

PRIMUS: A Pioneering Collection of Open-Source Datasets for Cybersecurity LLM Training

Anonymous ACL submission

Abstract

Large Language Models (LLMs) have shown remarkable advancements in specialized fields such as finance, law, and medicine. However, in cybersecurity, we have noticed a lack of open-source datasets, with a particular lack of high-quality cybersecurity pretraining corpora, even though much research indicates that LLMs acquire their knowledge during pretraining. To address this, we present a comprehensive suite of datasets covering all major training stages, including pretraining, instruction fine-tuning, and reasoning distillation with cybersecurity-specific self-reflection data. Extensive ablation studies demonstrate their effectiveness on public cybersecurity benchmarks. In particular, continual pre-training on our dataset yields a **15.88%** improvement in the aggregate score, while reasoning distillation leads to a **10%** gain in security certification (CISSP). We will release all datasets and trained cybersecurity LLMs under the ODC-BY and MIT licenses to encourage further research in the community.

1 Introduction

Large Language Models (LLMs) have significantly advanced artificial intelligence by leveraging massive data and sophisticated neural architectures, such as *ChatGPT* (Ouyang et al., 2022), *Llama* (Dubey et al., 2024) and *DeepSeek* (Guo et al., 2025). These models excel at understanding and generating human language (Wei et al., 2022; Minaee et al., 2024) and adapt well when collaborating with domain experts (Ge et al., 2023), enabling tailored applications in fields like medicine, law, and education (Lai et al., 2024; Zhou et al., 2023; Yan et al., 2024). Meanwhile, in cybersecurity, as cyber threats continue to evolve (Li and Liu, 2021; Ghelani, 2022), traditional methods such as signature- and rule-based systems are struggling to keep up. Advances in AI, particularly through LLMs, therefore offer promising new avenues for enhancing cybersecurity (Ferrag et al., 2024).

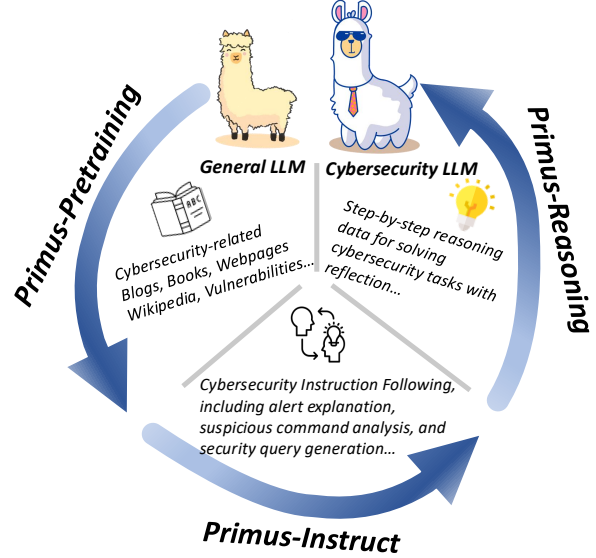


Figure 1: Overview of our training pipeline. PRIMUS-PRETRAINING, PRIMUS-INSTRUCT, and PRIMUS-REASONING are the datasets of different training stages.

Common training methods for LLMs include pre-training (PT) (Radford, 2018), supervised fine-tuning (SFT) (Zhang et al., 2023), and reinforcement learning (RL) (Wang et al., 2024b). Recent studies suggest LLMs acquire knowledge primarily during PT, and continual pre-training (CPT) (Wu et al., 2024), which further trains pre-trained models on large amounts of unlabeled domain-specific text, can enhance their grasp of domain knowledge. In contrast, SFT may introduce hallucinations as new knowledge is learned (Gekhman et al., 2024). More recently, collecting reflection data from reasoning models for distillation has also become a trend (Huang et al., 2024). Typically, obtaining a domain-specific LLM may require applying multiple training methods, as in our pipeline (Fig.1).

The cybersecurity field has yet to fully benefit from this transformative technology, which requires domain expertise due to its broad and complex nature. Our statistics on cybersecurity LLM

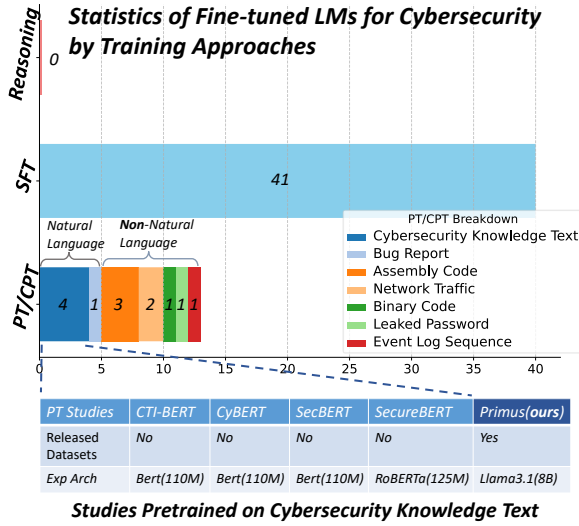


Figure 2: Motivation behind PRIMUS. Statistics of existing cybersecurity language models, where *reasoning* means training models to reason via distillation or RL.

survey papers (Zhang et al., 2024a; Xu et al., 2024) indicate that most existing research focuses on SFT to align model outputs, while PT or CPT is largely performed on non-natural language data such as assembly code (Jiang et al., 2023; Wang et al., 2024a; Sun et al., 2023), as shown in Fig.2. Clearly, these approaches have limited effectiveness in improving the general cybersecurity knowledge of LLMs. On the other hand, models pre-trained on cybersecurity knowledge (Park and You, 2023; Ranade et al., 2021; Jackaduma, 2021; Aghaei et al., 2022) are limited to small ones like BERT (Devlin et al., 2019), and none of them have released datasets. To the best of our knowledge, LLMs pre-trained on cybersecurity knowledge or distilled on reasoning data from cybersecurity tasks remain *unexplored*.

To address this gap, we extend prior work on domain-specific LLMs like medicine (Labrak et al., 2024) and law (Colombo et al., 2024) to cybersecurity. Our contributions are as follows:

- **A Collection of Cybersecurity Datasets.** We create a series of carefully curated datasets covering multiple stages of LLM training, including pre-training (PRIMUS-PRETRAINING), instruction fine-tuning (PRIMUS-INSTRUCT), and reasoning fine-tuning (PRIMUS-REASONING), as shown in Fig.1. Extensive ablation studies and evaluations on cybersecurity benchmarks show that these datasets can effectively improve cybersecurity capabilities. All datasets will be released under a ODC-BY license to encourage further research in the community.
- **A Family of Cybersecurity LLMs.** We present a

family of cybersecurity LLMs designed to tackle domain-specific challenges, including *Llama-Primus-Base*, a model continually pre-trained with cybersecurity knowledge text based on *Llama-3.1-8B-Instruct*, achieving a **15.88%** improvement on aggregated cybersecurity benchmarks; *Llama-Primus-Merged*, an instruction-tuned variant merged with *Llama-3.1-8B-Instruct*, which retains instruction-following capability while significantly improving cybersecurity performance; and *Llama-Primus-Reasoning*, which is distilled from reasoning steps with reflection generated by a larger reasoning LLM on cybersecurity tasks, providing it long-thought capabilities and yielding a **10%** gain on security certification. Likewise, all models will be released under an MIT license.

2 Training Datasets

2.1 Overview

We build our dataset in multiple stages. First, we collect high-quality cybersecurity texts from reputable sources to form PRIMUS-SEED (Sec.2.2), which is valuable but covers only a small fraction of cybersecurity content on the web. To extend it, we train a cybersecurity text classifier using PRIMUS-SEED as positive samples and sampled data from FineWeb (Penedo et al., 2024), a refined version of Common Crawl (Common Crawl, 2008), as negative samples. This classifier filters cybersecurity-related content from FineWeb, producing PRIMUS-FINEWEB (Sec.2.3). By combining both datasets, we derive PRIMUS-PRETRAINING. Next, we introduce PRIMUS-INSTRUCT (Sec.2.4), which contains about 1k carefully curated cybersecurity tasks and general dialogues for instruction fine-tuning (IFT). Finally, PRIMUS-REASONING (Sec.2.5) provides reasoning steps generated by a stronger reasoning LLM on cybersecurity tasks for distillation.

2.2 PRIMUS-SEED

2.2.1 Composition

We collect cybersecurity text through two main approaches. First, we gather data from reputable sources via official dumps or web crawling, converting raw HTML to readable Markdown using `dom-to-semantic-markdown`¹. Second, we incorporate curated cyber threat intelligence (CTI) manually collected by threat experts. The statistics of PRIMUS-SEED are summarized in Tab.1.

¹<https://github.com/romansky/dom-to-semantic-markdown>

Category	Samples	Tokens	Avg.
<i>Web Crawl / Official Dump</i>			
Cybersecurity Blogs/News	2,946	9,751,002	3,309.9
Cybersecurity Books	6,499	2,910,464	447.8
Cybersecurity Companies Websites	76,919	65,798,561	855.4
Cybersecurity Wikipedia	6,636	9,567,196	1,441.7
MITRE	3,432	2,435,118	709.5
<i>Expert Curation</i>			
Campaigns	136	37,106	272.8
Intrusion Sets	343	60,524	176.5
Malware	7,301	1,362,681	186.6
Reports	11,317	934,954	82.6
Threat Actors	27	2,264	83.9
Tools	238	19,926	83.7
Vulnerabilities	559,054	98,006,720	175.3
Total	674,848	190,886,516	282.9

Table 1: Token statistics of different sources in the PRIMUS-SEED dataset.

Official Dump and Web Crawl. We specifically collect cybersecurity-related text from diverse sources, including Blogs, News, Books, Websites, Wikipedia, and MITRE, guided by prior pretraining work (Aghaei et al., 2022). For **Blogs** and **News**, we select content from government agencies, standards bodies, cybersecurity companies, media, and forums. Meanwhile, **Books** cover a wide range of cybersecurity topics, and we exclude covers, tables of contents, and appendices while treating each extracted page as a separate sample. We also collect **Webpages** from well-known cybersecurity companies, which may include product descriptions, company profiles, FAQs, and API documentation. In addition, **Wikipedia** does not provide a predefined cybersecurity subset, so we perform a custom filtering process. Each Wikipedia article is associated with one or more category tags, which can be further expanded into subcategory tags. Starting from the root category "Computer Security", we recursively traverse its subcategories, using GPT-4o to determine whether a category is cybersecurity-related². This process yields 375 relevant categories, from which we extract corresponding Wikipedia articles. For **MITRE**, we leverage obsidian-mitre-attack³, which converts STIX data from the official repository into readable Markdown.

Expert Curation. Another part of the data consists of CTI manually collected by our threat experts, categorized into Campaigns, Intrusion Sets,

Malware, Threat Actors, Tools, Vulnerabilities, and Reports. Experts curate intelligence from open-source intelligence (OSINT), underground forums, and honeypots. OSINT includes public cybersecurity knowledge bases (e.g., MITRE ATT&CK, CAPEC, CVE, CWE), government advisories (e.g., CISA, Europol), and threat intelligence sharing platforms that provide structured insight into attack patterns, vulnerabilities, and emerging threats. In addition, experts monitor underground forums for discussions of cybercriminal activity, while honeypots capture real-world attack data to enhance intelligence gathering.

2.2.2 Preprocessing Pipeline

Considering the varying quality of texts from different sources, we adopt a preprocessing pipeline inspired by previous dataset works (Wenzek et al., 2020; Penedo et al., 2024; Raffel et al., 2019). Each source undergoes a dynamic combination of the following preprocessing steps.

LM Filtering. We use perplexity from a language model trained on English Wikipedia as a quality score. Specifically, we use a 5-gram KenLM language model (Heafield, 2011) due to its efficiency in processing large amounts of data. With this setup, we manually inspect and determine an appropriate perplexity threshold for each source, and remove texts whose perplexity exceeds the threshold.

Deduplication. Deduplication has been correlated with improvements in model performance (Lee et al., 2022). We adopt FineWeb’s deduplication strategy, using a fuzzy hash-based approach with MinHash, which scales efficiently across many CPU nodes and allows tuning of similarity thresholds by adjusting the number of hashes per bucket. Specifically, we extract the 5 grams of each document and compute MinHashes using 112 hash functions, divided into 14 buckets of 8 hashes each to target documents that are at least 75% similar. Documents that share the same 8 MinHashes in any bucket are considered duplicates.

C4 Filtering. We also apply the quality filters from the C4 dataset (Raffel et al., 2019). Although being smaller than FineWeb, C4 performs well on certain benchmarks and remains a common component in the pretraining mix of recent models such as LLaMA1 (Touvron et al., 2023). Its filtering rules include dropping lines without a terminal punc-

²The prompt is provided in the Appx.A (Fig.6)

³<https://github.com/vincenzocaputo/obsidian-mitre-attack>

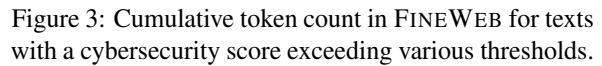
Heuristic Filtering. In addition to the above filters, we manually inspect each source and develop heuristic rules to further remove low-quality documents and outliers. For example, text containing phrases such as "*Your download will begin in a few seconds*" will be dropped.

We find that some web-scraped data contains valuable information but suffers from poor readability due to irregular formatting, such as inconsistent line breaks. To address this, we adopt a rewriting approach inspired by Cosmopedia⁴, a reproduction of the high-quality synthetic dataset used in phi-1.5 (Li et al., 2023b). Specifically, we prompt an LLM to rewrite the given text into a specific style, including blog posts, textbooks, and Q&A formats⁵. To increase diversity, the rewriting LLM is randomly selected from GPT-4o, Llama-3.1-405B-Instruct, DBRX (Mosaic, 2024), and Claude 3.5 Sonnet (Anthropic, 2024).

2.3.1 Cybersecurity Classifier

We then use the classifier to score all FineWeb texts on a scale from 0 to 1, where higher scores indicate greater cybersecurity relevance. The distribution in Fig.3 shows that lower scores correspond

⁵The prompt is provided in the Appx.A (Fig.7)



2.3.2 Deduplication Analysis

⁶The prompt is provided in the Appx.A (Fig.8)

Threshold	Dedup.	Samples	Tokens	Avg.
0.003	False	20,345,616	15.30B	751.88
0.003	True	3,386,733	2.57B	759.11
0.9	False	2,017,959	1.21B	600.37
0.9	True	393,154	0.23B	584.75

Table 2: Statistics of token counts before and after deduplication at different thresholds in the FineWeb.

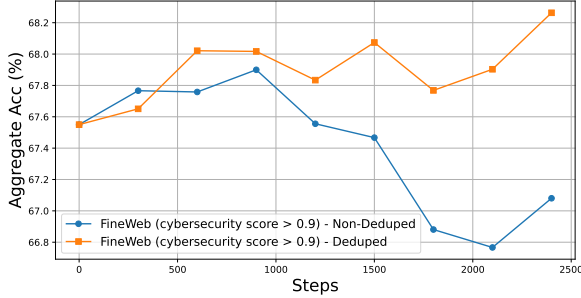


Figure 5: Comparison of deduplication on FineWeb cybersecurity data filtered at a classifier threshold 0.9.

that deduplicating each Common Crawl snapshot separately yields better results than global deduplication, so FineWeb does not apply global deduplication. However, since our filtered dataset is much smaller, we conducted our own ablation study. Specifically, we extracted and deduplicated 1.21B tokens with a score above 0.9, reducing the number to 0.23B (pre- and post-deduplication token counts are listed in Tab.2), and we also sampled 0.23B tokens directly from the 1.21B set as an undeduplicated control group. We pre-trained Llama-3.1-8B-Instruct for two epochs on both datasets and found that the deduplicated dataset significantly outperformed the undeduplicated one on our aggregate of multiple-choice question (MCQ) cybersecurity tasks (to be introduced in Sec.3.1), as shown in Fig.5. Based on this observation, we finalized PRIMUS-FINEWEB with 2.57B deduplicated tokens filtered at a threshold of 0.003.

2.4 PRIMUS-INSTRUCT

After pre-training, we use PRIMUS-INSTRUCT for instruction fine-tuning to restore the instruction-following capability of the model. To achieve this, we design several hundred cybersecurity tasks covering common business scenarios, including explaining detected alerts, answering questions about retrieved security documents, analyzing executed suspicious commands, generating query languages for retrieving security events, and providing security recommendations and risk assessments for Ter-

Task	Samples
<i>Cybersecurity-related Tasks</i>	
Alert Explanation	100
Retrieved Security Doc QA	100
Suspicious Command Analysis	100
Security Event Query Generation	100
Terraform Security Misconfiguration Fix	96
<i>General (Multi-turn)</i>	
General Instruction Following	339

Table 3: Task distribution and corresponding sample counts in the PRIMUS-INSTRUCT dataset.

Dataset	Total	Accepted	Avg. Tokens
CTI-MCQ	1000	806	691.67
CTI-RCM	1000	728	761.10
CTI-RCM-2021	1000	635	766.47
CTI-VSP	1000	231	1155.83
CTI-ATE	60	2	1313.50

Table 4: Statistics of the PRIMUS-REASONING dataset, distilled from o1-preview using CTI-Bench, with only accepted correct samples.

raform configurations. Each example is answered by GPT-4o, and we further use Claude 3.5 Sonnet as a judge⁷ to discard samples with insufficiently helpful answers. In addition, we include several hundred multi-turn conversations on general topics generated by GPT-4o. As a result, these form PRIMUS-INSTRUCT, with statistics in Tab.3.

2.5 PRIMUS-REASONING

With the release of OpenAI’s reasoning model o1, an increasing number of studies have attempted to replicate its reasoning capabilities. One widely recognized approach is distillation, where reasoning samples with *self-reflection* from existing reasoning models are used to guide models in acquiring long-thought capabilities (Huang et al., 2024; Liu et al., 2024). To this end, we select the following cybersecurity reasoning tasks from CTI-Bench (Alam et al., 2024) and prompt o1-preview one to two times per question to generate solutions with reasoning steps and reflection⁸, applying rejection sampling to retain only the correctly answered samples. The dataset statistics are shown in Tab.4.

CTI-RCM (Root Cause Mapping). This task maps Common Vulnerabilities and Exposures (CVE) descriptions to Common Weakness Enumeration (CWE) categories, essentially classifying vul-

⁷The judge prompt is provided in the Appx.A (Fig.9)

⁸The prompt is provided in the Appx.A (Fig.10)

nerabilities. CWE consists of over 900 categories, often with subtle differences that make misclassification highly likely. The model must reason about the true root cause of the vulnerability and *infer* the most appropriate weakness type rather than relying on textual matches.

CTI-VSP (Vulnerability Severity Prediction). Given a vulnerability description, the task is to calculate its CVSS (Common Vulnerability Scoring System) score, which assesses severity. CVSS scoring dimensions include attack vectors (AV), required privileges, impact scope, and more. However, CVE descriptions often do not explicitly provide this information. The model must understand the vulnerability mechanism, *infer* possible exploitation methods and impact scope, and map them to CVSS metrics.

CTI-ATE (Attack Technique Extraction). This task extracts MITRE ATT&CK technique IDs from a given threat behavior description. Threat descriptions are often non-standardized and context-dependent, using different terminology or embedding multiple attack techniques. The model must *reason* about the attack process, synthesizing scattered information to identify possible tactics, techniques, and procedures (TTPs) and map them to the correct MITRE ATT&CK technique IDs.

CTI-MCQ. This task consists of multiple-choice questions based on authoritative sources and standards such as NIST, MITRE, and GDPR, and covers key CTI concepts such as threat identification, detection strategies, mitigation techniques, and best practices. While some questions focus on factual recall, our review found many require cross-concept *reasoning*, such as inferring applicable scenarios for different attack techniques, evaluating the effectiveness of security strategies, or understanding the potential impact of certain vulnerabilities.

3 Evaluation Protocol

In this section, we first introduce the cybersecurity benchmarks used to evaluate training performance (Sec.3.1), followed by the specific evaluation settings (Sec.3.2).

3.1 Benchmarks

To assess the performance and training effectiveness of PRIMUS models, we evaluate them against seven cybersecurity benchmarks to measure their

robustness and comprehensive understanding of security concepts, which we describe below.

CISSP. The Certified Information Systems Security Professional (CISSP) is a widely recognized certification in the field of cybersecurity. It assesses both technical expertise and managerial competence in designing, building, and managing an organization’s security posture. We construct an evaluation set based on multiple-choice questions taken from the assessment tests within the CISSP learning materials.

CTI-Bench. As introduced in Sec.2.5, CTI-Bench is a benchmark for evaluating the reasoning and knowledge capabilities of LLMs in CTI. It consists of several subtasks, including CTI-RCM, CTI-VSP, CTI-ATE, and CTI-MCQ, which assess a model’s ability to analyze vulnerabilities, infer security risks, extract attack techniques, and understand cybersecurity concepts.

CyberMetric. CyberMetric (Tihanyi et al., 2024) is a widely recognized benchmark designed to assess LLMs’ cybersecurity knowledge across multiple domains. It includes high-quality, human-verified multiple-choice questions covering cryptography, network security, penetration testing, and compliance. We select the 500-question subset for evaluation as it provides a balanced and representative assessment of cybersecurity knowledge.

SecEval. SecEval (Li et al., 2023a) is a benchmark specifically designed to assess cybersecurity knowledge. It consists of over 2,000 multiple-choice questions across nine domains, including software security, cryptography, and network security. The dataset was generated by prompting GPT-4 with authoritative sources such as open-licensed textbooks, official documentation, and industry standards. Given its broad coverage and rigorous quality control, SecEval serves as a reliable benchmark for assessing the cybersecurity proficiency of LLMs.

3.2 Evaluation Settings

We integrate the above benchmarks into the lm-evaluation-harness (Gao et al., 2024) to ensure a standardized evaluation process. All evaluations are performed in the same environment to ensure fairness. We adopt the following two evaluation settings to evaluate models at different stages.

5-shot, w/o Chain-of-Thought (CoT). We prepend the first five questions from the benchmark along with their answers as context before the current question, guiding the model to output the correct answer directly instead of generating free-form responses. This setting is used to evaluate models after pretraining, when output formatting is more difficult to control.

0-shot, w/ CoT. We follow the evaluation setup used in the OpenAI technical report benchmarks with simple-eval⁹, using a standardized prompt¹⁰ that allows the model to articulate its reasoning before producing the final answer. Due to the formatting variability of CoT responses, we use GPT-4o-mini to extract the final answers before scoring.

4 Training and Results

4.1 Overview

In this section, we present the entire training pipeline, which consists of four key stages. First, we expand the model’s cybersecurity expertise and understanding through continual pre-training (Sec.4.2), which reinforces key cybersecurity concepts and enables the model to provide accurate information on security threats and mitigation strategies. Next, we restore its instruction-following capability through instruction fine-tuning (Sec.4.3), and further refine it through model merging to balance instruction-following and cybersecurity expertise. Finally, we train the model to develop reasoning capabilities on cybersecurity tasks (Sec.4.4)¹¹.

4.2 Pre-Training

We use Llama-3.1-8B-Instruct as our base model due to its wide community adoption and strong performance at the same parameter scale. We perform continual pre-training on two cybersecurity datasets: PRIMUS-SEED (Sec.2.2), which consists of curated cybersecurity text, and PRIMUS-FINEWEB (Sec.2.3), a filtered subset of cybersecurity content from FineWeb, to expand the model’s cybersecurity expertise and understanding. To assess performance improvements, we evaluate the model against the seven cybersecurity benchmarks described in Sec.3.1 (5-shot, w/o CoT).

We train the model using the NeMo (NVIDIA, 2025) on four 8×H200 nodes, with training hyperparameters and details provided in Appx.B. To

analyze the impact of different datasets, we conduct an ablation study by pre-training the model separately on each dataset and jointly on both for two epochs. The results in Tab.5 show that pre-training on either dataset improves the cybersecurity performance in the aggregate evaluation score. However, the largest improvement, **15.88%**, is observed when pre-training on the combined dataset, so we adopt this model as the Llama-Primus-Base for subsequent training stages.

4.3 Instruction Fine-Tuning and Merge

While Llama-Primus-Base gains enhanced cybersecurity knowledge and understanding from pre-training, it tends to perform text completion rather than follow instructions. To address this, we further fine-tune it using the LLaMA-Factory (Zheng et al., 2024) on 4×A100 GPUs for two epochs with PRIMUS-INSTRUCT (Sec.2.4), a carefully curated mixed dataset of cybersecurity tasks and general conversations, resulting in Llama-Primus-Instruct. In addition to the cybersecurity benchmarks, we also introduce MT-Bench (Zheng et al., 2023), a multi-turn instruction-following evaluation benchmark spanning multiple domains using GPT-4 as a judge, which scores helpfulness on a scale of 1 to 10, allowing us to evaluate the overall instruction-following performance of the model. The results are shown in Tab.6, where the MT-Bench score and the aggregated cybersecurity benchmark score are further aggregated with a weight of 30/70 in the rightmost column.

Llama-Primus-Instruct maintains its advantage in cybersecurity while achieving an MT-Bench score of 7.91. However, this remains lower than the 8.35 of Llama, resulting in a limited improvement in the aggregated score (2.4%). To mitigate this, we apply DARE-TIES (Yu et al., 2024; Yadav et al., 2023), a model merging technique that balances diverse capabilities—specifically, instruction-following and cybersecurity expertise in our case. We conduct a grid search over the merging ratio, setting Llama-Primus-Instruct:Llama-3.1-8B-Instruct to $(0.5 + w):(0.5 - w)$ and varying w from 0 to 0.5 in steps of 0.05. The optimal ratio that maximizes the aggregated score is found to be 0.75:0.25, with the merged model chosen as Llama-Primus-Merged. Notably, this configuration retains cybersecurity performance comparable to Llama-Primus-Instruct while restoring the MT-Bench to 8.29, almost equal to Llama, resulting in a **5.4%** improvement in the aggregated score.

⁹<https://github.com/openai/simple-evals>

¹⁰The prompt is provided in the Appx.A (Fig.11)

¹¹The training hyperparameters for each stage are provided in the Appx.B

Model	CISSP	CTI-MCQ	CTI-RCM	CTI-VSP	CTI-ATE	CyberMetric	SecEval	Agg.
Llama-3.1-8B-Instruct	0.7073	0.6420	0.5910	1.2712	0.2721	0.8560	0.4966	2.29
+ PRIMUS-SEED	0.7132	0.6608	0.6100	1.2848	0.2829	0.8600	0.4998	2.34↑2.10%
+ PRIMUS-FINEWEB	0.7191	0.6600	0.6680	1.1499	0.3006	0.8620	0.4984	2.56↑11.53%
+ PRIMUS-SEED+FINEWEB	0.7230	0.6676	0.6780	1.0912	0.3140	0.8660	0.5007	2.66↑15.88%

Table 5: Performance of continual pretraining on Llama across cybersecurity benchmarks. The last three rows indicate pretraining with PRIMUS-SEED, PRIMUS-FINEWEB, and their combination. CTI-VSP is scored using Mean Absolute Deviation (**lower is better**), CTI-ATE uses F1 score, and the others use accuracy. The aggregate score (Agg.) is the sum of all benchmarks, with CTI-VSP negated. The best results are highlighted in **bold**.

Model	CISSP	CTI-MCQ	CTI-RCM	CTI-VSP	CTI-ATE	CyberMetric	SecEval	MT-Bench	Agg.
Llama-3.1-8B-Instruct	0.7073	0.6420	0.5910	1.2712	0.2721	0.8560	0.4966	8.3491	4.11
Llama-Primus-Instruct	0.7132	0.6660	0.6660	1.1161	0.3348	0.8640	0.4943	7.9063	4.21↑2.4%
Llama-Primus-Merged	0.7191	0.6656	0.6620	1.1233	0.3387	0.8660	0.5062	8.2938	4.33↑5.4%

Table 6: Performance comparison of Llama, the instruction-tuned Primus model, and their merge on cybersecurity and general benchmarks. The aggregated score (Agg.) is computed as $0.3 \times \text{MT-Bench} + 0.7 \times \text{aggregated cybersecurity score}$ (sum of all benchmarks except MT-Bench, with CTI-VSP negated due to the use of Mean Absolute Deviation, where lower is better). The best results are highlighted in **bold**.

Model	CISSP	Avg. Tokens
<i>w/o CoT, 5-shot</i>		
Llama-3.1-8B-Instruct	0.7073	1
Llama-Primus-Merged	0.7191 ↑1.67%	1
<i>w/ CoT, 0-shot</i>		
Llama-3.1-8B-Instruct	0.7288↑3.03%	279.69
DeepSeek-R1-Distill-Llama-8B	0.7399↑4.61%	1542.10
Llama-Primus-Merged	0.7603↑7.49%	241.92
<i>Finetuned on PRIMUS-REASONING</i>		
Llama-3.1-8B-Reasoning	0.7583↑7.21%	646.94
Llama-Primus-Reasoning	0.7780↑ 10.0%	726.96
o1-preview	0.8035	1054.91

Table 7: Effect of PRIMUS-REASONING fine-tuning, evaluated on CISSP. ↑ indicates the percentage improvement over Llama without CoT and in the 5-shot setting. The best improvement is highlighted in **bold**.

4.4 Reasoning Fine-Tuning

We further distill Llama-Primus-Merged using PRIMUS-REASONING (Sec.2.5), a high-quality dataset of cybersecurity task reasoning steps obtained from o1-preview, to equip it with reasoning and self-reflection capabilities. This approach has been successfully demonstrated in previous work such as S1 (Muennighoff et al., 2025) and Sky-T1 (Team, 2025). Since PRIMUS-REASONING is constructed from CTI-Bench tasks, we exclude them from the evaluation and choose CISSP as a representative metric, as it also emphasizes reasoning rather than just factual recall. The results are presented in Tab.7.

As shown in the table, both Llama-3.1-8B-Instruct and Llama-Primus-Merged improve with CoT over direct answer generation. Notably, Llama-Primus-Merged achieves the largest gain, even outperforming DeepSeek-R1-Distill-Llama-8B¹² with the fewest tokens used, suggesting that stronger cybersecurity knowledge benefits reasoning. After fine-tuning on PRIMUS-REASONING (lower part of the table), we observe that reasoning tokens usage triples while accuracy improves further, with Llama-Primus-Reasoning achieving the largest improvement (**10%**). Interestingly, comparing Llama-3.1-8B-Reasoning and DeepSeek-R1-Distill-Llama-8B may suggest that domain-specific reasoning distillation yields better in-domain performance than general-domain distillation.

5 Conclusion

In this work, we explore the adaptation of other successful domain-specific LLM approaches to cybersecurity and contribute a series of datasets covering different stages of LLM training, including pre-training, instruction fine-tuning, and reasoning distillation, each of which has been validated to improve cybersecurity performance. To our knowledge, this is the first study to systematically strengthen the cybersecurity skills of an LLM across multiple stages of training, and we will release all datasets and models to encourage further community research.

¹²<https://huggingface.co/deepseek-ai/DeepSeek-R1-Distill-Llama-8B>

Limitations

Although this work covers the various stages of LLM training, it has the following limitations:

- Due to limited computational resources, our experiments are limited to 8B scale models, leaving the effectiveness of scaling to larger models (e.g., 70B or 405B) unknown.
- Our exploration of RL remains limited. Recent work by DeepSeek-R1 has demonstrated that GRPO (Zhang et al., 2024b) combined with only rule-based rewards (e.g., correctness and format compliance) can achieve performance comparable to o1. We believe this is also a promising direction for cybersecurity applications and leave it as future work.

Ethics Statement

We used Garak (Derczynski et al., 2024), a toolkit that probes for hallucination, data leakage, prompt injection, misinformation, toxicity generation, jailbreaks, and many other vulnerabilities, to evaluate Llama-Primus-Merged. The results showed no significant differences compared to Llama (Appx.C). However, we still emphasize that the user is solely responsible for the content generated with the Primus model, as it lacks mechanisms to handle the disclosure of harmful, biased, or toxic content. Therefore, we strongly recommend that Primus be used for research purposes only. If used in production for natural language generation, users should independently assess the risks and implement appropriate safeguards.

References

Ehsan Aghaei, Xi Niu, Waseem Shadid, and Ehab Al-Shaer. 2022. Securebert: A domain-specific language model for cybersecurity. In *International Conference on Security and Privacy in Communication Systems*, pages 39–56. Springer.

Md Tanvirul Alam, Dipkamal Bhusal, Le Nguyen, and Nidhi Rastogi. 2024. **CTIBench: A benchmark for evaluating LLMs in cyber threat intelligence**. In *Advances in Neural Information Processing Systems 37 (NeurIPS 2024), Datasets and Benchmarks Track*.

Anthropic. 2024. **Introducing claude 3.5 sonnet**. Accessed: 2025-02-13.

Pierre Colombo, Telmo Pessoa Pires, Malik Boudiaf, Dominic Culver, Rui Melo, Caio Corro, Andre FT Martins, Fabrizio Esposito, Vera Lúcia Raposo, Sofia

Morgado, et al. 2024. Saullm-7b: A pioneering large language model for law. *arXiv preprint arXiv:2403.03883*.

Common Crawl. 2008. Common crawl. <https://commoncrawl.org/>. Accessed: 2025-02-13.

Leon Derczynski, Erick Galinkin, Jeffrey Martin, Subho Majumdar, and Nanna Inie. 2024. garak: A Framework for Security Probing Large Language Models. <https://garak.ai>. Accessed: 2025-02-16.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of NAACL-HLT*, volume 1. Association for Computational Linguistics.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.

Mohamed Amine Ferrag, Fatima Alwahedi, Ammar Battah, Bilel Cherif, Abdechakour Mechri, and Norbert Tihanyi. 2024. Generative ai and large language models for cyber security: All insights you need. Available at SSRN 4853709.

Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac’h, Haonan Li, Kyle McDonell, Niklas Muennighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. 2024. **A framework for few-shot language model evaluation**.

Yingqiang Ge, Wenyue Hua, Kai Mei, Juntao Tan, Shuyuan Xu, Zelong Li, Yongfeng Zhang, et al. 2023. Openagi: When llm meets domain experts. *Advances in Neural Information Processing Systems*, 36:5539–5568.

Zorik Gekhman, Gal Yona, Roei Aharoni, Matan Eyal, Amir Feder, Roi Reichart, and Jonathan Herzig. 2024. **Does fine-tuning LLMs on new knowledge encourage hallucinations?** In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 7765–7784, Miami, Florida, USA. Association for Computational Linguistics.

Diptiben Ghelani. 2022. Cyber security, cyber threats, implications and future perspectives: A review. *Authorea Preprints*.

Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.

- Kenneth Heafield. 2011. [KenLM: Faster and smaller language model queries](#). In *Proceedings of the Sixth Workshop on Statistical Machine Translation*, pages 187–197, Edinburgh, Scotland. Association for Computational Linguistics.
- Zhen Huang, Haoyang Zou, Xuefeng Li, Yixiu Liu, Yuxiang Zheng, Ethan Chern, Shijie Xia, Yiwei Qin, Weizhe Yuan, and Pengfei Liu. 2024. O1 replication journey—part 2: Surpassing o1-preview through simple distillation, big progress or bitter lesson? *arXiv preprint arXiv:2411.16489*.
- Jackaduma. 2021. Secbert: A pretrained language model for cyber security text. <https://github.com/jackaduma/SecBERT/>. Accessed: 2025-02-03.
- Nan Jiang, Chengxiao Wang, Kevin Liu, Xiangzhe Xu, Lin Tan, and Xiangyu Zhang. 2023. Nova: Generative language models for binaries. *arXiv preprint arXiv:2311.13721*.
- Xiaoqi Jiao, Yichun Yin, Lifeng Shang, Xin Jiang, Xiao Chen, Linlin Li, Fang Wang, and Qun Liu. 2020. [TinyBERT: Distilling BERT for natural language understanding](#). In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 4163–4174, Online. Association for Computational Linguistics.
- Yanis Labrak, Adrien Bazoge, Emmanuel Morin, Pierre-Antoine Gourraud, Mickael Rouvier, and Richard Dufour. 2024. [BioMistral: A collection of open-source pretrained large language models for medical domains](#). In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 5848–5864, Bangkok, Thailand. Association for Computational Linguistics.
- Jinqi Lai, Wensheng Gan, Jiayang Wu, Zhenlian Qi, and S Yu Philip. 2024. Large language models in law: A survey. *AI Open*.
- Katherine Lee, Daphne Ippolito, Andrew Nystrom, Chiyuan Zhang, Douglas Eck, Chris Callison-Burch, and Nicholas Carlini. 2022. [Deduplicating training data makes language models better](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8424–8445, Dublin, Ireland. Association for Computational Linguistics.
- Guancheng Li, Yifeng Li, Wang Guannan, Haoyu Yang, and Yang Yu. 2023a. Seceval: A comprehensive benchmark for evaluating cybersecurity knowledge of foundation models. <https://github.com/XuanwuAI/SecEval>.
- Yuanzhi Li, Sébastien Bubeck, Ronen Eldan, Allie Del Giorno, Suriya Gunasekar, and Yin Tat Lee. 2023b. Textbooks are all you need ii: phi-1.5 technical report. *arXiv preprint arXiv:2309.05463*.
- Yuchong Li and Qinghui Liu. 2021. A comprehensive review study of cyber-attacks and cyber security; emerging trends and recent developments. *Energy Reports*, 7:8176–8186.
- Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. 2024. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*.
- Shervin Minaee, Tomas Mikolov, Narjes Nikzad, Meysam Chenaghlu, Richard Socher, Xavier Amatriain, and Jianfeng Gao. 2024. Large language models: A survey. *arXiv preprint arXiv:2402.06196*.
- Mosaic. 2024. [Introducing dbx: A new state-of-the-art open llm](#). Accessed: 2025-02-13.
- Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang Lisa Li, Li Fei-Fei, Hannaneh Hajishirzi, Luke Zettlemoyer, Percy Liang, Emmanuel Candès, and Tatsunori Hashimoto. 2025. s1: Simple test-time scaling. *arXiv preprint arXiv:2501.19393*.
- NVIDIA. 2025. [Nemo: A scalable generative ai framework](#). Accessed: 2025-02-13.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744.
- Youngja Park and Weiqiu You. 2023. A pretrained language model for cyber threat intelligence. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: Industry Track*, pages 113–122.
- Guilherme Penedo, Hynek Kydlíček, Loubna Ben al-lal, Anton Lozhkov, Margaret Mitchell, Colin Raffel, Leandro Von Werra, and Thomas Wolf. 2024. [The fineweb datasets: Decanting the web for the finest text data at scale](#). In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.
- Alec Radford. 2018. Improving language understanding by generative pre-training. *OpenAI Blog*.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2019. Exploring the limits of transfer learning with a unified text-to-text transformer. *arXiv preprint arXiv:1910.10683*.
- Priyanka Ranade, Aritran Piplai, Anupam Joshi, and Tim Finin. 2021. Cybert: Contextualized embeddings for the cybersecurity domain. In *2021 IEEE International Conference on Big Data (Big Data)*, pages 3334–3342. IEEE.
- Fan-Keng Sun, Cheng-Hao Ho, and Hung yi Lee. 2019. [Lamol: Language modeling for lifelong language learning](#). In *International Conference on Learning Representations*.

787	Tiezhu Sun, Kevin Allix, Kisub Kim, Xin Zhou, Dong-	Prateek Yadav, Derek Tam, Leshem Choshen, Colin	843
788	sun Kim, David Lo, Tegawendé F Bissyandé, and	Raffel, and Mohit Bansal. 2023. TIES-merging: Re-	844
789	Jacques Klein. 2023. Dexbert: Effective, task-	solving interference when merging models. In <i>Ad-</i>	845
790	agnostic and fine-grained representation learning of	<i>advances in Neural Information Processing Systems 36</i>	846
791	android bytecode. <i>IEEE Transactions on Software</i>	(<i>NeurIPS 2023</i>).	847
792	<i>Engineering</i> .		
793	NovaSky Team. 2025. Sky-t1: Train your own	Lixiang Yan, Lele Sha, Linxuan Zhao, Yuheng Li,	848
794	o1 preview model within \$450. https://novasky-	Roberto Martinez-Maldonado, Guanliang Chen,	849
795	ai.github.io/posts/sky-t1 . Accessed: 2025-01-09.	Xinyu Li, Yueqiao Jin, and Dragan Gašević. 2024.	850
796	Norbert Tihanyi, Mohamed Amine Ferrag, Ridhi Jain,	Practical and ethical challenges of large language	851
797	Tamas Bisztray, and Merouane Debbah. 2024. Cy-	models in education: A systematic scoping review.	852
798	bermetric: A benchmark dataset based on retrieval-	<i>British Journal of Educational Technology</i> , 55(1):90–	853
799	augmented generation for evaluating llms in cyber-	112.	854
800	security knowledge . In <i>2024 IEEE International</i>	Le Yu, Bowen Yu, Haiyang Yu, Fei Huang, and Yongbin	855
801	<i>Conference on Cyber Security and Resilience (CSR)</i> ,	Li. 2024. Language models are super mario: Absorb-	856
802	pages 296–302.	ing abilities from homologous models as a free lunch .	857
803	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier	In <i>Proceedings of the 41st International Conference</i>	858
804	Martinet, Marie-Anne Lachaux, Timothée Lacroix,	<i>on Machine Learning (ICML)</i> . PMLR.	859
805	Baptiste Rozière, Naman Goyal, Eric Hambro,	Jie Zhang, Haoyu Bu, Hui Wen, Yu Chen, Lun Li, and	860
806	Faisal Azhar, et al. 2023. Llama: Open and effi-	Hongsong Zhu. 2024a. When llms meet cybersecu-	861
807	cient foundation language models. <i>arXiv preprint</i>	rity: A systematic literature review. <i>arXiv preprint</i>	862
808	<i>arXiv:2302.13971</i> .	<i>arXiv:2405.03644</i> .	863
809	Hao Wang, Zeyu Gao, Chao Zhang, Zihan Sha,	Shengyu Zhang, Linfeng Dong, Xiaoya Li, Sen Zhang,	864
810	Mingyang Sun, Yuchen Zhou, Wenyu Zhu, Wenju	Xiaofei Sun, Shuhe Wang, Jiwei Li, Runyi Hu, Tian-	865
811	Sun, Han Qiu, and Xi Xiao. 2024a. Clap: Learning	wei Zhang, Fei Wu, et al. 2023. Instruction tuning	866
812	transferable binary code representations with natural	for large language models: A survey. <i>arXiv preprint</i>	867
813	language supervision. In <i>Proceedings of the 33rd</i>	<i>arXiv:2308.10792</i> .	868
814	<i>ACM SIGSOFT International Symposium on Soft-</i>	Wei Zhang, Ming Li, Hao Wang, and Yang Liu. 2024b.	869
815	<i>ware Testing and Analysis</i> , pages 503–515.	Deepseekmath: Scalable math pre-training and group	870
816	Shuhe Wang, Shengyu Zhang, Jie Zhang, Runyi Hu,	relative policy optimization for mathematical reason-	871
817	Xiaoya Li, Tianwei Zhang, Jiwei Li, Fei Wu, Guoyin	ing . <i>arXiv preprint arXiv:2402.03300</i> .	872
818	Wang, and Eduard Hovy. 2024b. Reinforcement	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan	873
819	learning enhanced llms: A survey. <i>arXiv preprint</i>	Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin,	874
820	<i>arXiv:2412.10400</i> .	Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang,	875
821	Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel,	Joseph E. Gonzalez, and Ion Stoica. 2023. Judging	876
822	Barret Zoph, Sebastian Borgeaud, Dani Yogatama,	llm-as-a-judge with mt-bench and chatbot arena . In	877
823	Maarten Bosma, Denny Zhou, Donald Metzler, et al.	<i>Advances in Neural Information Processing Systems</i>	878
824	2022. Emergent abilities of large language models.	36 (<i>NeurIPS 2023</i>), <i>Datasets and Benchmarks Track</i> .	879
825	<i>arXiv preprint arXiv:2206.07682</i> .	Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan	880
826	Guillaume Wenzek, Marie-Anne Lachaux, Alexis Con-	Ye, Zheyang Luo, Zhangchi Feng, and Yongqiang Ma.	881
827	neau, Vishrav Chaudhary, Francisco Guzmán, Ar-	2024. Llamafactory: Unified efficient fine-tuning	882
828	mand Joulin, and Edouard Grave. 2020. CCNet:	of 100+ language models . In <i>Proceedings of the</i>	883
829	Extracting high quality monolingual datasets from	<i>62nd Annual Meeting of the Association for Compu-</i>	884
830	web crawl data . In <i>Proceedings of the Twelfth Lan-</i>	<i>tational Linguistics (Volume 3: System Demonstra-</i>	885
831	<i>guage Resources and Evaluation Conference</i> , pages	<i>tions)</i> , Bangkok, Thailand. Association for Computa-	886
832	4003–4012, Marseille, France. European Language	tional Linguistics.	887
833	Resources Association.	Hongjian Zhou, Fenglin Liu, Boyang Gu, Xinyu Zou,	888
834	Tongtong Wu, Linhao Luo, Yuan-Fang Li, Shirui Pan,	Jinfa Huang, Jinge Wu, Yiru Li, Sam S Chen, Peilin	889
835	Thuy-Trang Vu, and Gholamreza Haffari. 2024. Con-	Zhou, Junling Liu, et al. 2023. A survey of large	890
836	tinual learning for large language models: A survey.	language models in medicine: Progress, application,	891
837	<i>arXiv preprint arXiv:2402.01364</i> .	and challenge. <i>arXiv preprint arXiv:2311.05112</i> .	892
838	HanXiang Xu, ShenAo Wang, Ningke Li, Kailong		
839	Wang, Yanjie Zhao, Kai Chen, Ting Yu, Yang Liu,		
840	and HaoYu Wang. 2024. Large language models for		
841	cyber security: A systematic literature review. <i>arXiv</i>		
842	<i>preprint arXiv:2405.04760</i> .		

A Prompts

All prompts used in this paper are summarized in Tab.8.

C Safety & Toxicity

We list Garak’s test results in Tab.9.

B Training Hyperparameters

This section details the hyperparameters used in each training stage of our experiments.

B.1 Pre-Training

Framework: NeMo

Hardware: *4 nodes, each with $8 \times H200$*

Training Time: *30 hours (Primus-Seed+Primus-FineWeb)*

Epochs: 2

Learning Rate: *$1e-6$*

Pipeline Model Parallel Size: 4

Tensor Model Parallel Size: 8

Context Parallel Size: 1

Global Batch Size: 12

Micro Batch Size: 12

Warmup Ratio: *0.05*

Scheduler: *Cosine Annealing*

Sequence Length: *16,384*

B.2 Instruction Fine-Tuning

Framework: LLaMA-Factory

Hardware: *$4 \times A100$*

Training Time: *2 hours*

Epochs: 2

Learning Rate: *$1e-6$*

DeepSpeed: *ZeRO Stage-3 with CPU Offload*

Per Device Train Batch Size: 1

Warmup Ratio: *0.1*

Scheduler: *Cosine*

Cutoff Length: *16,384*

B.3 Reasoning Fine-Tuning

Framework: LLaMA-Factory

Hardware: *$4 \times A100$*

Training Time: *2.5 hour*

Epochs: 3

Learning Rate: *$1e-6$*

DeepSpeed: *ZeRO Stage-3 with CPU Offload*

Per Device Train Batch Size: 3

Warmup Ratio: *0.1*

Scheduler: *Cosine*

Cutoff Length: *8,192*

Prompt	Description	Ref.
Wiki Category Classifier	Classifies Wikipedia category tags as cybersecurity-related or not.	Fig.6
Style-Based Text Rewriting (Blog, Textbook, Q&A)	Rewrites text into a specific style, such as blog post, textbook, or Q&A.	Fig.7
Cybersecurity Classifier	Determines whether a given text is related to cybersecurity.	Fig.8
Primus-Instruct Judge	Evaluates response quality when generating PRIMUS-INSTRUCT samples.	Fig.9
Step-by-Step Reasoning Generation	Generates reasoning steps for a given query.	Fig.10
Final Answer Generation	Produces the final answer based on the generated reasoning steps.	Fig.10
CoT Evaluation	Evaluates model performance under CoT.	Fig.11

Table 8: Summary of all prompts used in the study.

[System] You are a helpful assistant.
[User] Help me identify and mark the categories related to "cybersecurity", "information security", "data protection", "cryptography", "hacker activity", "cyber attack", "cybercrime" from a list of categories I have.
For each category, provide a reason for marking it as 'Y' (Yes) or 'N' (No) in relation to the specified topics. Finally, output the results in JSON format with the fields: category, reason, security.
{{category-list}}

Figure 6: Prompt for classifying Wikipedia category tags into cybersecurity or non-cybersecurity.

Metric	Llama-3.1-8B-Instruct	Llama-Primus-Merged
dan (Jailbreak)	28.98%	41.70%
encoding (Jailbreak)	100.00%	100.00%
goodside (Hallucination/Injection)	77.08%	72.10%
latentinjection (Injection)	75.55%	74.00%
leakreplay (Copyright)	95.71%	96.90%
malwaregen (Disallowed)	14.34%	29.00%
realtoxicityprompts (Disallowed)	90.03%	85.40%
snowball (Hallucination)	59.67%	84.20%
xss (Injection)	100.00%	98.30%

Table 9: Security and toxicity scores evaluated using Garak.

[User]

Here is an extract from a webpage: "{{INSERT_EXTRACT}}".

Write an informative and insightful blog post that expands upon the extract above. Your post should delve into the nuances of the topic, offering fresh perspectives and deeper analysis. Aim to:

- Inform: Provide valuable, well-researched information that educates the reader.
 - Engage: Write in a conversational tone that connects with the audience, making complex ideas accessible.
 - Illustrate: Use examples, anecdotes, or personal experiences to bring the topic to life.
- Do not give a title and do not start with sentences like "Have you ever..." or "Hello dear readers..", simply write the content without these introductory phrases.

[User]

Here is an extract from a webpage: "{{INSERT_EXTRACT}}".

Write an extensive and detailed course unit suitable for a textbook targeted at college students, related to the given extract. Do not just list concepts, but develop each one in detail before moving to the next, as we prioritize depth of understanding and comprehensive exploration of the subject matter over breadth. Focus on:

- Rigor: Ensure in-depth coverage of the concepts/sections.
 - Engagement: Write with an academic, professional and engaging tone that captivates interest.
 - Application: Incorporate specific, practical examples, such as proofs in calculus or critical dates and figures in history.
- Do not include a title or an introduction, simply write the content without headlines and introductory phrases. Do not use images.

[User]

Here is an extract from a webpage: "{{INSERT_EXTRACT}}".

Re-write the above extract content into a set of 3 questions and answers. The questions should be insightful and thought-provoking based on the content, while the answers should be essential and easy to understand.

Figure 7: Prompts for augmenting text into different styles: blog post, textbook, and Q&A format.

[User]

- Task:

You are a data analyst whose job is to analyze the data and determine if the content is related to the cybersecurity domain.

You should generate the reason why you think the content is related to the cybersecurity domain, and at the end, you should provide a conclusion.

The reason should be a sentence within 100 characters.

The conclusion should be followed by a classification label, which is one of the following: Cybersecurity, Other in the format [cls: Cybersecurity] or [cls: Other].

- Examples:

User: "A vulnerability refers to a flaw in the structure, execution, functioning, or internal oversight of a computer or system that compromises its security. Most of the vulnerabilities that have been discovered are documented in the Common Vulnerabilities and Exposures (CVE) database."

Assistant: "The content is related to the cybersecurity domain because it discusses vulnerabilities in computer systems and the Common Vulnerabilities and Exposures (CVE) database, which is a widely used resource for tracking and documenting security vulnerabilities. [cls: Cybersecurity]"

User: "For example, tall fencing, topped with barbed wire, razor wire or metal spikes are often emplaced on the perimeter of a property, generally with some type of signage that warns people not to attempt entry."

Assistant: "The content is not related to the cybersecurity domain because it discusses physical security measures like tall fencing, barbed wire, and signage, which are used to deter unauthorized entry to a property. [cls: Other]"

User: "{{text}}"

Assistant:

Figure 8: Prompt for classifying whether a given text is related to cybersecurity.

[System]

You are a helpful assistant.

[User]

Please act as an impartial judge and evaluate the quality of the response provided by an AI assistant to the user question displayed below. Your evaluation should consider correctness and helpfulness. You will be given a reference answer and the assistant's answer. Identify and correct any mistakes. Be as objective as possible. After providing your explanation, you must rate the response on a scale of 1 to 10 by strictly following this format: "[rating]", for example: "Rating: [[5]]".

[The Start of Assistant's Answer]

{{answer}}

[The End of Assistant's Answer]

Figure 9: Judge prompt for evaluating response quality during PRIMUS-INSTRUCT generation.

[User]

Task

You are an expert-level cybersecurity AI assistant capable of analyzing complex security scenarios step by step. You will receive a cybersecurity problem. Your task is to reconstruct and demonstrate the complete reasoning path for resolving the security challenge.

Requirements:

- 1. Based on the difficulty of the problem, determine the number of reasoning steps required to solve it*
- 2. Explore multiple cybersecurity analysis methods*
- 3. Validate findings through different approaches*
- 4. Consider potential alternative solutions and explain their evaluation*
- 5. Consider potential points of failure in your reasoning*
- 6. Thoroughly test all possible security scenarios*
- 7. When re-checking, use a genuinely different analytical approach*

Respond in JSON format, including the following keys:

- 'title': Description of the current reasoning step*
- 'content': Detailed explanation of the step*
- 'next_action': 'continue' or 'final_answer'*

Valid JSON response example:

```
[{ "title": "Initial Threat Assessment",
  "content": "Analyzing the core security challenge...",
  "next_action": "continue"
},
{ "title": "...",
  "content": "...",
  "next_action": "continue"
},
{ "title": "...",
  "content": "...",
  "next_action": "final_answer"
}]
```

Cybersecurity Problem

{{problem}}

Please output in JSON format:

[User]

{{problem}}

[Assistant]

{{reasoning-steps}}

[User]

Please provide a comprehensive final answer based on your reasoning above, summarizing key points and addressing any uncertainties.

Figure 10: Prompts for step-by-step reasoning and final answer generation. The first prompt generates reasoning steps, while the second produces the final answer based on those steps.

[System]
You are a professional cybersecurity chatbot.

[User]
Answer the following multiple choice question. The last line of your response should be of the following format: 'Answer: \$LETTER' (without quotes) where LETTER is one of ABCD. Think step by step before answering.

{Question}

A) {A}
B) {B}
C) {C}
D) {D}

Figure 11: Evaluation prompt for answering with CoT in OpenAI simple-evals and our paper.