# Detecting ChatGPT: A Survey of the State of Detecting ChatGPT-Generated Text

> https://arxiv.org/pdf/2309.07689.pdf

这是一篇聚焦辨别chatgpt生成文本的文章，选择看这篇survey是因为ChatGPT没开源，对它只有API level的access，所以很多方法用不了，想去看一看对于这种强但神秘的模型应该怎么去甄别

- Dataset

| Dataset (name) | Domain | Public | OOD | Size and Setup |
|---|---|---|---|---|
| Guo et al. 2023 (HC3-English) | Multi-domain | ✓ | ✗ | *Q&A*<br>Questions: 24,322<br>Human-A: 58,546<br>ChatGPT-A: 26,903 |
| Guo et al. 2023 (HC3-Chinese) | Multi-domain | ✓ | ✗ | *Q&A*<br>Questions: 12,853<br>Human-A: 22,259<br>ChatGPT-A: 17,522 |
| Yu et al. 2023 (CHEAT) | Scientific | ✗ | ✓ | *Abstracts*<br>Human: 15,395<br>ChatGPT: 35,304 |
| He et al. 2023 (MGTBench) | General | ✓ | ✗ | *Q&A pairs*<br>Human: 2,817<br>ChatGPT: 2,817 |
| Liu et al. 2023 (ArguGPT) | Education | ✓ | ✗ | *Essays*<br>Human: 4,115<br>ChatGPT: 4,038 |
| Vasilatos et al. 2023 | Education | Human* | ✗ | *Q&A*<br>Questions: 320<br>Human-A: 960<br>ChatGPT-A: 960 |
| Mitrović et al. 2023 | General | Human* | ✓ | *Reviews*<br>Human: 1,000<br>ChatGPT-query: 395<br>ChatGPT-rephrase: 1,000 |
| Weng et al. 2023 | Scientific | Human | ✗ | *Title-Abstract pairs*<br>Human: 59,232<br>ChatGPT: 59,232 |
| Antoun et al. 2023a | General | ✓ | ✓ | *Q&A*<br>HC3-English<br>OOD-ChatGPT: 5,969 |
| Liao et al. 2023 | Medical | Human | ✗ | *Abstracts and records*<br>Human: 2,200<br>ChatGPT: 2,200 |

- Method：基本都需要fine-tune模型

| Paper | Dataset | Approaches | Explainability | Code |
|---|---|---|---|---|
| Mitrović et al. 2023 | Mitrović et al. 2023 | DistilBERT<br>PBC | SHAP | ✗ |
| Liao et al. 2023 | Liao et al. 2023 | BERT<br>PBC<br>XGBoost<br>CART | transformer-interpret | ✗ |
| Liu et al. 2023 | Liu et al. 2023 (ArguGPT) | RoBERTa-large<br>SVM | ✗ | ✓* |
| Guo et al. 2023 | Guo et al. 2023 (HC3) | GLTR<br>RoBERTa-single<br>RoBERTa-QA | ✗ | ✓ |
| Antoun et al. 2023a | Antoun et al. 2023a<br>Guo et al. 2023 (HC3) | CamemBERT<br>CamemBERTa<br>RoBERTa<br>ELECTRA<br>XLM-R | ✗ | ✓ |
| Vasilatos et al. 2023 | Ibrahim et al. 2023 | PBC | ✗ | ✗ |

感觉最大的问题还是可解释性，没有可解释性的预测很可能不鲁棒，稍微rewrite一下，变一变prompt可能就会出问题