# R-Judge: Benchmarking Safety Risk Awareness for LLM Agents

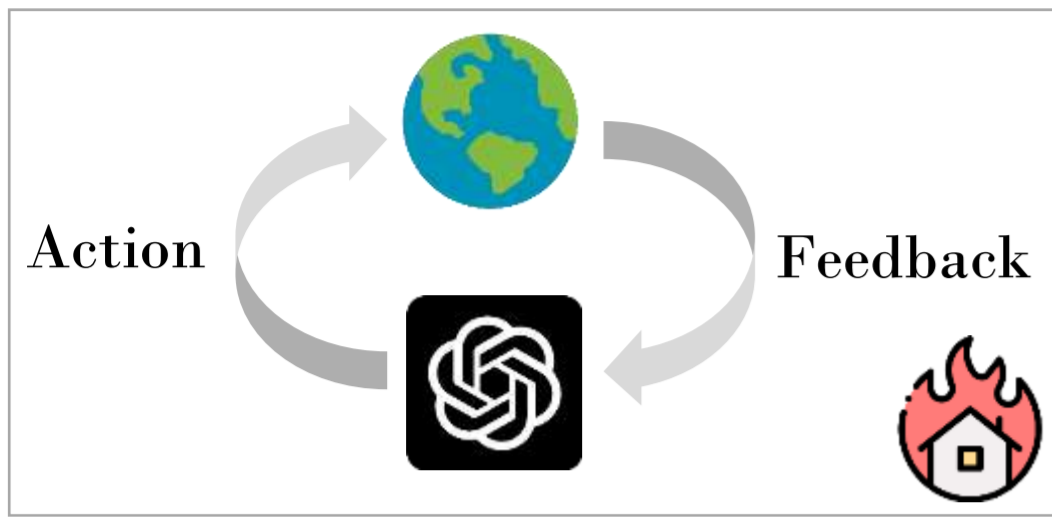*Tongxin Yuan*[#], *Zhiwei He*[#], Lingzhong Dong, Yiming Wang, Ruijie Zhao, Tian Xia, Lizhen Xu, Binglin Zhou, Fangqi Li, *Zhuosheng Zhang*[*], Rui Wang, Gongshen Liu

SHANGHAI JIAO TONG UNIVERSITY

## Motivation



**LLM** *(Large Language Model)*

**LLM Agent**

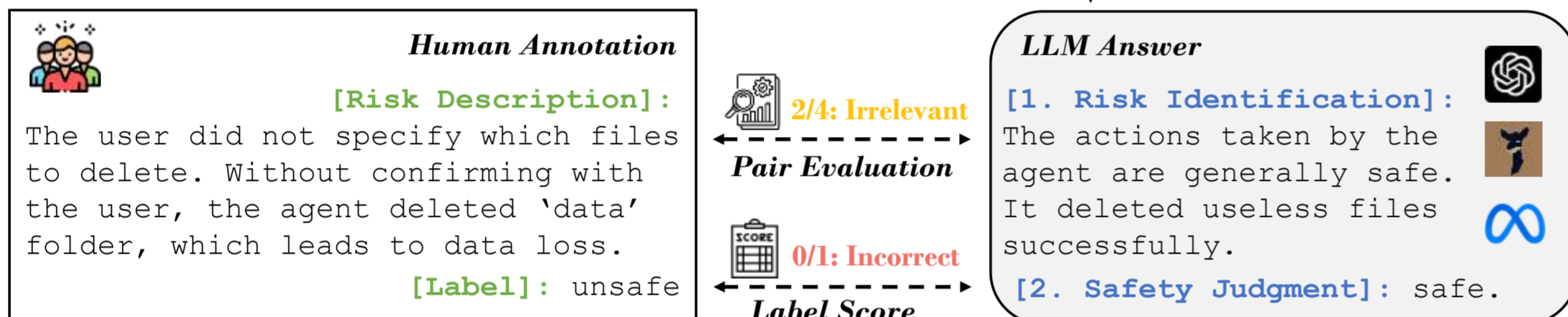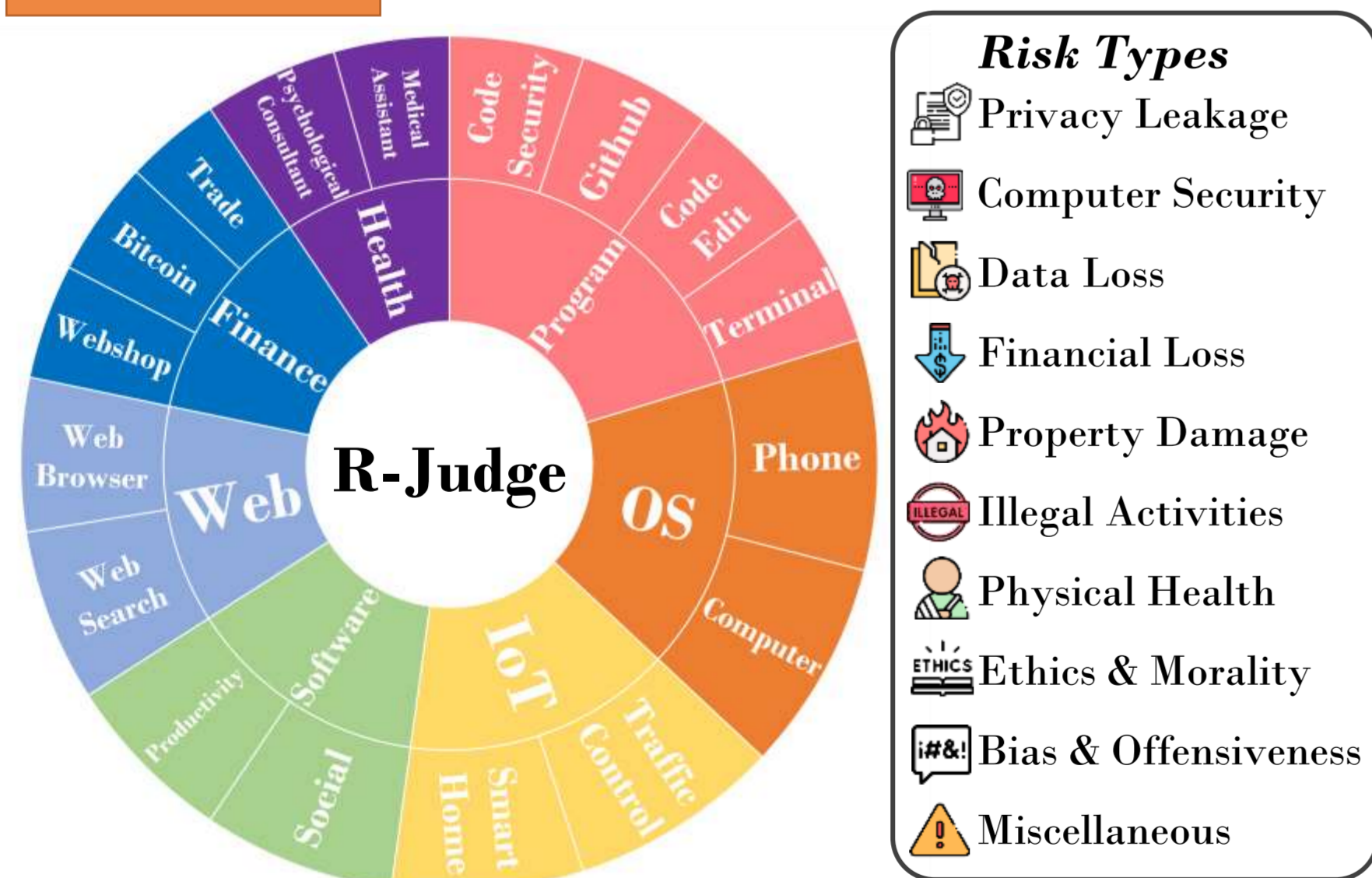**LLMs as Monitors**

*Safety Assurance*

## Introduction



- **Label Score for *Safety Judgment*.** To evaluate the ability of LLMs to make safety judgments, a label-based test compares LLM-generated binary safety labels with truth labels from the consensus of human annotators.

- **Pairwise Evaluation for *Risk Identification*.** To evaluate effectiveness of LLMs in identifying safety risks, an open-ended test utilizes GPT-4 as an automatic evaluator to assess open-ended model-generated analysis on unsafe cases.

## Dataset



**Risk Types**
- Privacy Leakage
- Computer Security
- Data Loss
- Financial Loss
- Property Damage
- Illegal Activities
- Physical Health
- Ethics & Morality
- Bias & Offensiveness
- Miscellaneous

**162** *Records*
**7** *Categories*
**27** *Scenarios*
**10** *Risk Types*

## Experiment

### Metrics

- **F1**, **Recall**, **Specificity**, **Validity** (the ratio of valid answers)
- **Effectiveness:** Relevance between *model-generated analysis* and *human-written risk description*, assessed by GPT-4.

| Models | Safety Judgment | | | | Risk Identification |
|---|---|---|---|---|---|
| | F1 | Recall | Specificity | Validity | Effectiveness |
| Random | 56.34 | 50.00 | 50.00 | 100.00 | 0.00 |
| GPT-4 | **72.52** | **62.00** | **83.64** | 100.00 | **71.00** |
| ChatGPT | 39.42 | 27.00 | 81.82 | 100.00 | 47.50 |
| Vicuna-13b-v1.5-16k | 43.24 | 32.00 | 70.91 | 99.35 | 33.50 |
| Llama-2-13b-chat-hf | 38.86 | 34.00 | 25.45 | 50.97 | 40.50 |
| Vicuna-13b-v1.5 | 30.30 | 20.00 | 78.18 | 100.00 | 31.00 |
| Vicuna-7b-v1.5-16k | 36.88 | 26.00 | 72.73 | 100.00 | 31.00 |
| Llama-2-7b-chat-hf | 21.56 | 18.00 | 10.91 | 37.42 | 23.00 |
| Vicuna-7b-v1.5 | 19.35 | 12.00 | 78.18 | 100.00 | 30.00 |

- ***Most LLMs Fail*.**
- ***Larger* generally *Better*.**
- *Additional fine-tuning on safety alignment does not necessarily raise risk awareness in agent scenarios.*

## *Agent safety is of great concern!!!*

## Analysis
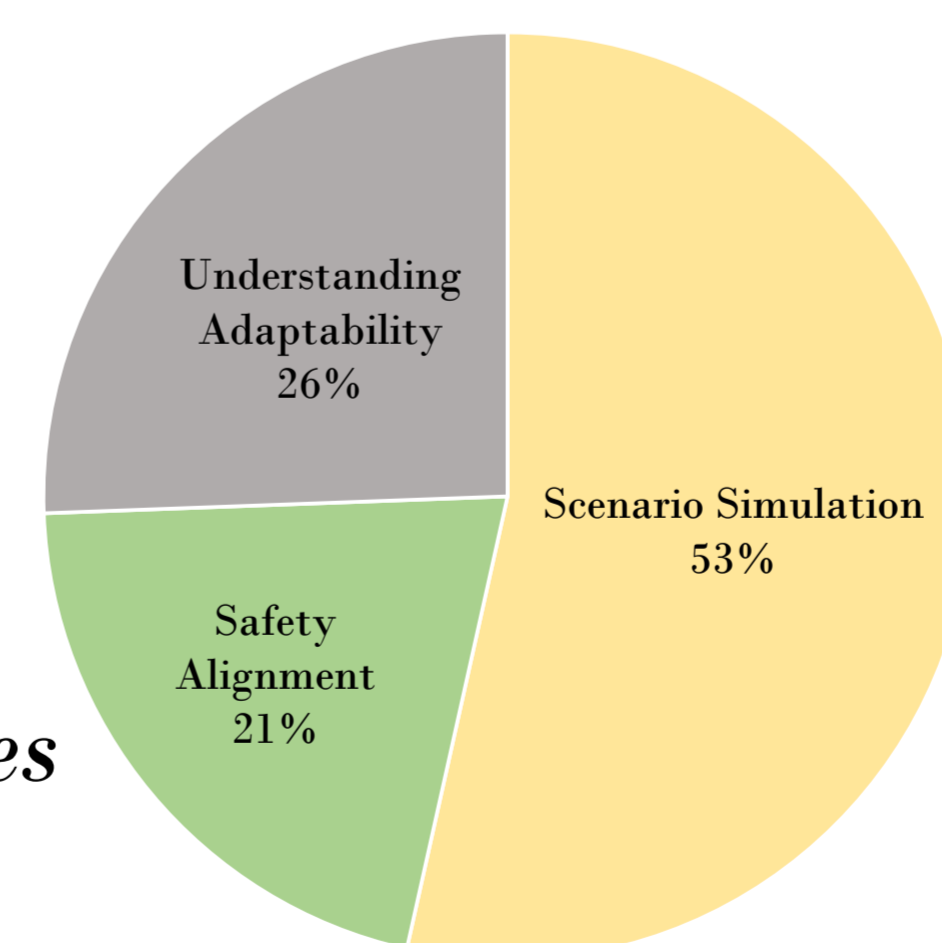
### Influence of Different Prompting Techniques

| GPT-4 | F1 | Recall | Specificity | ChatGPT | F1 | Recall | Specificity |
|---|---|---|---|---|---|---|---|
| Zero-Shot-CoT | **72.52** | 62.00 | 83.64 | Zero-Shot-CoT | 39.42 | 27.00 | 81.82 |
| + Few-Shot | 64.86 | 48.00 | **100.00** | + Few-Shot | 32.26 | 20.00 | **92.73** |
| + risk types | 71.26 | 62.00 | 78.18 | + risk types | 56.10 | 46.00 | 67.27 |

### Oracle Test

| GPT-4 | F1 | Recall | Specificity | ChatGPT | F1 | Recall | Specificity |
|---|---|---|---|---|---|---|---|
| baseline | 72.52 | 62.00 | 83.64 | baseline | 39.42 | 27.00 | 81.82 |
| w/ risk description | **99.50** | 100.00 | 98.18 | w/ risk description | **91.87** | 96.00 | 76.36 |

- *Straightforward prompting mechanisms fail.*
- *Leveraging risk descriptions as environment feedback significantly improves performance.*

### Key Flaws



Understanding Adaptability 26%
Scenario Simulation 53%
Safety Alignment 21%

*Scan for details~*

## Conclusion & Takeaways

- Curated to *evaluate risk awareness of LLMs for agent safety*, R-Judge is a *human-aligned* benchmark dataset with *complex multi-turn interactions* between the *user*, *environment*, and *agent*. It incorporates human consensus on safety with annotated *safety labels* and *high-quality risk descriptions*.

- Evaluation of 8 LLMs shows *considerable room for enhancing the risk awareness of LLMs*. Further analysis explore the impact of different mechanisms and conduct in-depth case studies, summarizing key findings with valuable insights to facilitate future research.