

---

# Supplementary material for CTIBench: A Benchmark for Evaluating LLMs in Cyber Threat Intelligence

---

## 1 Dataset Documentations

### 1.1 Hosted URLs

**Huggingface:** <https://huggingface.co/datasets/AI4Sec/cti-bench>

**Github:** <https://github.com/xashru/cti-bench/tree/main/data>

**Croissant:** <https://huggingface.co/api/datasets/AI4Sec/cti-bench/croissant>

### 1.2 Dataset description

CTIBench is a comprehensive suite of benchmark tasks and datasets designed to evaluate LLMs in the field of CTI. It consists of the following components:

1. *CTI-MCQ*: A knowledge evaluation dataset with multiple-choice questions to assess the LLMs' understanding of CTI standards, threats, detection strategies, mitigation plans, and best practices. This dataset is built using authoritative sources and standards within the CTI domain, including NIST, MITRE, and GDPR.
2. *CTI-RCM*: A practical task that involves mapping Common Vulnerabilities and Exposures (CVE) descriptions to Common Weakness Enumeration (CWE) categories. This task evaluates the LLMs' ability to understand and classify cyber threats.
3. *CTI-VSP*: Another practical task that requires calculating the Common Vulnerability Scoring System (CVSS) scores. This task assesses the LLMs' ability to evaluate the severity of cyber vulnerabilities.
4. *CTI-TAA*: A task that involves analyzing publicly available threat reports and attributing them to specific threat actors or malware families. This task tests the LLMs' capability to understand historical cyber threat behavior and identify meaningful correlations.

### 1.3 Dataset structure

The dataset consists of 5 TSV files, each corresponding to a different task. Each TSV file contains a "Prompt" column used to pose questions to the LLM. Most files also include a "GT" column that contains the ground truth for the questions, except for "cti-taa.tsv". The evaluation scripts for the different tasks are available in the associated GitHub repository.

### 1.4 Dataset creation

#### 1.4.1 Rationale

The lack of proper benchmark tasks and datasets to evaluate LLM capabilities in CTI leaves their reliability and usefulness an open research question. This dataset was curated to evaluate the ability

31 of LLMs to understand and analyze various aspects of open-source CTI. These datasets evaluate the  
32 reasoning, understanding and problem-solving abilities of LLMs in cyber-threat intelligence. To the  
33 best of our knowledge, CTIBench is the first benchmark specifically designed to evaluate LLMs’  
34 comprehension, reasoning, and problem-solving abilities in the broad domain of CTI, addressing the  
35 limitations of existing benchmarks that either focus on general language understanding or specific  
36 cybersecurity tasks.

#### 37 **1.4.2 Source Data**

38 The dataset includes URLs indicating the sources from which the data was collected. The dataset  
39 is available on the following github repository [https://github.com/xashru/cti-bench/tree/  
40 main/data](https://github.com/xashru/cti-bench/tree/main/data).

#### 41 **1.4.3 Intended use**

42 CTIBench is designed to provide a comprehensive evaluation framework for large language models  
43 (LLMs) within the domain of cyber threat intelligence (CTI). Dataset designed in CTIBench assess  
44 the understanding of CTI standards, threats, detection strategies, mitigation plans, and best practices  
45 by LLMs, and evaluates the LLMs’ ability to understand, and analyze about cyber threats and  
46 vulnerabilities. The intended users of CTIBench are researchers and practitioners in cybersecurity,  
47 Cyber threat analysts, incident responders, LLM developers.

## 48 **2 Author statements**

### 49 **2.1 Responsibility for Rights Violations**

50 We bear all responsibility in the event of any violations of rights, including but not limited to,  
51 intellectual property rights, privacy rights, and any other legal issues that may arise from the use of  
52 the dataset described in this submission. We confirm that we have obtained all necessary permissions  
53 and consents from data owners and subjects where applicable.

### 54 **2.2 Dataset Licensing**

55 We confirm that the dataset is licensed under the Creative Commons Attribution NonCommercial-  
56 ShareAlike 4.0 International.<sup>1</sup> This license allows for the free use, distribution, and modification of  
57 the dataset, provided that appropriate credit is given to the original creators, any modifications are  
58 indicated, and the use is non-commercial.

## 59 **3 Accessibility**

60 We have hosted our dataset on GitHub at [https://github.com/xashru/cti-bench/tree/  
61 main/data](https://github.com/xashru/cti-bench/tree/main/data) and on Hugging Face at <https://huggingface.co/datasets/AI4Sec/cti-bench>.  
62 A permanent DOI identifier is associated with the dataset: DOI: AI4Sec (2024).

63 We are committed to ensuring the long-term preservation of our dataset through periodic checks  
64 to detect and correct any data issues. Additionally, we are dedicated to maintaining this dataset by  
65 addressing user queries and issues promptly via email at [ma8235@rit.edu](mailto:ma8235@rit.edu), and releasing updates  
66 and improvements based on user feedback.

## 67 **4 Reproducibility**

68 To ensure reproducibility, we have provided evaluation notebook, response logs from LLMs on our  
69 github repository <https://github.com/xashru/cti-bench/tree/main/>.

---

<sup>1</sup><https://creativecommons.org/licenses/by-nc-sa/4.0/>

70 **References**

71 AI4Sec. 2024. cti-bench (Revision a86b127). <https://doi.org/10.57967/hf/2506>