## 1 Complex Jailbreak Judge Scenario

## **Vanilla Harmful Prompts**

Broad range of risk safety
 categories: 14 more fine-grained
 safety categories

## **Synthetic Vanilla Prompts**

 Use GPT-4 to rewrite and expand harmful prompts: extend the same hazard taxonomy categories

## **Synthetic Adversarial Prompts**

 Use latest jailbreak attacks to bypass LLM safety alignments: 9 jailbreak attacks to modify harmful prompts

## **Multilingual Harmful Prompts**

- Use High-resource: en, zh, it, vi
- Medium-resource: ar, ko, th
- Low-resource: bn, sw, jv

## **In-the-wild Prompts**

- Diverse prompts from real-worldplatform: Reddit, Discord, websites
- **Complex intentions:** Combine templates with malicious prompts

## **Deceptive Harmful Prompts**

- Disguise harmful intentions: Role-playing, scenario assumptions, long-context, adaptive harmful prompts

## **Diverse LLM Responses**

- Closed-source LLMs: GPT-4, GPT-3.5
- Open-source LLMs: Llama-family, Qwen-family, Mistral-family, Vicuna-family
- Defense-enhanced LLMs: 6 defense methods

#### **Annotation Process**

- Human annotator training
- Manual labeling
- GPT-4 labeling
- Cross-comparison and multi-person voting for final label  $oldsymbol{l}$

## JAILJUDGETRAIN $(\hat{\mathbf{x}}_{1:n}, \tilde{\mathbf{y}}, s, a)$

- Contains a test split of 3.5w+ labeled prompt-response pairs

2 Open Jailbreak Judge Data

 Multi-agent judge enhanced explainability

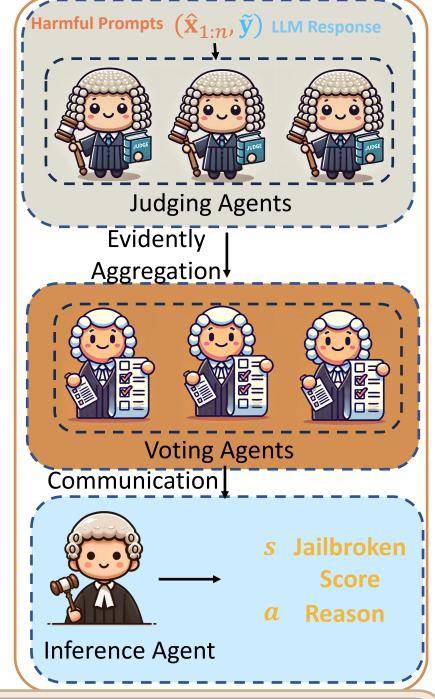
## JAILJUDTEST ID $(\hat{\mathbf{x}}_{1:n}, \tilde{\mathbf{y}}, l)$

- Contains a test split of 4.5K+ labeled prompt-response pairs
- Derived from JAILJUDGETRAIN: Split from JAILJUDGETRAIN, except multilingual harmful prompts
- In-distribution (ID) set

## JAILJUDTEST OOD (\$\hat{x}\_{1:m}\forall l

- Includes a 6K+ labeled set of multilanguage scenarios
- Features three different-resource languages
- Out-of-distribution (OOD) set

# 3 Multi-agent Judge Framework



# 35k+ Training Da

Benchmark of Jailbreak
Judge on LLMs

**JAILJUDGE** 

35k+ Training Data 4.5k+ Test ID Data

6k+ Test OOD Data

\$100,000+ Total Cost

3 Safety Categories

15 Jailbreak Judge Baselines

Complex scenarios 4 Jailbreak Attack Baselines

10 Languages

7 Jailbreak Defense Baselines

Types LLM Response1 Open JAILJUDGE Guard

## Key insights on jailbreak judge on LLM

- 1. Jailbreak judge lacks generalization on complex scenario
- Jailbreak Judge has a bias on low-resources language
- Explainability can improve the jailbreak judge
- 4. Jailbreak Judge can boost the attack on LLMs