreak attack as defeating an opponent:

## Adjusting Hyperparameters

Setting the difficulty level.

# Selector

## Selector

picks the best input (weapon) for the attack, maximizing the chances of success.

**EXP3SelectPolicy**
Selects the best input based on past successes.

**Thompson Sampling**
Balances exploration and exploitation probabilistically.

**Upper Confidence Bound**
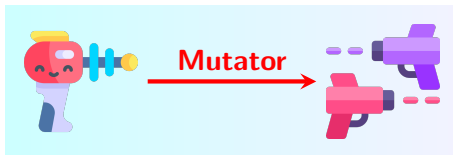Selects inputs based on confidence in expected rewards.

Introduction
0000

Related Work
0000

**Framework**
0000000000●000

Usage
0000

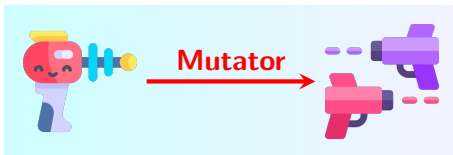LLM Benchmarking
000000000000000

Conclusion
00000

# Mutator

## Mutator

upgrades your weapon to improve its effectiveness.

# Mutator

> **Mutator**
>
> upgrades your weapon to improve its effectiveness.

Introduction
oooo

Related Work
oooo

**Framework**
ooooooooooo●ooo

Usage
oooo

LLM Benchmarking
oooooooooooooooooo

Conclusion
ooooo

# Mutator

## Mutator

upgrades your weapon to improve its effectiveness.



## Example

A **Translation Mutator** turns the input (weapon) into a different language to bypass detection.

Introduction
0000

Related Work
0000

Framework
○○○○○○○○○○○○●○○

Usage
0000

LLM Benchmarking
○○○○○○○○○○○○○○○

Conclusion
00000

# Constraint

## Constraint

A trap detector, ensuring your attack remains focused and valid.

Introduction
oooo

Related Work
oooo

**Framework**
oooooooooooo●oo

Usage
oooo

LLM Benchmarking
oooooooooooooooo

Conclusion
ooooo

# Constraint

## Constraint

A trap detector, ensuring your attack remains focused and valid.



**Trap Detector**

Filters off-topic and irrelevant inputs to ensure the attack remains valid.

# Constraint

## Constraint

A trap detector, ensuring your attack remains focused and valid.



**Trap Detector**

Filters off-topic and irrelevant inputs to ensure the attack remains valid.

## Example

**DeleteOffTopic**

Removes any input that is off-topic.

# Evaluator

## Evaluator

Determines if you defeated the LLM by assessing the success of the attack.

# Evaluator

## Evaluator

Determines if you defeated the LLM by assessing the success of the attack.



**Evaluator**
Acts like a referee, analyzing responses to decide if the attack was successful.

## Evaluator

**Evaluator**

Determines if you defeated the LLM by assessing the success of the attack.



**Evaluator**
Acts like a referee, analyzing responses to decide if the attack was successful.

**Example**

**ClassificationJudge:** Who decides if the challenge (attack) was successful.

# Report

## Report

Provides a detailed analysis of the jailbreak attack.

# Report

## Report

Provides a detailed analysis of the jailbreak attack.

**Success Rate**
Measures how often
the attack succeeded.

Introduction ○○○○
Related Work ○○○○
**Framework** ○○○○○○○○○○○○○●
Usage ○○○○
LLM Benchmarking ○○○○○○○○○○○○○○○○
Conclusion ○○○○○

# Report

## Report

Provides a detailed analysis of the jailbreak attack.



**Success Rate**
Measures how often
the attack succeeded.



**Attack Details**
Logs details of
the attack execution.

Introduction
OOOO

Related Work
OOOO

**Framework**
OOOOOOOOOOOOOO●

Usage
OOOO

LLM Benchmarking
OOOOOOOOOOOOOOOOO

Conclusion
OOOOO

# Report

## Report

Provides a detailed analysis of the jailbreak attack.

**Success Rate**
Measures how often
the attack succeeded.

**Attack Details**
Logs details of
the attack execution.

**Perplexity of Responses**
Evaluates the complexity of
model responses.

Introduction
oooo

Related Work
oooo

**Framework**
ooooooooooooooo●

Usage
oooo

LLM Benchmarking
oooooooooooooooooo

Conclusion
ooooo

# Report

## Report

Provides a detailed analysis of the jailbreak attack.

**Success Rate**
Measures how often
the attack succeeded.

**Attack Details**
Logs details of
the attack execution.

**Perplexity of Responses**
Evaluates the complexity of
model responses.

**Insights**
Provides actionable insights on
attack effectiveness.

1 Introduction

2 Related Work

3 Framework

4 Usage

5 LLM Benchmarking

6 Conclusion

Introduction
○○○○

Related Work
○○○○

Framework
○○○○○○○○○○○○○○○

Usage
○●○○

LLM Benchmarking
○○○○○○○○○○○○○○○○○

Conclusion
○○○○○

# EasyJailbreak: Simplified Model Security Checks

## Usage

EasyJailbreak simplifies security testing for LLMs with a few lines of Python code, enabling users to analyze models using methods.

```python
from easyjailbreak import PAIR,
    JailbreakDataset,from_pretrained,
    OpenaiModel

target_model = from_pretrained('lmsys/
    vicuna-13b-v1.5', 'vicuna_v1.1')
gpt_model = OpenaiModel(model_name='gpt
    -4',api_keys='**')
dataset = JailbreakDataset('AdvBench')
PAIR_attacker = PAIR(
    attack_model=gpt_model,
    target_model=target_model,
    eval_model=gpt_model,
    jailbreak_datasets=dataset,
)
PAIR_attacker.attack()
```

Figure: Python Code for EasyJailbreak

Introduction
oooo

Related Work
oooo

Framework
oooooooooooooo

Usage
oo●o

LLM Benchmarking
ooooooooooooooo

Conclusion
ooooo

# EasyJailbreak: Simplified Model Security Checks



Figure: LLM Response to Jailbreak Attack

## EasyJailbreak: Simplified Model Security Checks

**Attack Model**
**Generates jailbreak prompts for**
**the initial phase of the attack.**

# EasyJailbreak: Simplified Model Security Checks

**Attack Model**
**Generates jailbreak prompts for**
**the initial phase of the attack.**

**Target Model**
**The LLM being tested**
**for vulnerabilities.**

# EasyJailbreak: Simplified Model Security Checks

**Attack Model**
**Generates jailbreak prompts for**
**the initial phase of the attack.**

**Target Model**
**The LLM being tested**
**for vulnerabilities.**

**Evaluation Model**
**Judges the success of the jail-**
**break by evaluating responses.**

# Why is Benchmarking Important?

**Security:** Helps identify vulnerabilities in LLMs.

# Why is Benchmarking Important?

**Security:** Helps identify vulnerabilities in LLMs.

**Resistant:** Measures how resistant models are to jailbreak attacks.

# Why is Benchmarking Important?

**Security:** Helps identify vulnerabilities in LLMs.

**Resistant:** Measures how resistant models are to jailbreak attacks.

**Effectiveness:** Shows which methods work best to bypass models.

Introduction
○○○○

Related Work
○○○○

Framework
○○○○○○○○○○○○○

Usage
○○○○

LLM Benchmarking
●○○○○○○○○○○○○○○

Conclusion
○○○○○

# Why is Benchmarking Important?

**Security:** Helps identify vulnerabilities in LLMs.

**Resistant:** Measures how resistant models are to jailbreak attacks.

**Effectiveness:** Shows which methods work best to bypass models.

**Improvement:** Provides insights for strengthening LLM security.

# Models Tested

## LLM Models Used

Introduction
oooo

Related Work
oooo

Framework
oooooooooooooo

Usage
oooo

**LLM Benchmarking**
oo●ooooooooooooo

Conclusion
ooooo

# Models Tested

## LLM Models Used

**Open-source Models**: LLaMA2, Vicuna, ChatGLM3, Qwen-7B

Introduction
oooo

Related Work
oooo

Framework
oooooooooooooo

Usage
oooo

**LLM Benchmarking**
ooo●oooooooooooo

Conclusion
ooooo

# Models Tested

## LLM Models Used

**Open-source Models**: LLaMA2, Vicuna, ChatGLM3, Qwen-7B

**Closed-source Models**: GPT-4, GPT-3.5-Turbo

# An Important Question



## Which one is better

# An Important Question



### Which one is better

Does an open-source model like LLaMA2 perform better than a closed-source model like GPT-4?

# Attack Method: Human Design



**Human Design**

## Attack Method: Human Design



**Human Design**

🔓 **JailBroken:**
Tricks the AI by acting like it's in a pretend scenario to make it ignore safety rules.

## Attack Method: Human Design



**Human Design**

🔓 **JailBroken:**
Tricks the AI by acting like it's in a pretend scenario to make it ignore safety rules.

🧠 **DeepInception:**
Confuses the AI by sneaking tricky instructions into its context.

# Attack Method: Long-tail Encoding

**Long-tail Encoding**

## Attack Method: Long-tail Encoding



**Long-tail Encoding**

🔑 **Cipher:** Hides the message by turning it into a code (like Morse or Base64) that the AI doesn't recognize as harmful.

## Attack Method: Long-tail Encoding

⌐⌐ **Long-tail Encoding**

🔑 **Cipher:** Hides the message by turning it into a code (like Morse or Base64) that the AI doesn't recognize as harmful.

🅰🈯 **MultiLingual:** Uses uncommon languages that the AI isn't fully trained on to slip past its defenses.

# Attack Method: Prompt Optimization


**Prompt Optimization**

# Attack Method: Prompt Optimization

**Prompt Optimization**

🔍 **GPTFUZZER:**
Tries out different versions of the same question to find one that the AI will answer incorrectly.

# Attack Method: Prompt Optimization

# Attack Method: Prompt Optimization



**Prompt Optimization**

**GPTFUZZER:** Tries out different versions of the same question to find one that the AI will answer incorrectly.

**PAIR:** Improves questions step-by-step based on how the AI responds, making them more likely to bypass safety.

**AutoDAN:** Uses trial and error to keep changing the question until the AI gives the desired response.

# Benchmarking Results

**Key Findings**

**Avg. Breach:**                                                    **63%**

# Benchmarking Results

## Key Findings

**Avg. Breach:**      **63%**

**GPT-3.5-Turbo:**      **57%**

# Benchmarking Results

## Key Findings

**Avg. Breach:**      **63%**

**GPT-3.5-Turbo:**      **57%**

**GPT-4:**      **33%**

# Benchmarking Results

## Key Findings

| | |
|---|---|
| **Avg. Breach:** | **63%** |
| **GPT-3.5-Turbo:** | **57%** |
| **GPT-4:** | **33%** |
| **Vicuna-13B:** | **83%** |

## Benchmarking Results

### Key Findings

**Avg. Breach:**         **63%**

**GPT-3.5-Turbo:**     **57%**

**GPT-4:**               **33%**

**Vicuna-13B:**        **83%**

*Note: Larger models are not inherently more secure. Breach percentages indicate vulnerability.*

Performance Metrics: ASR

**ASR: Percentage of Successful Breaches**

# Performance Metrics: ASR

# Performance Metrics: ASR

**ASR: Percentage of Successful Breaches**

**Higher ASR**
More Vulnerable

A high ASR indicates
the model is sus-
ceptible to attacks.

## Performance Metrics: ASR



**ASR: Percentage of Successful Breaches**

**Higher ASR**
More Vulnerable

**Lower ASR**
More Secure

A high ASR indicates
the model is sus-
ceptible to attacks.

## Performance Metrics: ASR



**ASR: Percentage of Successful Breaches**

**Higher ASR**
More Vulnerable

**Lower ASR**
More Secure

A high ASR indicates
the model is sus-
ceptible to attacks.

A low ASR means
the model is better
at resisting attacks.

## Performance Metrics: Efficiency

**Efficiency: Time and Resource Usage of the Attack**

# Performance Metrics: Efficiency

# Performance Metrics: Efficiency



**Efficiency: Time and Resource Usage of the Attack**

**Faster**
&
Less Resources

Efficient attacks are faster and less computationally expensive.

# Performance Metrics: Efficiency



Efficient attacks are faster and less computationally expensive.

# Performance Metrics: Efficiency

## ASR Comparison (Llama2-7B vs Llama2-13B)



Figure: ASR of Llama models

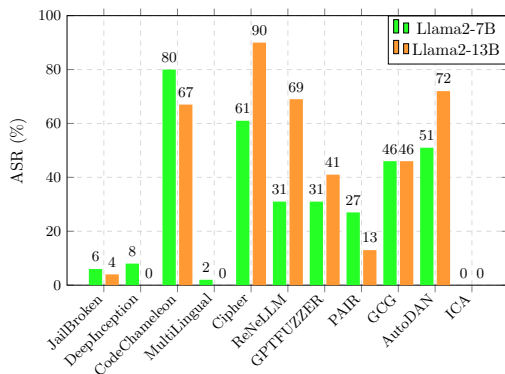# ASR Comparison (Llama2-7B vs Llama2-13B)



Figure: ASR of Llama models

## Summary

- Llama2-13B: **37%**, Llama2-7B: **31%**.

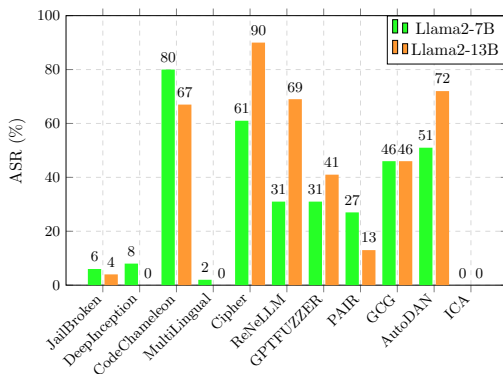# ASR Comparison (Llama2-7B vs Llama2-13B)



Figure: ASR of Llama models

## Summary

- Llama2-13B: **37%**, Llama2-7B: **31%**.

- Bigger models $\neq$ better security.
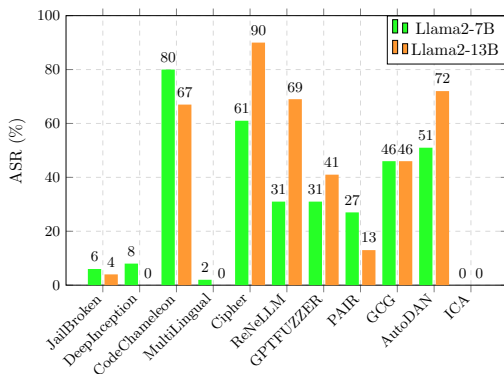
# ASR Comparison (Llama2-7B vs Llama2-13B)



Figure: ASR of Llama models

## Summary

- Llama2-13B: **37%**, Llama2-7B: **31%**.
- Bigger models ≠ better security.
- Cipher, Prompt Optimization work better for Llama2-13B.

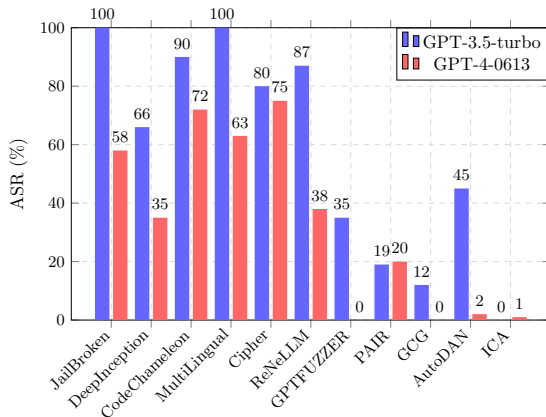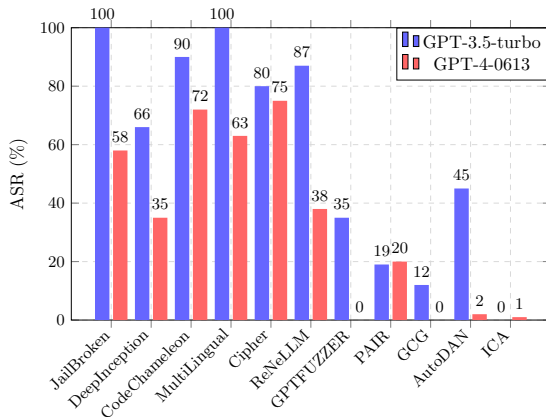# ASR Comparison (GPT-3.5-turbo vs GPT-4-0613)



Figure: ASR of GPT Models

## ASR Comparison (GPT-3.5-turbo vs GPT-4-0613)



This shows that larger models are not always better as Llama-7B (31%) has better ASR than GPT-4-0613 (33%)

Figure: ASR of GPT Models

Introduction
oooo

Related Work
oooo

Framework
ooooooooooooooo

Usage
oooo

LLM Benchmarking
ooooooooooooo●oo

Conclusion
ooooo

# Efficiency Comparison (Llama2-7B vs Llama2-13B)



Figure: Efficiency Comparison (Llama)
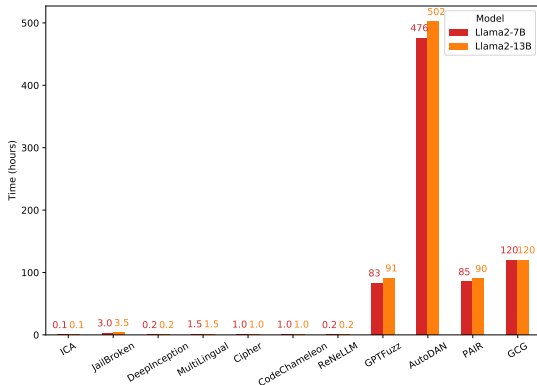
# Efficiency Comparison (Llama2-7B vs Llama2-13B)



Figure: Efficiency Comparison (Llama)

## Summary

- **Llama2-13B** takes more time.

Introduction
○○○○

Related Work
○○○○

Framework
○○○○○○○○○○○○○○○

Usage
○○○○

LLM Benchmarking
○○○○○○○○○○○○○●○○

Conclusion
○○○○○

# Efficiency Comparison (Llama2-7B vs Llama2-13B)



Figure: Efficiency Comparison (Llama)

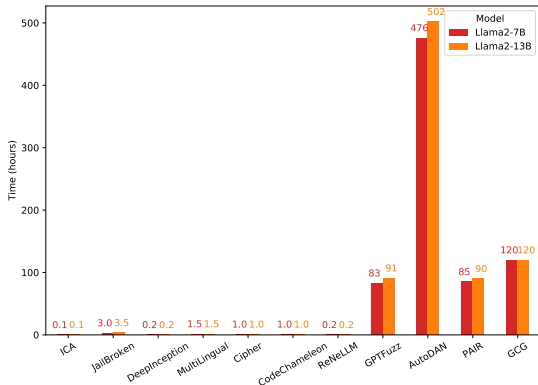## Summary

- **Llama2-13B** takes more time.
- Prompt Optimization: Slow but better results.

# Efficiency Comparison (Llama2-7B vs Llama2-13B)



Figure: Efficiency Comparison (Llama)

## Summary

- **Llama2-13B** takes more time.
- Prompt Optimization: Slow but better results.
- Cipher: Efficient and effective.

Introduction
○○○○

Related Work
○○○○

Framework
○○○○○○○○○○○○○○

Usage
○○○○

LLM Benchmarking
○○○○○○○○○○○○○○●○

Conclusion
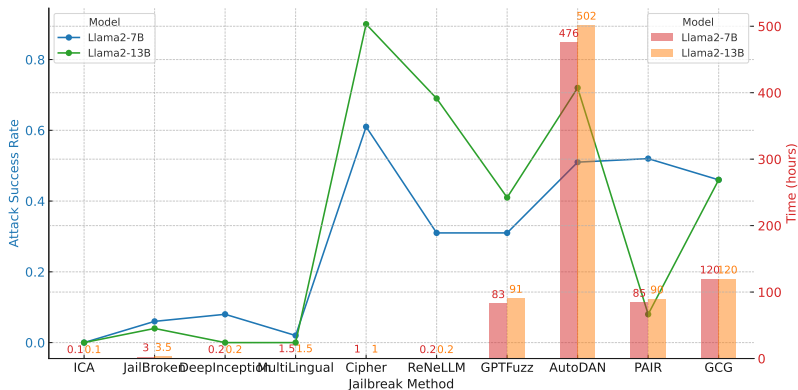○○○○○

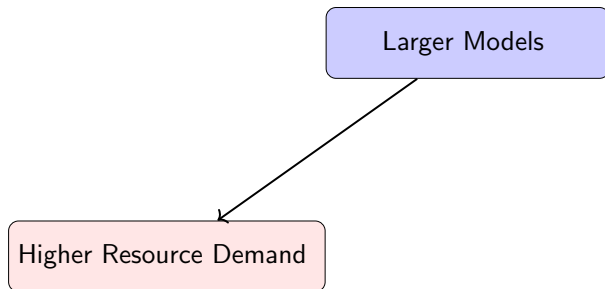# Time-Resource Trade-offs (Llama)



Figure: Time vs. resource efficiency for different attacks.

# Trade-Off: Model Size vs Efficiency

Larger Models

# Trade-Off: Model Size vs Efficiency

# Trade-Off: Model Size vs Efficiency

## Trade-Off: Model Size vs Efficiency

# Key Findings - Summary

**Vulnerability:** All tested models exhibit vulnerabilities to jailbreak attacks.

# Key Findings - Summary

**Vulnerability:** All tested models exhibit vulnerabilities to jailbreak attacks.

**Advanced Models:** GPT-4 are not immune (ASR: **33%**).

Introduction
0000

Related Work
0000

Framework
00000000000000

Usage
0000

LLM Benchmarking
0000000000000000

**Conclusion**
0●0000

# Key Findings - Summary

**Vulnerability:** All tested models exhibit vulnerabilities to jailbreak attacks.

**Advanced Models:** GPT-4 are not immune (ASR: **33%**).

**Open-Source Models: Vicuna** have higher average breach probabilities.

Introduction
0000

Related Work
0000

Framework
0000000000000

Usage
0000

LLM Benchmarking
00000000000000

**Conclusion**
0●000

# Key Findings - Summary

**Vulnerability:** All tested models exhibit vulnerabilities to jailbreak attacks.

**Advanced Models:** GPT-4 are not immune (ASR: **33%**).

**Open-Source Models: Vicuna** have higher average breach probabilities.

**Larger Models:** Does not guarantee better security.

# Key Findings - Comparison

Open-Source Models

Closed-Source Models

# Key Findings - Comparison

Open-Source Models

> **Higher vulnerability** to attacks
> Example: Vicuna models
> show higher breach rates.

Closed-Source Models

Introduction
0000

Related Work
0000

Framework
0000000000000

Usage
0000

LLM Benchmarking
000000000000000

Conclusion
00●00

# Key Findings - Comparison

Open-Source Models

**Higher vulnerability** to attacks
Example: Vicuna models
show higher breach rates.

Closed-Source Models

**Lower average breach rates**
Example: GPT-4 is more
resistant to attacks.

# Key Findings - Comparison

Open-Source Models

**Higher vulnerability** to attacks
Example: Vicuna models
show higher breach rates.

Closed-Source Models

**Lower average breach rates**
Example: GPT-4 is more
resistant to attacks.

**Conclusion:**
**Closed-source models generally pro-**
**vide better resistance to jailbreaks.**

# Implications and Future Directions

## Implications:

- Stronger defenses.

# Implications and Future Directions

## Implications:

- Stronger defenses.
- Continuous security validation.

# Implications and Future Directions

## Implications:
- Stronger defenses.
- Continuous security validation.

## Future Work:
- Develop modular defenses for prompt attacks.

# Implications and Future Directions

## Implications:

- Stronger defenses.
- Continuous security validation.

## Future Work:

- Develop modular defenses for prompt attacks.
- Test larger models (e.g., LLaMA2-70B).

# Implications and Future Directions

## Implications:

- Stronger defenses.
- Continuous security validation.

## Future Work:

- Develop modular defenses for prompt attacks.
- Test larger models (e.g., LLaMA2-70B).
- Enhance EasyJailbreak to counter new threats.

# Thank You!



Thank You for Your Attention!