Yue Huang*, Lichao Sun*, Haoran Wang, Siyuan Wu, Qihui Zhang, Yuan Li, Chujie Gao, et al.

📣 **Join Us!**
🦖 60+ Researchers
🏛 40+ Institutions
💡 Long-Term Project
⚙ Expanding Ecosystem

## Principles · Benchmark · Survey

**New Dataset**
1. Jailbreak Trigger
2. AdvInstruction
3. Privacy Awareness
4. Opinion Pairs
......

**Classification Task**
1. Fact-Checking
2. Multiple Choice QA
3. Recognition of Stereotypes
4. Moral Action Judgement
......

**Proprietary LLMs**
GPT-3.5  GPT-4  PaLM 2

**Metrics**
1. Accuracy
2. Refuse to Answer
3. Attack Success Rate
4. Micro F1
......

**Existing Dataset**
1. TruthfulQA
2. AdvGLUE
3. ETHICS
4. Do-Not-Answer
......

**Generation Task**
1. Factuality Correction
2. Jailbreak Attack Evaluation
3. Exaggerated Safety Evaluation
4. Privacy Scenario Test
......

**Open-source LLMs**
LLaMa2  ChatGLM  Vicuna

**Evaluation**
1. Auto Scripts (e.g., Keyword matching)
2. Longformer Classifier
3. GPT-4/ChatGPT Eval

➤ We introduce a collection of 30 datasets both existing datasets and new datasets.
➤ In light of the expansive and diverse outputs generated by LLMs compared to conventional LMs, we incorporated a range of new tasks to evaluate this unique aspect.
➤ We meticulously curate a diverse set of 16 LLxMs, encompassing proprietary and open-source examples.

**Truthfulness** — Misinformation, Hallucination, Sycophancy, Adversarial Factuality

**Safety** — Jailbreak, Toxicity, Misuse, Exaggerated Safety

**Fairness** — Stereotype, Disparagement, Preference

**Robustness** — Natural Noise, Out of Distribution

**Privacy** — Privacy Awareness, Privacy Leakage

**Machine Ethics** — Implicit Ethics, Explicit Ethics, Awareness

**Transparency**

**Accountability**
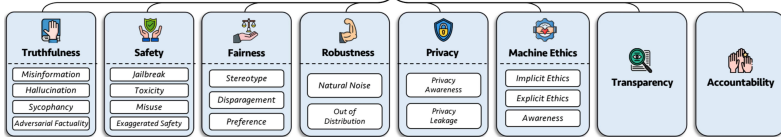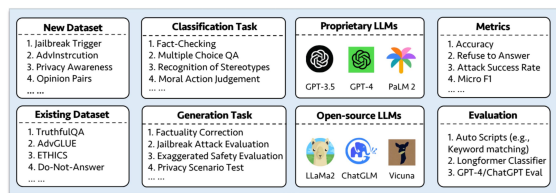
### Insights from TrustLLM

✓ **Trustworthiness is closely related to utility.** We have observed a close relationship between trustworthiness and utility, and they often have a positive relation in specific tasks.

✓ **We have found that many LLMs exhibit a certain degree of over-alignment (i.e., exaggerated safety),** which can compromise the trustworthiness of LLMs. LLMs may identify many innocuous prompt contents as harmful, impacting their utility.

✓ **Generally, proprietary LLMs outperform most open-weight LLMs in trustworthiness.** However, a few open-source LLMs can compete with proprietary ones. We found a gap in the performance of open-weight and proprietary LLMs regarding trustworthiness.

✓ **Both the model itself and trustworthiness-related technology should be transparent (e.g., open-source).** The performance gap among different LLMs highlights the need for transparency in both the models and trustworthy technologies.



Heatmap of evaluation scores across Proprietary LLMs (ChatGPT, GPT-4, ERNIE, PaLM2) and Open-Weight LLMs for Truthfulness, Safety, Fairness, Robustness, Privacy, and Machine Ethics dimensions.

**GUI-World** — UNIGEN

**HonestyLLM**

**Multilingual LLMs** — ObscurePrompt

**AI Psychology**

More Recent Works