BENCHMARKING LLM-ASSISTED BLUE TEAMING VIA STANDARDIZED THREAT HUNTING

Anonymous authorsPaper under double-blind review

000

001

002003004

010 011

012

013

014

016

018

019

021

025

026

027

028029030

031

033

034

037

040

041

042

043

044

046

047

048

051

052

ABSTRACT

As cyber threats continue to grow in scale and sophistication, blue team defenders increasingly require advanced tools to proactively detect and mitigate risks. Large Language Models (LLMs) offer promising capabilities for enhancing threat analysis. However, their effectiveness in real-world blue team threat-hunting scenarios remains insufficiently explored. This paper presents CYBERTEAM, a benchmark designed to guide LLMs in blue teaming practice. CYBERTEAM constructs a standardized workflow in two stages. First, it models realistic threat-hunting workflows by capturing the dependencies among analytical tasks from threat attribution to incident response. Next, each task is addressed through a set of operational modules tailored to its specific analytical requirements. This transforms threat hunting into a structured sequence of reasoning steps, with each step grounded in a discrete operation and ordered according to task-specific dependencies. Guided by this framework, LLMs are directed to perform threat-hunting tasks through modularized steps. Overall, CYBERTEAM integrates 30 tasks and 9 operational modules to guide LLMs through standardized threat analysis. We evaluate both leading LLMs and state-of-the-art cybersecurity agents, comparing CYBERTEAM against open-ended reasoning strategies. Our results highlight the improvements enabled by standardized design, while also revealing the limitations of open-ended reasoning in real-world threat hunting.

1 Introduction

The increasing frequency and sophistication of cyber threats continue to pose significant challenges to organizational security. In 2024 alone, over 11,000 more (38% increase!) vulnerabilities were reported compared to 2023, as evidenced by the MITRE CVE database (The MITRE Corporation, n.d.). Defenders, commonly known as the **blue team** (Diogenes & Ozkaya, 2018; Rajendran et al., 2011), are under increasing pressure to identify, analyze, and respond to malicious activities in a timely and accurate manner, a process termed **threat hunting**.

Recent advances in Large Language Models (LLMs) have demonstrated impressive potential to augment cybersecurity practices, including malware analysis (Abusitta et al., 2021; Al-Karaki et al., 2024; Qian et al., 2025; Devadiga et al., 2023), penetration testing (Deng et al., 2023; 2024; Happe & Cito, 2023; Muzsai et al., 2024), and fuzzing (Zhang et al., 2025; Oliinyk et al., 2024; Black et al., 2024). Building on this progress, there is growing interest in leveraging LLMs to automate or assist in threat hunting, enabling blue team defenders to scale their investigations across complex threat landscapes and respond to incidents more effectively. However, despite this momentum, the application of LLMs in blue team threat hunting remains underdeveloped. Existing frameworks tend to focus on isolated analytical tasks (Sehgal & Thymianis, 2023; Faghihi et al., 2023; Dash et al., 2022), such as generating advisory recommendations without integrating earlier steps like threat group attribution. This fragmented design limits our understanding of how LLMs perform within complex, interdependent threat-hunting workflows.

To address this gap, we introduce CYBERTEAM, a practical benchmark designed to rigorously evaluate and guide the use of LLMs in blue team threat hunting. CYBERTEAM supports blue team threat-hunting workflows through the following aspects:

Broader Coverage. CYBERTEAM is constructed from a diverse and large-scale repository of threat intelligence data sourced from 23 vulnerability databases, including MITRE (MITRE Corporation,

Table 1: Comparison between CYBERTEAM and other LLM-oriented cybersecurity benchmarks.

Benchmark	Focus	#Data	#Task	#Source	Coverage	Unique Feature
CWE-Bench-Java (Li et al., 2025)	Java vulnerability	120	4	1	Four CWE classes	Large-scale Java codes
CTIBench (Alam et al., 2024)	Cyber Threat Intelligence	2,500	3	6	CVE, CWE, CVSS, ATT&CK	Multi-choice questions (MCQ)
SevenLLM-Bench (Ji et al., 2024)	Report understanding	91,401	28	N/A	Bilingual instruction corpus	Synthetic Data, MCQ, QA
SWE-Bench (Jimenez et al., 2023)	Software bug fixing	2,294	12	1	GitHub issues	Python repository
CYBERTEAM (Ours)	Blue team threat hunting	452,293	30	23	Threat-hunting lifecycle (3.1)	Open Generation, Standardized Reasoning Env

Cyber Threat Log

On Dec. 10, 2024, our SIEM system flagged multiple anomalous outbound DNS requests from internal host host-192-168-10-21.local to dns-update.evilcorp.net. Investigation revealed that the host had received a suspicious email containing an attachment named Invoice_April2025.doc, which, when opened, triggered a connection to a known C2 domain via an obfuscated PowerShell script. The initial vector appears to be a phishing campaign exploiting. The attacker leveraged PowerShell to execute a memory-resident payload that conducted system reconnaissance, credential harvesting (via LSASS dump), and lateral movement using SMB.

Detected IOCs include: C2 Domains: clns-update.evilcorp.cn, smbauth.c2redir.net. IP Addresses: 185.100.87.21, 192.168.10.22

1. Threat Attribution

2. Behavior Analysis

3. Prioritization 4. Response & Mitig

A												A A												
1. Threat Attribution						2	2. Behavior Analysis					3. Prioritization 4.						. Response & Mitigation						
	NER	-	REX	-	RAG	-	MAP		NE	R →	RAG	-	MAP		+	SUM	-	CLS	_	-	RAG	-	SUM	
				Ŧ								7					4					₩		
Evidence Actor				Ob	Observation TTPs					Severity High					Response Action									
C2	Domair	n: si	mbauth.c2	redi	r.net	CVE-	2024-216	78	Pow	PowerShell Obfuscation			T1059.0	001	Zero-day exploitation					Apply Microsoft patch KB5000871				
Ma	Malware dns-update.evilcorp.cn APT41 or TA505		LSASS Memory Dump			T1003.0	001	Credential theft detected				ł	Block connections to *.evilcrop.net											
de					Spear	Spearphishing Attachment T1566.001			001	Internal host comprimised					Isolate affected host 192.169.10.21									
$\overline{}$															$\overline{}$				_	$\overline{}$				

Figure 1: A CYBERTEAM threat hunting example equipped with operational modules. Module names: NER-named entity recognition, REX-regex parsing, MAP-text mapping, RAG-retrieval-augmented generation, CLS-classification, SUM-summarization.

2024), NVD (National Institute of Standards and Technology (NIST), 2024), and CISE (CISE Program, 2024), as well as security platforms such as Red Hat Bugzilla (Red Hat, Inc., 2024), Oracle Security Alerts (Oracle Corporation, 2024), and IBM X-Force (IBM Corporation, 2024). In addition, CYBERTEAM presents a larger number of tasks and samples than existing cybersecurity benchmarks (Jimenez et al., 2023; Li et al., 2025; Alam et al., 2024; Ji et al., 2024), as summarized in Table 1. This extensive coverage allows for a more comprehensive and nuanced evaluation of LLM performance across a wide range of threat-hunting scenarios.

Standardized Workflow. An important feature of CYBERTEAM is its structured, modular workflow for guiding LLMs within a standardized reasoning environment (Yang et al., 2024; Cheng et al., 2025). This design is inspired by blue team practices, where analysts typically follow standardized procedures to interpret threat metadata and conduct investigations (Sehgal & Thymianis, 2023; Diogenes & Ozkaya, 2018; Brotherston et al., 2024). However, strict adherence to such procedures can limit flexibility when analyzing unstructured threat logs or addressing emerging, zero-day threats. To balance standardization and flexibility, CYBERTEAM integrates a set of operational modules that regulate LLM behavior while allowing for open-ended reasoning where needed. As illustrated in Figure 1, CYBERTEAM first models the dependency structure among threat-hunting objectives (e.g., attribution, behavior analysis, mitigation) as a task chain, and then maps this chain into a corresponding sequence of operational modules. In this process, functions such as NER enforce structured outputs (e.g., extracting threat actor entities), while functions like RAG support more flexible reasoning (e.g., summarizing relevant patching strategies).

Evaluation Strategy. CYBERTEAM incorporates agent-based evaluation strategies tailored to each threat-hunting objective. We benchmark leading LLMs and state-of-the-art (SOTA) cybersecurity agents, comparing CYBERTEAM 's modularized approach with popular open-ended reasoning strategies such as In-Context Learning (ICL) (Dong et al., 2022), Chain-of-Thought (CoT) (Wei et al., 2022), Tree-of-Thought (ToT) (Yao et al., 2023). Our evaluation provides insights into the actionable threat hunting across 30 tasks.

In summary, this paper makes the following contributions: (1) We introduce CYBERTEAM, a practice-informed, broadly scoped benchmark that enables rigorous evaluation of LLMs for blue team threat hunting, (2) we construct a standardized reasoning workflow that models the dependencies among

threat-hunting tasks and guides LLMs through standardized yet flexible reasoning workflow, (3) we conduct comprehensive evaluations and provide insights to improve LLM performance among threat-hunting scenarios. To facilitate future research, we release codes at: https://anonymous.4open.science/r/LLM-Cyberteam-7433/.

2 RELATED WORK

LLMs for Cybersecurity. Recently, LLMs have shown promise in enhancing cybersecurity tasks such as malware classification (Abusitta et al., 2021; Al-Karaki et al., 2024; Qian et al., 2025; Devadiga et al., 2023), code vulnerability detection (Russell et al., 2018; Lu et al., 2024; Sheng et al., 2024), penetration testing (Happe & Cito, 2023; Muzsai et al., 2024; Shen et al., 2024), phishing detection (Kulkarni et al., 2024; Greco et al., 2024), and incident report generation (Bernardi et al., 2024; Sufi, 2024; McGregor et al., 2025). These applications leverage the language understanding and reasoning capabilities of LLMs to analyze technical data, recommend solutions, or simulate attacker behaviors. However, existing applications typically target isolated tasks without considering broader analyst workflows. Additionally, their open-ended reasoning often results in hallucinations and inconsistencies (Mündler et al., 2023; Simhi et al., 2025; Shrivastava), raising concerns about reliability in high-stakes defensive scenarios.

Cybersecurity Benchmarks. Recent benchmarks have focused on static analysis (Reinhold et al., 2024; Higuera et al., 2020; Braga et al., 2017), software vulnerabilities (Hossen et al., 2024; Sawant et al., 2024), and threat report generation (Tihanyi et al., 2024; Perrina et al., 2023; Čupka et al., 2023). These benchmarks evaluate predefined tasks such as identifying CWE categories, matching CVEs, or summarizing intelligence reports (Alam et al., 2024; Aghaei et al., 2020; Branescu et al., 2024; Hemberg et al., 2020). While helpful for reproducibility, they often cover narrow domains and lack the complexity and task interdependencies inherent in real-world threat investigations. In contrast, benchmarks from other high-stakes fields (e.g., law, medicine, finance) increasingly include complex, multistep tasks requiring diverse reasoning skills (Fei et al., 2023; Wang et al., 2024; Choshen et al., 2024; Lucas et al., 2024; Zhou et al.). Inspired by these efforts, we introduce CYBERTEAM to emphasize structured reasoning and realistic interdependencies for blue teaming scenarios.

Operation-Guided Agents. Recent research has proposed agents with operational modules to structure LLM reasoning into modular, interpretable steps (Driess et al., 2023; Dongre et al., 2024; Hu et al., 2024). Such frameworks have achieved notable success in robotics (Jeong et al., 2024; Akkaladevi et al., 2021), database querying (Kadir et al., 2024; Dar et al., 2019), and scientific reasoning tasks (Abate et al., 2020; Vaesen & Houkes, 2021). However, their use in cybersecurity, especially defensive operations, remains underexplored despite the need for structured workflows. Our work addresses this gap by introducing a modular environment aligned with blue team practices, enabling procedural reasoning within a structured analytical pipeline.

3 CYBERTEAM

In this section, we provide a detailed introduction of CYBERTEAM regarding the collected threat hunting tasks (3.1), data sources (3.2), and the modularized strategy (3.3).

3.1 THREAT HUNTING TASKS

As shown in Table 2, CYBERTEAM reflects the full lifecycle of threat hunting tasks. Specifically, CYBERTEAM systematizes analytical tasks into four categories: **Threat Attribution**, **Behavior Analysis**, **Prioritization**, and **Response & Mitigation**. Each category captures a stage in the threat-hunting workflow from investigating cyber threats to identifying countermeasures. Specifically:

Threat Attribution aims at uncovering the origins and nature of a threat. This includes tasks such as extracting infrastructure artifacts (e.g., domains, IPs, URLs), classifying malware families based on observed behaviors, matching known threat signatures, and linking activities to known campaigns or actor groups (e.g., APT or MITRE ATT&CK (MITRE Corporation, 2024)). Further granularity is achieved through geographic and temporal pattern analysis, as well as victimology and affiliation linking, all of which help analysts contextualize incidents in terms of their broader threat landscape.

Table 2: Threat hunting tasks, description of targets, corresponding modularized operations, number of instances, and evaluation metrics. Details of implemented 9 modules and involved metrics are in Appendix B and C, respectively.

Task	Analytical Target	Function	#Data	Metri
	Threat Attribution			
Malware Identification	Malware delivery or toolset	NER, SUM	15,742	F1
Signature Matching	Techniques from known threat groups	NER, SIM	5,166	F1
Temporal Pattern Matching	Known work schedules	REX	4,203	Sim
Affiliation Linking	Source organizations	NER, MAP	17,583	F1
Geographic Analysis	Geographic or cultural indicators	NER, SIM	6,164	F1
Victimology Profiling	Targeted victims or attacker motives	NER, REX	18,612	F1
Infrastructure Extraction	Domains, IPs, URLs, or file hashes	NER, REX, SUM	24,129	F1
Actor Identification	The threat group or actor (e.g., APT28)	NER, RAG, MAP	17,823	F1
Campaign Correlation	Threat campaigns or incidents	NER, MAP	27,762	F1
	Behavior Analysis			
File System Activity Detection	Suspicious file creation, deletion, or access	SPA, NER, SUM	4,653	Sim
Network Behavior Profiling	Patterns of external communication (e.g., C2)	SPA, NER, SUM	2,617	Sim
Credential Access Detection	Theft or misuse of credentials	SPA, NER, SUM	2,492	Sim
Execution Context Analysis	Execution behaviors by user or process	SPA, NER, SUM	23,888	Sim
Command & Script Analysis	Suspicious commands or scripts	SPA, NER, SUM	20,232	F1
Privilege Escalation Inference	Privilege escalation attempts	SPA, NER, SUM	15,953	Sim
Evasion Behavior Detection	Evasion or obfuscation techniques	SPA, NER, SUM	8,973	Sim
Event Sequence Reconstruction	Timeline of attack-related events	SUM	23,265	Sim
TTP Extraction	Tactics, techniques, and procedures	RAG, MAP	28,292	F1
	Prioritization			
Attack Vector Classification	Exploitation vectors (e.g., network, local, physical)	SUM, CLS	17,448	Acc
Attack Complexity Classification	Level of hurdles required to carry out the attack	SUM, CLS	17,116	Acc
Privileges Requirement Detection	Level of access privileges an attacker needs	SUM, CLS	18,030	Acc
User Interaction Categorization	If exploitation requires user participation	SUM, CLS	17,075	Acc
Attack Scope Detection	If the vulnerability affects one/multiple components	SUM, CLS	18,570	Acc
Impact Level Classification	Consequences on confidentiality, integrity, and availability	SUM, CLS	18,736	Acc
Severity Scoring	A numerical score indicating the overall attack severity	SUM, MATH	11,507	Dist
	Response & Mitigation			
Playbook Recommendation	Relevant response actions based on threat type	RAG, SUM	10,718	Hit
Security Control Adjustment	Firewall rules, EDR settings, or group policies	RAG, SUM	9,929	Sim
Patch Code Generation	Code snippets to patch the vulnerability	RAG, SUM	11,341	Pass
Patch Tool Suggestion	Security tools or utilities	RAG, SUM	9,763	Hit
Advisory Correlation	Security advisories or best practices	RAG, SUM	24,511	Hit

Subsequently, **Behavior Analysis** focuses on understanding how adversaries interact with systems over time. Tasks in this category include mapping unusual file system activities, profiling network behaviors (e.g., Monitoring outbound traffic), detecting credential access, and analyzing the use of commands and scripts. Analysts aim to reconstruct sequences of attack events and associate them with specific execution contexts or behavioral patterns.

When multiple threats emerge simultaneously, **Prioritization** assesses their relative urgency and associated risk. This involves analyzing the attack vector and complexity, identifying privilege requirements and user interaction dependencies, and estimating potential impact. These factors are then synthesized into impact labels and severity scores (e.g., CVSS (FIRST, a)) to guide effective triage. Finally, **Response & Mitigation** focus on generating actionable defense strategies. This includes recommending response playbooks, generating patch code, correlating relevant security advisories, and suggesting appropriate tools or configuration changes to neutralize the threat.

3.2 Data Sources

CYBERTEAM collects threat metadata from two primary sources: (1) vulnerability databases, which offer authoritative structural and non-structural information about threats, and (2) threat intelligence platforms, which report event-driven, context-rich threat data.

Vulnerability databases serve as foundational resources for automated threat hunting by providing machine-readable records of software flaws, exposure types, and critical contextual metadata. We aggregate threat entries from established sources such as NVD (National Institute of Standards and Technology (NIST), 2024), MITRE CVE (The MITRE Corporation, n.d.), ATT&CK (MITRE Corporation, 2024), CWE (MITRE Corporation, b), CAPEC (MITRE Corporation, a), D3FEND (MITRE Corporation, c), Exploit-DB (Offensive Security, 2024), and VulDB (VulDB Team, 2024).

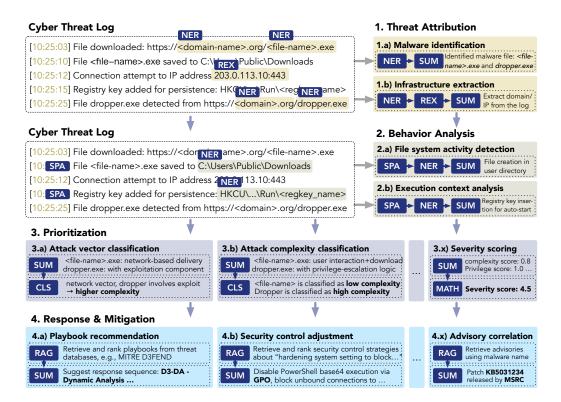


Figure 2: A threat hunting example demonstrating a dependency chain of analytical tasks, where each task is completed through a sequence of operational modules executed by LLMs autonomously.

These sources include detailed insights such as exploitability scores (EPSS (FIRST, b)), severity metrics (CVSS (FIRST, a)), and remediation guidance. Additionally, we incorporate data from vendor-maintained repositories (e.g., the Microsoft Security Update Guide (Microsoft Corporation), IBM X-Force (IBM Corporation, 2024)) to capture fine-grained details on affected systems, attack vectors, and patch methods.

Threat intelligence platforms complement these databases by providing behavioral and contextual signals linked to adversary activity. Platforms such as VirusTotal (VirusTotal (Google Chronicle), 2024), AlienVault OTX (AlienVault (AT&T Cybersecurity), 2024), and MISP (MISP Project, 2024) contribute indicators of compromise (IOCs), behavioral patterns, and telemetry that enable tasks like campaign correlation, infrastructure extraction, and actor attribution. Furthermore, industry threat reports—from sources, such as Mandiant (Mandiant (Google Cloud), 2024), Recorded Future (Recorded Future, 2024), Palo Alto Unit 42 (Palo Alto Networks, 2024), and Apache (The Apache Software Foundation, 2024), offer semi-structured intelligence, including incident timelines, IOC lists, and narrative analyses, which are essential for modeling multi-stage attack sequences and evaluating blue team responses.

Additional details on how these databases and platforms are used are provided in Appendix A.

3.3 STANDARDIZED THREAT HUNTING WITH OPERATIONAL MODULES

Task Dependency. Threat hunting is inherently a multi-stage analytical process (Sauerwein et al., 2019; Caltagirone et al., 2013; Hillier & Karroubi, 2022), where downstream actions, such as incident response and mitigation, rely on outcomes derived from upstream analytical steps. For example, recommending an effective response playbook requires accurate attribution of the threat actor and thorough behavioral analysis of the compromise. To explicitly model this structured workflow, CYBERTEAM formulates threat hunting as a *Dependency Chain*. As illustrated in Figure 2, all analytical tasks (e.g., 1.a: Malware Identification or 2.a: File System Activity Detection) are organized into a pipelined workflow that reflects their inherent dependencies. For example,

attack complexity classification relies on prior analyses of file system activity and execution context. Meanwhile, tasks within the same category (e.g., malware identification and infrastructure extraction under threat attribution) can often be performed in parallel, as they address distinct dimensions of the threat and do not exhibit direct interdependencies.

Highlight \(\varphi\). Instead of enumerating all tasks listed in Table 2, LLMs are asked to determine which tasks to perform at each stage, opening reasoning flexibility in threat hunting.

Operational Modules. Within each threat hunting task, CYBERTEAM invokes a set of operations (functional modules) designed to produce actionable threat analyses and progressively address the current threat hunting target (e.g., incident response). Specifically, each threat hunting task t_i is associated with a corresponding set of operational modules $\mathcal{F}_i = \{f_i^1, f_i^2, \ldots\}$. Each task t_i involves executing a sequence $f_i^* \in \mathcal{F}_i$, as detailed in the third column of Table 2. The resulting output $y_i = f_i^*(x)$ is subsequently passed to dependent downstream tasks. For instance, the task of TTP Extraction involves invoking both Retrieval-Augmented Generation (RAG) and Mapping (MAP) functions to identify relevant tactics and techniques from unstructured logs. Subsequently, a downstream task such as Tool Suggestion utilizes RAG and summarization (SUM) functions to map these identified TTPs to suitable defensive tools.

Highlight Q. These modules provide broad coverage of threat hunting practices (as shown in Table 2), while retaining flexibility (e.g., in SUM, RAG) for LLM reasoning to adapt across diverse scenarios, thereby balancing flexibility with standardization in blue team threat hunting.

Due to space constraints, we defer implementation details and design rationales to Appendix B.

4 EXPERIMENT

CYBERTEAM aims to empirically address the following research questions: \mathbf{RQ}_1 : How effective is standardization compared to open-ended reasoning for threat-hunting tasks? \mathbf{RQ}_2 : Can LLMs accurately solve individual threat-hunting tasks? \mathbf{RQ}_3 : How robust are LLMs, under the guidance of CYBERTEAM, when analyzing noisy inputs?

LLMs. We evaluate a range of industry-leading large language models, including GPT-40 (G40), GPT-04-mini (G04), Qwen3-32B (QW), Gemini-2.5 (GM), Claude-Sonnet-4 (CD), Llama-3.1-405B (L3.1), Llama-4-Scout-17B (L4), and Gemma-3-27b (GA). In addition, we assess state-of-the-art cybersecurity-focused LLM agents, including Lily-Cybersecurity-7B (LY) (Segolily Labs, 2025), DeepHat-7B (DH) (DeepHat, 2025), and SevenLLM-7B (SL) (Ji et al., 2024).

Open-ended Reasoning. In open-ended reasoning, we consider three widely used prompting structures: (1) In-Context Learning (ICL) (Dong et al., 2022) – including basic task instructions along with five (or ten) illustrative examples to demonstrate the desired solution format. (2) Chain-of-Thought (CoT) (Wei et al., 2022) – encouraging the model to generate "step-by-step" reasoning results before producing the final answer. (3) Tree-of-Thought (ToT) (Yao et al., 2023) – guiding LLMs to explore multiple reasoning paths and select the most plausible one.

Metrics. Table 2 lists evaluation metrics tailored to each task. For information extraction tasks (e.g., malware identification), we use the **F1 score** to balance precision and recall. For classification tasks (e.g., privilege escalation inference), we adopt **accuracy** among well-defined categories. Generation or summarization tasks (e.g., behavioral profiling) are evaluated using **BERTScore** (Zhang* et al., 2020) to measure semantic similarity. Tasks involving ranking (e.g., security playbook recommendation) utilize $\mathbf{Hit@k}$ (default k=10), measuring whether correct choices appear in the top-k outputs. For programmatic outputs (e.g., patch code generation), we apply **Pass** rate using UNITEST in Python to assess functional correctness. Numeric estimation tasks (e.g., severity scoring) are evaluated using **Normalized Distance** to quantify similarity to ground truth values. All metrics are scaled to the range [0, 1]. We explain the rationale for those metrics in Appendix C.

4.1 Standardized Threat Hunting vs. Open-Ended Reasoning (RQ_1)

Ultimately, CYBERTEAM is designed to generate actionable responses and mitigation strategies against cyber threats. We begin by evaluating the overall quality of LLM-generated responses and

Table 3: Results of LLMs threat-hunting performance (scaled to 100%) on CYBERTEAM, using corresponding metrics tailored to each analytical target as detailed in Table 2. We use **boldface** to indicate the best results and <u>underline</u> to denote the second-best results.

Metho	d	Cyber	security	Agent	Industry-Leading LLM								
1120110	-	LY	DH	SL	G4o	Go4	QW	GM	CD	L3.1	L4	GA	
	Playbook Recommend												
Open-ended	ICL5 ICL10 CoT ToT	42.3 44.1 51.6 48.1	54.2 52.5 50.6 53.3	54.7 55.3 50.5 54.3	64.5 65.2 78.3 75.2	73.1 74.5 89.2 85.1	52.8 53.6 67.5 71.4	79.4 80.2 80.1 83.5	63.7 64.9 <u>81.4</u> 77.2	65.8 66.4 77.3 82.1	55.8 56.4 67.3 72.1	54.9 55.5 66.4 71.2	
Standardized		67.2	58.4	66.8	84.6	91.4	71.4 79.3	91.8	89.3	89.7	79.7	78.8	
				Sec	urity Co	ntrol Adj	ust						
Open-ended	ICL5 ICL10 CoT ToT	51.5 53.2 60.3 66.7	66.3 68.4 70.5 <u>72.1</u>	43.9 45.6 <u>68.4</u> 61.6	61.8 62.7 70.3 75.9	70.3 71.8 80.2 85.6	50.6 51.2 59.8 66.3	65.8 66.4 79.2 73.6	79.2 80.1 77.2 73.1	61.5 62.3 77.9 72.8	51.5 52.3 <u>67.9</u> 62.8	50.6 51.4 <u>63.0</u> 61.9	
Standardized	l (Ours)	74.2	77.6	80.1	82.1	89.7	74.7	88.5	86.5	86.4	76.4	75.5	
Patch Code Generation													
Open-ended	ICL5 ICL10 CoT ToT	10.8 12.6 24.5 25.3	49.8 51.2 <u>54.7</u> 50.9	29.2 31.5 55.1 58.3	56.2 57.8 58.4 61.8	58.4 59.1 <u>76.3</u> 72.5	39.3 40.1 <u>54.7</u> 50.2	63.7 64.9 65.3 69.8	47.5 48.6 66.3 61.4	49.2 50.1 <u>67.4</u> 62.9	39.2 40.1 <u>57.4</u> 52.9	38.3 39.2 51.5 52.2	
Standardized		29.7	63.4	60.2	72.5	87.4	65.4	82.6	79.2	80.6	70.6	69.7	
				Pat	tch Tool	Suggestic	on						
Open-ended	ICL5 ICL10 CoT ToT	48.2 49.1 53.6 56.5	65.2 64.7 70.1 71.8	61.5 63.1 <u>77.2</u> 68.1	68.9 69.7 79.2 75.8	79.4 80.6 <u>90.1</u> 86.3	59.2 60.3 70.3 74.5	74.1 74.9 81.7 86.3	68.5 69.8 79.1 83.7	70.3 71.4 79.6 84.2	60.3 61.4 69.6 74.2	59.4 60.5 <u>68.7</u> 67.3	
Standardized	l (Ours)	69.1	76.5	77.7	87.4	96.9	83.6	93.2	91.2	92.1	82.1	81.2	
Advisory Correlation													
Open-ended Standardized	ICL5 ICL10 CoT ToT I (Ours)	21.7 22.9 49.5 46.8 73.4	57.5 59.1 71.4 73.2 78.8	63.8 64.7 69.5 67.2 77.1	64.7 65.9 67.2 70.8 80.3	67.2 68.1 80.5 84.2 92.3	48.5 49.2 61.7 64.8 76.5	62.4 63.2 77.5 73.1 86.9	56.8 58.1 76.2 72.5 84.5	58.7 59.5 76.3 71.8 84.9	48.7 49.5 66.3 61.8 74.9	47.8 48.6 65.4 60.9 74.0	

mitigation outputs on CYBERTEAM. Table 3 presents the results, using task-specific metrics detailed in Table 2.

Effectiveness of Standardization. From Table 3, we observe that using operational modules (**Ours**) outperforms typical open-ended reasoning methods. For instance, modular operations enable GPT-o4 to achieve over 90% Hit@10 in playbook recommendation and over 92% in advisory correlation. In contrast, open-ended reasoning achieves only secondary effectiveness, with a significant performance gap observed (e.g., in the security control adjustment task of SevenLLM). This demonstrates the effectiveness of combining standardized guidance with the inherent flexibility of LLMs.

Gains from Standard Operating Procedures. Notably, while ICL, CoT, and ToT have been shown to improve generation quality for general-purpose tasks (Dong et al., 2022; Yu et al., 2023; Wang et al., 2022), they provide limited guidance for domain-specific problems that require precise procedural knowledge and structured analytical workflows. By contrast, standardized threat hunting workflows help LLMs follow standard operating procedures by decomposing complex tasks into modular steps. This reduces hallucination and enforces structure. In tasks requiring strict sequencing (e.g., threat actor identification followed by response planning), workflow-based methods ensure the correct order and information flow, outperforming ICL, CoT, and ToT, which often lack such control.

Case Study I (Failure Case). When using CoT to generate a response plan for LockBit (a ransomware), GPT-40 offers generic recommendations "... the first step is to isolate affected machines. Next, the system should assess backup availability and notify stakeholders ..." without tailoring to LockBit and ignoring unique traits like double extortion tactics or known exploits.

By contrast, operations in CYBERTEAM constrain LLM reasoning to resolve correct analytical sequences, ensuring outputs remain aligned with operational goals:

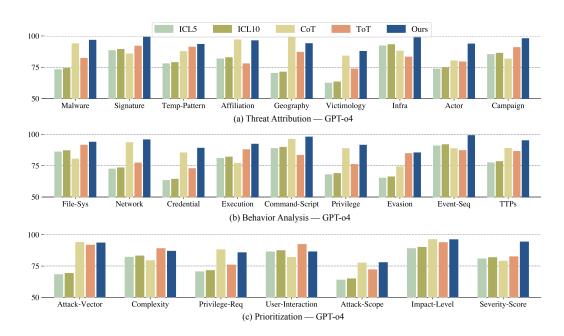


Figure 3: Threat-hunting performance (scaled to 100%) on individual tasks, evaluating under GPT-o4-mini. Results for additional LLMs are provided in Appendix E.

Case Study II (Successful Case). The modular operation framework guides GPT-40 to explicitly invoke **RAG** and **SUM** modules. Specifically, RAG retrieves up-to-date security advisories (e.g., *CISA Alert AA23-325A*) specific to LockBit, while SUM outlines mitigation strategies with *double extortion prevention* and *air-gapped offline backups*.

These results suggest that in cybersecurity, particularly in threat-hunting scenarios, structured elicitation methods are necessary for reliably leveraging LLM capabilities.

Operational Interpretability. Notably, the modular approach enhances interpretability for analysts, as outputs can be traced back to specific operations (e.g., RAG for evidence retrieval, SUM for summarization). In contrast, open-ended prompts produce opaque reasoning chains that are harder to audit what real-world evidence is integrated.

Case Study III (Interpretability). For the MOVEit vulnerability (CVE-2023-34362), an openended Qwen prompt returned only a vague recommendation ("apply vendor patches and monitor suspicious traffic"). In contrast, our pipeline invoked the **RAG** module to retrieve Progress Software's advisory and the **NER** module to extract SQL injection IOCs. This modular trace improved accuracy and enabled analysts to audit advisory steps.

Due to space constraints, we provide additional evaluation of the trade-off between latency and reliability in Appendix E.1. Our results show that the standardized threat hunting method achieves a more favorable balance compared with open-ended reasoning.

Design Insights \bigcirc **.** The evaluation provides two actionable insights for blue team practices: (1) Breaking threat analysis into smaller, modularized operations (e.g., IOC extraction, TTP mapping), each guided by distinct reasoning objectives; (2) Integrate LLMs into existing analytic pipelines where upstream outputs (e.g., extracted indicators) are fed into downstream modules rather than relying on single-pass generation.

4.2 THREAT-HUNTING PERFORMANCE FOR INDIVIDUAL TASKS (RQ₂)

Complementing Section 4.1, we also evaluate individual threat-hunting tasks prior to the response & mitigation stage, as outlined in Table 2. Figures 3 and Appendix E present the experimental results.

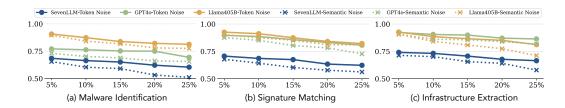


Figure 4: LLM performance (metrics corresponding to Table 2) when input threat logs are perturbed with token-level noise (solid line) or semantic-level noise (dashed line). X-axis shows the noise ratios.

Observe that using standardized threat hunting consistently achieves the highest performance across all intermediate tasks. However, **the magnitude of performance gains varies across task types**. For instance, in complex reasoning tasks (e.g., Event Sequence Construction), the standardized method yields substantial improvements over open-ended reasoning strategies like CoT and ToT, boosting accuracy by over 20% using GPT-o4-mini. These gains are most notable when task dependencies are strong. For example, generating effective responses depends on accurate upstream analysis. Module-guided models can preserve and pass critical context, while ICL/CoT/ToT often fail to coordinate such multi-stage reasoning reliably. This is largely because these tasks require multi-hop reasoning, evidence synthesis, and careful dependency tracking, which are capabilities that general prompting methods struggle to coordinate effectively. In contrast, for narrower, classification-focused tasks (e.g., attack vector categorization or privilege escalation inference), the performance gap between operational modules and standard prompting is smaller. Here, the tasks are more self-contained, and models can often arrive at correct predictions even without explicit task decomposition or function integration.

Design Insights ?. While standardized threat hunting offer general advantages, their relative benefit is particularly significant in **scenarios requiring structured reasoning** over interconnected steps. This demonstrates the importance of modular guidance in complex cybersecurity workflows.

Due to space constraints, we provide complementary results and analyses in Appendix E.2.

4.3 LLM Robustness against Noisy Inputs (RQ₃)

Experimental Setting. We also investigate LLM robustness when input threat logs contain noisy text. We introduce (i) token-level noise using TextAttack (Morris et al., 2020), which randomly injects or substitutes tokens, and (ii) semantic-level noise using BART-paraphraser (Lewis et al., 2019), which subtly introduces misleading or shifted context. Both noise types are applied at controlled levels (e.g., perturbing 10% of the input).

Results and Observations. From Figure 4, we observe that token-level noise has a smaller impact on LLM performance compared to semantic-level noise. For example, under 10% perturbation, random character insertions or deletions lead to less than 5% performance drop across tasks. In contrast, semantic-level noise (e.g., paraphrased or subtly altered context) causes a much larger decline. These findings suggest that while LLMs handle surface-level errors relatively well, they struggle with the semantic shifting, even when guided by CYBERTEAM. This highlights the importance of curating expert-level threat reports in threat hunting, as imprecise statements can unintentionally mislead blue team efforts and degrade overall analysis.

5 CONCLUSION

We present CYBERTEAM, a benchmark designed to evaluate the capabilities of LLMs in blue team threat-hunting workflows. By combining broad and diverse real-world datasets, a standardized workflow environment with modular function-guided reasoning, and detailed evaluation strategies, CYBERTEAM provides a comprehensive workflow for assessing LLM capabilities in realistic cyber defense scenarios. Our empirical findings offer actionable insights for integrating standardized operations into security workflows. We hope CYBERTEAM will serve as a valuable resource for the research community and practitioners alike, driving future innovations in AI-assisted cybersecurity.

ETHICS STATEMENT

This study is based solely on publicly accessible cybersecurity reports, vulnerability databases, and open-source intelligence platforms, each used in accordance with their copyright and licensing conditions. No proprietary, sensitive, or personally identifiable data were collected or processed.

REPRODUCIBILITY STATEMENT

To ensure reproducibility, we provide an anonymous GitHub repository containing benchmark construction details, operational module implementations, evaluation pipelines, and experiment configurations.

REFERENCES

- Tsedeke Abate, Kassa Michael, and Carl Angell. Assessment of scientific reasoning: Development and validation of scientific reasoning assessment tool. *Eurasia Journal of Mathematics, Science and Technology Education*, 16(12):em1927, 2020.
- Adel Abusitta, Miles Q Li, and Benjamin CM Fung. Malware classification and composition analysis: A survey of recent developments. *Journal of Information Security and Applications*, 59:102828, 2021.
- Ehsan Aghaei, Waseem Shadid, and Ehab Al-Shaer. Threatzoom: Cve2cwe using hierarchical neural network. *arXiv preprint arXiv:2009.11501*, 2020.
- Sharath Chandra Akkaladevi, Matthias Plasch, Michael Hofmann, and Andreas Pichler. Semantic knowledge based reasoning framework for human robot collaboration. *Procedia CIRP*, 97:373–378, 2021.
- Jamal Al-Karaki, Muhammad Al-Zafar Khan, and Marwan Omar. Exploring llms for malware detection: Review, framework design, and countermeasure approaches. *arXiv preprint arXiv:2409.07587*, 2024.
- Md Tanvirul Alam, Dipkamal Bhusal, Le Nguyen, and Nidhi Rastogi. Ctibench: A benchmark for evaluating llms in cyber threat intelligence. *arXiv preprint arXiv:2406.07599*, 2024.
- AlienVault (AT&T Cybersecurity). Alienvault open threat exchange (otx). https://otx.alienvault.com, 2024.
- Mario Luca Bernardi, Marta Cimitile, and Riccardo Pecori. Automatic job safety report generation using rag-based llms. In *2024 International Joint Conference on Neural Networks (IJCNN)*, pp. 1–8. IEEE, 2024.
- Gavin Black, Varghese Vaidyan, and Gurcan Comert. Evaluating large language models for enhanced fuzzing: An analysis framework for llm-driven seed generation. *IEEE Access*, 2024.
- Alexandre Braga, Ricardo Dahab, Nuno Antunes, Nuno Laranjeiro, and Marco Vieira. Practical evaluation of static analysis tools for cryptography: Benchmarking method and case study. In 2017 IEEE 28th International Symposium on Software Reliability Engineering (ISSRE), pp. 170–181. IEEE, 2017.
- Ioana Branescu, Octavian Grigorescu, and Mihai Dascalu. Automated mapping of common vulnerabilities and exposures to mitre att&ck tactics. *Information*, 15(4):214, 2024.
- Lee Brotherston, Amanda Berlin, and William F Reyor III. *Defensive security handbook*. "O'Reilly Media, Inc.", 2024.
- Sergio Caltagirone, Andrew Pendergast, and Christopher Betz. The diamond model of intrusion analysis. *Threat Connect*, 298(0704):1–61, 2013.
 - Zhili Cheng, Yuge Tu, Ran Li, Shiqi Dai, Jinyi Hu, Shengding Hu, Jiahao Li, Yang Shi, Tianyu Yu, Weize Chen, et al. Embodiedeval: Evaluate multimodal llms as embodied agents. *arXiv preprint arXiv:2501.11858*, 2025.

- Leshem Choshen, Ariel Gera, Yotam Perlitz, Michal Shmueli-Scheuer, and Gabriel Stanovsky.
 Navigating the modern evaluation landscape: Considerations in benchmarks and frameworks for large language models (Ilms). In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024): Tutorial Summaries*, pp. 19–25, 2024.
 - CISE Program. Cybersecurity information sharing environment (cise), 2024. URL https://www.cisa.gov/cybersecurity-information-sharing.
 - Ondrej Čupka, Ester Federlova, and Peter Vesely. Comparison of methodologies used in cybersecurity reports. In *Developments in Information and Knowledge Management Systems for Business Applications: Volume 7*, pp. 313–348. Springer, 2023.
 - Hafsa Shareef Dar, M Ikramullah Lali, Moin Ul Din, Khalid Mahmood Malik, and Syed Ahmad Chan Bukhari. Frameworks for querying databases using natural language: a literature review. *arXiv preprint arXiv:1909.01822*, 2019.
 - Tirtharaj Dash, Sharad Chitlangia, Aditya Ahuja, and Ashwin Srinivasan. A review of some techniques for inclusion of domain-knowledge into deep neural networks. *Scientific Reports*, 12(1): 1040, 2022.
 - DeepHat. Deephat-v1-7b. Model on Hugging Face, 2025.
 - Gelei Deng, Yi Liu, Víctor Mayoral-Vilches, Peng Liu, Yuekang Li, Yuan Xu, Tianwei Zhang, Yang Liu, Martin Pinzger, and Stefan Rass. Pentestgpt: An Ilm-empowered automatic penetration testing tool. *arXiv preprint arXiv:2308.06782*, 2023.
 - Gelei Deng, Yi Liu, Víctor Mayoral-Vilches, Peng Liu, Yuekang Li, Yuan Xu, Tianwei Zhang, Yang Liu, Martin Pinzger, and Stefan Rass. {PentestGPT}: Evaluating and harnessing large language models for automated penetration testing. In *33rd USENIX Security Symposium (USENIX Security 24)*, pp. 847–864, 2024.
 - Dharani Devadiga, Gordon Jin, Bisti Potdar, Hankyu Koo, Andrew Han, Anusha Shringi, Angad Singh, Kinjal Chaudhari, and Saurav Kumar. Gleam: Gan and Ilm for evasive adversarial malware. In 2023 14th International Conference on Information and Communication Technology Convergence (ICTC), pp. 53–58. IEEE, 2023.
 - Yuri Diogenes and Erdal Ozkaya. Cybersecurity-attack and defense strategies: Infrastructure security with red team and blue team tactics. Packt Publishing Ltd, 2018.
 - Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Jingyuan Ma, Rui Li, Heming Xia, Jingjing Xu, Zhiyong Wu, Tianyu Liu, et al. A survey on in-context learning. *arXiv preprint arXiv:2301.00234*, 2022.
 - Vardhan Dongre, Xiaocheng Yang, Emre Can Acikgoz, Suvodip Dey, Gokhan Tur, and Dilek Hakkani-Tür. Respact: Harmonizing reasoning, speaking, and acting towards building large language model-based conversational ai agents. *arXiv preprint arXiv:2411.00927*, 2024.
 - Danny Driess, Fei Xia, Mehdi SM Sajjadi, Corey Lynch, Aakanksha Chowdhery, Ayzaan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, Wenlong Huang, et al. Palm-e: An embodied multimodal language model. 2023.
 - Hossein Rajaby Faghihi, Aliakbar Nafar, Chen Zheng, Roshanak Mirzaee, Yue Zhang, Andrzej Uszok, Alexander Wan, Tanawan Premsri, Dan Roth, and Parisa Kordjamshidi. Gluecons: A generic benchmark for learning under constraints. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pp. 9552–9561, 2023.
 - Zhiwei Fei, Xiaoyu Shen, Dawei Zhu, Fengzhe Zhou, Zhuo Han, Songyang Zhang, Kai Chen, Zongwen Shen, and Jidong Ge. Lawbench: Benchmarking legal knowledge of large language models. *arXiv preprint arXiv:2309.16289*, 2023.
 - FIRST. Common vulnerability scoring system (cvss). https://www.first.org/cvss/, a.
 - FIRST. Exploit prediction scoring system (epss). https://www.first.org/epss/, b.

- Francesco Greco, Giuseppe Desolda, Andrea Esposito, Alessandro Carelli, et al. David versus goliath: Can machine learning detect llm-generated text? a case study in the detection of phishing emails. In *The Italian Conference on CyberSecurity*, 2024.
- Andreas Happe and Jürgen Cito. Getting pwn'd by ai: Penetration testing with large language models. In *Proceedings of the 31st ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, pp. 2082–2086, 2023.
- Erik Hemberg, Jonathan Kelly, Michal Shlapentokh-Rothman, Bryn Reinstadler, Katherine Xu, Nick Rutar, and Una-May O'Reilly. Linking threat tactics, techniques, and patterns with defensive weaknesses, vulnerabilities and affected platform configurations for cyber hunting. *arXiv* preprint arXiv:2010.00533, 2020.
- Juan R Bermejo Higuera, Javier Bermejo Higuera, Juan A Sicilia Montalvo, Javier Cubo Villalba, and Juan José Nombela Pérez. Benchmarking approach to compare web applications static analysis tools detecting owasp top ten security vulnerabilities. *Computers, Materials & Continua*, 64(3), 2020.
- Caroline Hillier and Talieh Karroubi. Turning the hunted into the hunter via threat hunting: Life cycle, ecosystem, challenges and the great promise of ai. arXiv preprint arXiv:2204.11076, 2022.
- Md Imran Hossen, Jianyi Zhang, Yinzhi Cao, and Xiali Hei. Assessing cybersecurity vulnerabilities in code large language models. *arXiv preprint arXiv:2404.18567*, 2024.
- Mengkang Hu, Pu Zhao, Can Xu, Qingfeng Sun, Jianguang Lou, Qingwei Lin, Ping Luo, and Saravan Rajmohan. Agentgen: Enhancing planning abilities for large language model based agent via environment and task generation. *arXiv preprint arXiv:2408.00764*, 2024.
- IBM Corporation. Ibm x-force exchange, 2024. URL https://exchange.xforce.ibmcloud.com/.
- Hyeongyo Jeong, Haechan Lee, Changwon Kim, and Sungtae Shin. A survey of robot intelligence with large language models. *Applied Sciences*, 14(19):8868, 2024.
- Hangyuan Ji, Jian Yang, Linzheng Chai, Chaoren Wei, Liqun Yang, Yunlong Duan, Yunli Wang, Tianzhen Sun, Hongcheng Guo, Tongliang Li, et al. Sevenllm: Benchmarking, eliciting, and enhancing abilities of large language models in cyber threat intelligence. *arXiv preprint arXiv:2405.03446*, 2024.
- Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik Narasimhan. Swe-bench: Can language models resolve real-world github issues? *arXiv preprint arXiv:2310.06770*, 2023.
- Rabiah Abdul Kadir, Ely Salwana Mat Surin, and Mahidur R Sarker. A systematic review of automated classification for simple and complex query sql on nosql database. *Computer Systems Science & Engineering*, 48(6), 2024.
- Aditya Kulkarni, Vivek Balachandran, Dinil Mon Divakaran, and Tamal Das. From ml to llm: Evaluating the robustness of phishing webpage detection models against adversarial attacks. *arXiv preprint arXiv:2407.20361*, 2024.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Ves Stoyanov, and Luke Zettlemoyer. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension, 2019.
- Ziyang Li, Saikat Dutta, and Mayur Naik. Iris: Llm-assisted static analysis for detecting security vulnerabilities. In *The Thirteenth International Conference on Learning Representations*, 2025.
- Guilong Lu, Xiaolin Ju, Xiang Chen, Wenlong Pei, and Zhilong Cai. Grace: Empowering llm-based software vulnerability detection with graph structure and in-context learning. *Journal of Systems and Software*, 212:112031, 2024.

- Mary M Lucas, Justin Yang, Jon K Pomeroy, and Christopher C Yang. Reasoning with large language models for medical question answering. *Journal of the American Medical Informatics Association*, 31(9):1964–1975, 2024.
 - Mandiant (Google Cloud). Mandiant threat intelligence reports. https://www.mandiant.com/resources/reports, 2024.
 - Sean McGregor, Allyson Ettinger, Nick Judd, Paul Albee, Liwei Jiang, Kavel Rao, William H Smith, Shayne Longpre, Avijit Ghosh, Christopher Fiorelli, et al. To err is ai: A case study informing llm flaw reporting practices. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pp. 28938–28945, 2025.
 - Microsoft Corporation. Microsoft security update guide. https://msrc.microsoft.com/update-guide.
 - MISP Project. Misp threat intelligence sharing platform. https://www.misp-project.org, 2024.
 - MITRE Corporation. Common attack pattern enumeration and classification (capec). https://capec.mitre.org/, a.
 - MITRE Corporation. Common weakness enumeration (cwe). https://cwe.mitre.org/, b.
 - MITRE Corporation. D3fend: A knowledge graph of cybersecurity countermeasures. https://d3fend.mitre.org/, c.
 - MITRE Corporation. Mitre att&ck framework, 2024. URL https://attack.mitre.org/.
 - John Morris, Eli Lifland, Jin Yong Yoo, Jake Grigsby, Di Jin, and Yanjun Qi. Textattack: A framework for adversarial attacks, data augmentation, and adversarial training in nlp. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pp. 119–126, 2020.
 - Niels Mündler, Jingxuan He, Slobodan Jenko, and Martin Vechev. Self-contradictory hallucinations of large language models: Evaluation, detection and mitigation. *arXiv preprint arXiv:2305.15852*, 2023.
 - Lajos Muzsai, David Imolai, and András Lukács. Hacksynth: Llm agent and evaluation framework for autonomous penetration testing. *arXiv preprint arXiv:2412.01778*, 2024.
 - National Institute of Standards and Technology (NIST). National vulnerability database (nvd), 2024. URL https://nvd.nist.gov/.
 - Offensive Security. Exploit database (exploit-db), 2024. URL https://www.exploit-db.com/.
 - Yaroslav Oliinyk, Michael Scott, Ryan Tsang, Chongzhou Fang, Houman Homayoun, et al. Fuzzing {BusyBox}: Leveraging {LLM} and crash reuse for embedded bug unearthing. In *33rd USENIX Security Symposium (USENIX Security 24)*, pp. 883–900, 2024.
 - Oracle Corporation. Oracle security alerts, 2024. URL https://www.oracle.com/security-alerts/.
 - Palo Alto Networks. Unit 42 threat research reports. https://unit42.paloaltonetworks.com, 2024.
 - Filippo Perrina, Francesco Marchiori, Mauro Conti, and Nino Vincenzo Verde. Agir: Automating cyber threat intelligence reporting with natural language generation. In *2023 IEEE International Conference on Big Data (BigData)*, pp. 3053–3062. IEEE, 2023.
 - Xingzhi Qian, Xinran Zheng, Yiling He, Shuo Yang, and Lorenzo Cavallaro. Lamd: Context-driven android malware detection and classification with llms. *arXiv* preprint arXiv:2502.13055, 2025.

- Jeyavijayan Rajendran, Vinayaka Jyothi, and Ramesh Karri. Blue team red team approach to hardware trust assessment. In 2011 IEEE 29th international conference on computer design (ICCD), pp. 285–288. IEEE, 2011.
 - Recorded Future. Recorded future threat intelligence reports. https://www.recordedfuture.com/research, 2024.
 - Red Hat, Inc. Red hat security advisories (rhsa). https://access.redhat.com/.
 - Red Hat, Inc. Red hat bugzilla, 2024. URL https://buqzilla.redhat.com/.
 - Ann Marie Reinhold, Brittany Boles, A Redempta Manzi Muneza, Thomas McElroy, and Clemente Izurieta. Surmounting challenges in aggregating results from static analysis tools. *Military Cyber Affairs*, 7(1):5–11, 2024.
 - Rebecca Russell, Louis Kim, Lei Hamilton, Tomo Lazovich, Jacob Harer, Onur Ozdemir, Paul Ellingwood, and Marc McConley. Automated vulnerability detection in source code using deep representation learning. In 2018 17th IEEE international conference on machine learning and applications (ICMLA), pp. 757–762. IEEE, 2018.
 - Clemens Sauerwein, Christian Sillaber, Andrea Mussmann, and Ruth Breu. A framework for cyber threat hunting. In ACM CCS Workshop on Security and Privacy Analytics, 2019.
 - Devesh Sawant, Manjesh K Hanawal, and Atul Kabra. Improving discovery of known software vulnerability for enhanced cybersecurity. *arXiv preprint arXiv:2412.16607*, 2024.
 - Segolily Labs. Lily-Cybersecurity-7B-v0.2. https://huggingface.co/segolilylabs/Lily-Cybersecurity-7B-v0.2, 2025.
 - Kunal Sehgal and Nikolaos Thymianis. *Cybersecurity Blue Team Strategies: Uncover the secrets of blue teams to combat cyber threats in your organization*. Packt Publishing Ltd, 2023.
 - Xiangmin Shen, Lingzhi Wang, Zhenyuan Li, Yan Chen, Wencheng Zhao, Dawei Sun, Jiashui Wang, and Wei Ruan. Pentestagent: Incorporating llm agents to automated penetration testing. *arXiv* preprint arXiv:2411.05185, 2024.
 - Ze Sheng, Fenghua Wu, Xiangwu Zuo, Chao Li, Yuxin Qiao, and Lei Hang. Lprotector: An llm-driven vulnerability detection system. *arXiv preprint arXiv:2411.06493*, 2024.
 - Aryan Shrivastava. Response inconsistency of large language models in high-stakes military decision making.
 - Adi Simhi, Itay Itzhak, Fazl Barez, Gabriel Stanovsky, and Yonatan Belinkov. Trust me, i'm wrong: High-certainty hallucinations in llms. *arXiv preprint arXiv:2502.12964*, 2025.
 - Fahim Sufi. An innovative gpt-based open-source intelligence using historical cyber incident reports. *Natural Language Processing Journal*, 7:100074, 2024.
 - The Apache Software Foundation. Apache security advisories. https://www.apache.org/security/, 2024.
 - The MITRE Corporation. Common Vulnerabilities and Exposures (CVE). https://cve.mitre.org/, n.d.
 - Norbert Tihanyi, Mohamed Amine Ferrag, Ridhi Jain, Tamas Bisztray, and Merouane Debbah. Cybermetric: a benchmark dataset based on retrieval-augmented generation for evaluating llms in cybersecurity knowledge. In 2024 IEEE International Conference on Cyber Security and Resilience (CSR), pp. 296–302. IEEE, 2024.
 - Krist Vaesen and Wybo Houkes. A new framework for teaching scientific reasoning to students from application-oriented sciences. *European journal for philosophy of science*, 11(2):56, 2021.
- VirusTotal (Google Chronicle). Virustotal: Analyze suspicious files and urls. https://www.virustotal.com, 2024.

- VulDB Team. Vuldb vulnerability database, 2024. URL https://vuldb.com/.
 - Boshi Wang, Sewon Min, Xiang Deng, Jiaming Shen, You Wu, Luke Zettlemoyer, and Huan Sun. Towards understanding chain-of-thought prompting: An empirical study of what matters. *arXiv* preprint arXiv:2212.10001, 2022.
 - Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming Ren, Aaran Arulraj, Xuan He, Ziyan Jiang, et al. Mmlu-pro: A more robust and challenging multitask language understanding benchmark. In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2024.
 - Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.
 - Yue Yang, Fan-Yun Sun, Luca Weihs, Eli VanderBilt, Alvaro Herrasti, Winson Han, Jiajun Wu, Nick Haber, Ranjay Krishna, Lingjie Liu, et al. Holodeck: Language guided generation of 3d embodied ai environments. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 16227–16237, 2024.
 - Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. *Advances in neural information processing systems*, 36:11809–11822, 2023.
 - Zihan Yu, Liang He, Zhen Wu, Xinyu Dai, and Jiajun Chen. Towards better chain-of-thought prompting strategies: A survey. *arXiv preprint arXiv:2310.04959*, 2023.
 - Jie Zhang, Haoyu Bu, Hui Wen, Yongji Liu, Haiqiang Fei, Rongrong Xi, Lun Li, Yun Yang, Hongsong Zhu, and Dan Meng. When llms meet cybersecurity: A systematic literature review. *Cybersecurity*, 8(1):1–41, 2025.
 - Tianyi Zhang*, Varsha Kishore*, Felix Wu*, Kilian Q. Weinberger, and Yoav Artzi. Bertscore: Evaluating text generation with bert. In *International Conference on Learning Representations*, 2020.
 - Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. Bertscore: Evaluating text generation with bert. In *International Conference on Learning Representations (ICLR)*, 2020. URL https://arxiv.org/abs/1904.09675. Official implementation: https://github.com/Tiiiger/bert_score.
 - Yuxuan Zhou, Xien Liu, Chen Ning, Xiao Zhang, Chenwei Yan, Xiangling Fu, and Ji Wu. Revisiting the scaling effects of llms on medical reasoning capabilities.

A DATA SOURCE AND METADATA COLLECTION

The MITRE CVE (Common Vulnerabilities and Exposures) system (The MITRE Corporation, n.d.) is a foundational database that provides unique identifiers for publicly disclosed cybersecurity vulnerabilities. Each CVE record includes an ID, a brief description, references to external resources, and associated vendors or platforms. This source allows for consistent naming and indexing of vulnerabilities across tools and reports. We collect structured metadata such as CVE IDs, descriptions, reference links, and related CWE classifications. CVE feeds (XML/JSON) are used for automated ingestion and linkage to other threat intelligence frameworks like CAPEC and ATT&CK.

- Maintained by NIST, the **NVD** (**National Vulnerability Database**) (National Institute of Standards and Technology (NIST), 2024) builds on MITRE CVE data by adding rich metadata, including CVSS scores (base, temporal, environmental), CWE mappings, configuration impacts, patch availability, and severity vectors. We extract metadata through the official JSON data feeds, parsing CVE-level risk metrics, impact sub-scores, and associated product configurations. This information is critical for prioritizing remediation and understanding the real-world impact of vulnerabilities.
- **Exploit-DB** (Offensive Security, 2024) is a curated collection of publicly available exploits and proof-of-concept code. Each entry includes exploit titles, CVE references, author information, platform

tags, and the actual code used in attacks. Unlike CVE/NVD, Exploit-DB provides practical insights into how vulnerabilities are weaponized in real environments. We extract titles, descriptions, exploit types (e.g., Local, Remote), and related CVEs using web scraping and NLP-based text classification.

CWE(Common Weakness Enumeration) (MITRE Corporation, b) is a taxonomy developed by MITRE to classify software and hardware weaknesses. Each CWE includes a unique ID, a detailed explanation, potential consequences, examples, and related patterns (e.g., CAPEC). We use CWE to enrich CVE data with root cause information, enabling fine-grained vulnerability clustering and defensive prioritization. The metadata includes weakness category, severity, and relationships with CAPEC and CVE entries.

CAPEC (Common Attack Pattern Enumeration and Classification) (MITRE Corporation, a) provides a standardized catalog of common attack strategies. Each pattern includes the attacker's objectives, prerequisites, execution flow, related weaknesses (CWE), and example scenarios. We extract attack pattern IDs, descriptions, related CWEs, and suggested mitigations. These data points enable us to map vulnerabilities to adversarial behaviors, enhancing our CTI behavioral modeling capabilities.

The MITRE ATT&CK (MITRE Corporation, 2024) framework systematically catalogs adversary tactics, techniques, and procedures (TTPs) observed in real-world incidents. Each entry includes tactic categories (e.g., Privilege Escalation), techniques, mitigations, detection suggestions, and threat actor mappings. We extract technique IDs, corresponding software, mitigation strategies, and detection methods. These are used to link CVEs and exploits to higher-level attacker behaviors, supporting advanced threat modeling.

D3FEND (MITRE Corporation, c) is a curated knowledge graph that maps defensive techniques to specific threat behaviors and artifacts. D3FEND complements the well-known ATT&CK framework by focusing on how defenders can detect, disrupt, and respond to adversarial actions. To integrate this resource into CYBERTEAM, we crawl D3FEND's publicly available ontology and extract metadata on detection, deception, and mitigation techniques, along with their associated digital artifacts (e.g., file paths, registry keys, network signatures). This metadata is then linked to relevant analytical tasks, such as behavioral profiling and response planning, providing a rich, standardized reference for grounding LLM outputs in practical defensive actions.

Oracle Security Alerts (Oracle Corporation, 2024) provides detailed security patch advisories for its product suite. Each alert includes the CVEs addressed, severity scores, and remediation timelines. We parse the advisories to gather product-specific vulnerability timelines, vendor patch statuses, and mitigation instructions, which complement the NVD and MITRE CVE datasets.

Red Hat Bugzilla (Red Hat, Inc., 2024) is a bug tracking system that includes detailed discussions and technical logs about software bugs, many of which are security-related. Entries often include CVE links, fix status, patch availability, and affected components. We scrape metadata such as Bug IDs, CVE references, affected packages, and resolution details to supplement our understanding of vulnerability lifecycle management.

The RHSA(Red Hat Security Advisories) (Red Hat, Inc.) portal lists all critical, important, and moderate security advisories affecting Red Hat products. Each advisory provides CVE mappings, severity scores, fixed packages, and risk summaries. Metadata extraction includes advisory IDs, publication dates, CVE linkages, and suggested upgrades or patches, enabling alignment with real-world remediation practices.

IBM X-Force Exchange (IBM Corporation, 2024) is a commercial threat intelligence sharing platform that provides in-depth reports on vulnerabilities, exploits, malware, and threat actors. Each CVE entry is enriched with exploitability status, malware connections, and actor attribution. We extract structured threat metadata such as exploit availability, indicators of compromise (IOCs), campaign tags, and actor profiling to complement CVE risk modeling.

CISE (Cybersecurity Information Sharing Environment) (CISE Program, 2024), maintained by CISA, promotes cybersecurity information exchange across government and private sector entities. The platform facilitates sharing of indicators of compromise (IOCs), analysis reports, and threat mitigation strategies through structured partnerships. We extract strategic-level threat metadata, including threat vectors, vulnerability trends, and response best practices from shared reports and alerts. This supports broader CTI tasks like attribution and risk contextualization.

VulDB (Vulnerability Database) (VulDB Team, 2024) is a commercial vulnerability intelligence service that provides insights into current exploits, threat actor behavior, and exploit trends. Entries often include exploitability scores, attack vectors, exploitation status, and tags related to malware or campaigns. We collect CVE mappings, vulnerability titles, exploitation timelines, and associated actors, enabling temporal and behavioral correlation with other sources like Exploit-DB and MITRE ATT&CK.

Apache's official security advisory page lists all disclosed vulnerabilities affecting Apache projects (e.g., HTTP Server, Tomcat, Struts) (The Apache Software Foundation, 2024). Each advisory includes CVE references, affected versions, and patch instructions. We extract CVE mappings, patch details, vulnerability types, and affected modules. These insights are cross-referenced with MITRE CVE and NVD entries to improve accuracy in software-specific threat tracking.

Mandiant Threat Intelligence Reports (Mandiant (Google Cloud), 2024), now part of Google Cloud, publishes in-depth research on nation-state APTs, malware campaigns, and threat actor tactics. Their reports include IOC lists, ATT&CK mappings, and campaign chronologies. We extract metadata on APT groups, attack stages, observed TTPs, and malware toolkits. These data points support the attribution and behavioral modeling dimensions of our threat intelligence corpus.

Recorded Future Threat Intelligence Reports (Recorded Future, 2024) publishes real-time, machine-readable threat intelligence covering threat actors, vulnerabilities, dark web chatter, and geopolitical cyber campaigns. Reports often include structured indicators, predictive analytics, and CVE exploitability assessments. We leverage this source to collect threat context, emerging trends, and exploit discussion patterns—enabling our system to associate vulnerabilities with evolving threat actor intent and capability.

Unit 42 Threat Research (Palo Alto Networks) (Palo Alto Networks, 2024) provides malware analysis, campaign forensics, and actor behavior insights from Palo Alto Networks' global threat intelligence platform. Their publications include links to malicious infrastructure, malware families, and ATT&CK references. We extract TTPs, CVE-to-malware correlations, and campaign data. This enhances our contextual metadata for linking specific vulnerabilities to real-world exploitation scenarios.

Microsoft's Security Update Guide (Microsoft Corporation) lists monthly updates across its software stack. Entries contain CVEs, severity ratings, exploitability assessments, patch availability, and affected platforms. Metadata extraction includes CVE linkage, threat vectors (e.g., local, remote), exploitation likelihood, and patch rollout status—enriching vendor-specific vulnerability intelligence.

CVSS (Common Vulnerability Scoring System) (FIRST, a) is a widely adopted scoring system developed by FIRST to assess the severity of software vulnerabilities. It breaks down risk into Base, Temporal, and Environmental components. We use this framework to interpret NVD scores, compare severity across platforms, and calibrate exploitability in relation to business-critical systems.

EPSS (Exploit Prediction Scoring System) (FIRST, b), also developed by FIRST, provides probabilistic predictions of whether a vulnerability is likely to be exploited in the wild. It integrates data from CVSS, Exploit-DB, and historical attack patterns. We ingest EPSS scores via API to prioritize vulnerabilities not just by severity, but by real-world exploitation likelihood—enabling dynamic risk-based vulnerability management.

MISP (Malware Information Sharing Platform) (MISP Project, 2024) is an open-source platform designed for structured threat intelligence sharing using STIX/TAXII formats. It facilitates sharing of IOCs, threat event correlations, and TTP mappings. We integrate MISP data via its API to ingest indicators (e.g., hashes, domains, IPs), related threat actors, and event metadata. These enrich our knowledge graph with actionable CTI feeds.

VirusTotal (VirusTotal (Google Chronicle), 2024) is a widely used threat intelligence platform that aggregates malware analysis and sandbox reports from multiple antivirus engines and security vendors. To support behavior analysis and attribution tasks, CYBERTEAM collects structured threat metadata from VirusTotal's public API, including file hashes (MD5, SHA-1, SHA-256), behavioral execution traces, contacted IPs/domains, dropped files, and detection labels. This information is linked to threat artifacts such as malware families, indicators of compromise (IOCs), and known campaign signatures. The extracted metadata enables CYBERTEAM to contextualize adversarial

behaviors and enrich analytical functions like malware classification, infrastructure extraction, and campaign correlation.

AlienVault Open Threat Exchange (OTX) (AlienVault (AT&T Cybersecurity), 2024) is a collaborative threat-sharing platform that provides community-contributed threat indicators and contextual threat intelligence. Cyberteam leverages the OTX API to collect threat pulses—curated collections of IOCs and metadata describing specific threat actors, campaigns, or vulnerabilities. These pulses include information such as associated IPs, domains, file hashes, CVEs, and targeted sectors. By integrating OTX data, Cyberteam enhances its ability to support tasks like actor attribution, TTP matching, and community correlation, allowing LLMs to reason over shared intelligence and align analysis with ongoing threat landscapes.

Data Ethics. All data used in CYBERTEAM are collected from publicly available vulnerability databases and open-source threat intelligence platforms. No sensitive personal information or proprietary organizational data are included.

B MODULARIZED OPERATIONS: BASIC COMPONENT OF STANDARDIZED THREAT HUNTING

To support modular and extensible capabilities within our CYBERTEAM, we decompose complex NLP workflows into discrete, modularized operations. This section detail the implementation of NLP modules as described in section 3.1. Each module corresponds to a specific operation type, described as follows:

B.1 NER (NAMED ENTITY RECOGNITION)

To identify and classify cybersecurity-relevant entities such as threat actors, malware names, vulnerabilities, and indicators of compromise (IOCs) in unstructured textual data, NER facilitates automated extraction for threat attribution and situational awareness. We employ prompt-based techniques that enable entity recognition without retraining, thus maintaining adaptability to emerging domain vocabulary.

Prompt 1. NER Prompt for Threat Attribution

System Prompt: You are a cybersecurity threat intelligence assistant specialized in named entity recognition. Your task is to extract and categorize all named entities relevant to threat attribution from the provided text. Focus on answering: "Who is responsible for the attack?", "How was the attack carried out?".

Instructions: Given a cybersecurity-related document or report excerpt, extract all relevant named entities and classify them into:

- Threat Actor: Individual(s) or groups suspected or known to conduct the activity.
- Malware/Tool: Names of malicious software, exploits, or hacking tools.
- **Vulnerability:** CVE identifiers or technical flaws exploited.
- **Infrastructure:** IPs, domains, file hashes, or URLs used.

Output: Return results as a structured JSON object.

Design Rationale: In real-world threat hunting, analysts are constantly overwhelmed by unstructured reports, logs, and advisories filled with technical jargon and entity references. Automating named entity recognition helps blue teams immediately isolate critical items such as threat actors, malware strains, or CVE identifiers without combing through entire reports manually. This reduces analyst workload, accelerates attribution, and ensures no important entity slips through, particularly when adversaries recycle or slightly modify names and indicators across campaigns.

B.2 REX (REGEX PARSING)

To extract structured indicators from cybersecurity logs or reports, REX employs predefined regular expressions to match patterns like IP addresses, domain names, file hashes, and timestamps. This

rule-based approach offers high precision in normalizing threat data for correlation and enrichment tasks.

Prompt 2. Regex Pattern Matching Prompt

System Prompt: You are a cybersecurity parsing assistant. Your task is to extract standard threat indicators from raw incident reports using predefined regex patterns.

Instructions: Parse the following document and extract any matches for:

· IP addresses

- File hashes (MD5, SHA1, SHA256)
- · Domain names
- Timestamps

Output: Return all matches grouped by type in structured JSON format.

Design Rationale: Regex parsing remains indispensable because many threat indicators—such as IP addresses, hashes, and domains—follow strict syntactic patterns. Blue team analysts often must quickly normalize raw log data or incident feeds into structured formats suitable for correlation across SIEM or TIP platforms. Automated regex-based extraction delivers high precision and avoids false alarms, providing reliable building blocks for pivoting investigations, linking disparate alerts, and enriching threat databases with verified observables.

B.3 SUM (SUMMARIZATION)

To enable analysts to quickly grasp key information from lengthy threat reports, SUM generates concise summaries while preserving critical details such as TTPs, IOCs, and incident timelines.

Prompt 3. Threat Report Summarization Prompt

System Prompt: You are a cybersecurity analyst assistant. Your task is to summarize the following threat report in 3–4 sentences, preserving the attack vector, affected systems, timeline, and any mentioned threat actors or IOCs.

Instructions: Summarize only the essential intelligence. Avoid generic phrases. Include dates, names, and tools where available.

Output: Return a plain-text summary paragraph.

Design Rationale: Threat reports and advisories are typically lengthy, verbose, and include redundant or irrelevant details. In time-critical investigations, analysts need condensed yet accurate snapshots that retain attack vectors, key actors, and affected assets. Automated summarization provides blue teams with quick situational awareness, enabling them to brief stakeholders or prioritize triage without missing essential context. It also helps align tactical actions with strategic threat intelligence by stripping away noise and surfacing the essentials.

B.4 SIM (TEXT SIMILARITY MATCHING)

To determine semantic equivalence between pairs of threat indicators—particularly geographic or cultural references (e.g., "Eastern European" vs. "Russian-speaking")—the SIM function applies LLM-based textual similarity matching. This is critical for normalizing contextual descriptions found in incident reports or threat assessments that use varied, informal, or aliasing terms to describe similar threat origin profiles. Rather than relying on surface-level keyword overlap, SIM leverages the LLM's contextual understanding to judge whether two descriptions refer to the same underlying group or region. This helps unify disparate threat intelligence entries that may use different terminology for the same adversarial origin.

1027 1028

1029 1030 1031

1032 1033

1034 1035 1036

1039 1040

1041

1043 1044 1045

1046

1051 1052 1053

1055 1056 1057

1054

1058 1059

1061 1062 1063

1064

1066 1067

1068

1074 1075

1076

1077 1078

1079

Prompt 4. Text Similarity Matching Prompt for Geocultural Indicators

System Prompt: You are a cybersecurity assistant that helps analysts determine whether two geolocation or cultural indicators refer to the same threat origin. Use contextual reasoning to decide whether the two phrases describe the same group or region in a cyber threat context. Instructions: Given two input phrases describing threat origin (e.g., "Russian-affiliated" vs. "Eastern Bloc actor"), determine whether they semantically refer to the same group or geopolitical background.

Answer the following questions:

- Do both descriptions point to the same cultural, linguistic, or geopolitical region?
- Are the expressions used interchangeably in threat intelligence contexts?

Output: Return a JSON object with:

- "match": Boolean (true/false)
- "confidence": A float score from 0.0 to 1.0
- "justification": One or two sentences explaining the decision

Design Rationale: Threat hunting often suffers from inconsistent terminology—analysts and vendors may describe the same adversary group or region in different ways. By applying semantic similarity matching, blue teams can unify aliases, regional descriptions, or contextual cues, thereby avoiding fragmented analysis. For example, detecting that "Eastern European actors" and "Russian-speaking threat groups" likely refer to the same set of adversaries allows more coherent attribution and prevents intelligence silos that adversaries can exploit.

B.5 MAP (TEXT MAPPING)

To visualize and semantically relate named entities and key concepts extracted from cybersecurity documents, the MAP function supports construction of structured representations such as knowledge graphs or threat maps. These representations help uncover infrastructure relationships, campaign patterns, and geotemporal dynamics in threat activity. When powered by large language models, MAP enables flexible and context-aware extraction of relational triples from unstructured threat reports.

Prompt 7. Threat Knowledge Mapping Prompt

System Prompt: You are a cybersecurity knowledge graph assistant. Extract and relate key entities from the given threat report to form subject-predicate-object triples.

Instructions: Identify entities (e.g., threat actors, tools, organizations, IP addresses) and the relationships between them (e.g., "uses", "targets", "associated with").

Output: Return a list of triples in the format: [subject, predicate, object] Include a confidence score (0-1) if applicable.

Design Rationale: Attack campaigns rarely consist of isolated events—they are orchestrated through complex infrastructures and actor-tool relationships. Mapping extracted entities into structured knowledge graphs helps analysts visualize these relationships and trace adversary activity across time and geography. This capability supports detection of infrastructure reuse, identification of campaign evolution, and discovery of hidden connections that might otherwise remain unnoticed, enabling more proactive defense strategies and long-term threat tracking.

B.6 RAG (RETRIEVAL-AUGMENTED GENERATION)

To enhance generation with accurate and recent data, RAG combines LLM output with real-time retrieval from external threat intelligence APIs or databases. It is particularly useful for describing evolving threats or identifying actor affiliations.

Prompt 4. Structured Query for Retrieval

System Prompt: You are a cybersecurity assistant. Formulate a concise search query to retrieve current information about the topic specified below.

Instructions: Based on the topic "Recent activity by APT29 involving phishing attacks", generate a query such as:

"APT29 phishing campaign 2024 indicators, tools, and targets site:mitre.org OR site:virustotal.com"

Output: Return the final query string and optionally list key evidence passages from results.

Design Rationale: Adversary tactics evolve daily, and static LLMs quickly become outdated if disconnected from real-time sources. Retrieval-augmented generation enables blue teams to ground LLM outputs with fresh, authoritative information from trusted CTI feeds, vulnerability databases, or public repositories. This ensures that generated insights remain both accurate and timely, supporting decisions during live incidents such as phishing outbreaks or zero-day exploitation campaigns where stale intelligence could lead to ineffective responses.

B.7 SPA (TEXT SPAN LOCALIZATION)

To precisely extract actionable phrases—such as indicators of compromise or technique descriptions—from long-form cybersecurity text, **Text Span Localization** (SPA) models are used.

Two key metrics evaluate SPA effectiveness:

• Exact Match (EM):

$$EM = \frac{Number\ of\ exact\ matches}{Total\ predictions}$$

• Intersection over Union (IoU):

$$IoU = \frac{|S_p \cap S_t|}{|S_p \cup S_t|}$$

These metrics assess both strict and partial correctness, aiding in accurate downstream processing such as relation extraction or automated summarization.

Prompt 5. Span Extraction Prompt

System Prompt: You are a cybersecurity span identification assistant. Extract the text span that describes the primary technique used in the attack.

Instructions: Given a report excerpt, locate and return the sentence or phrase that directly describes how the attacker compromised the system (e.g., phishing, lateral movement, privilege escalation).

Output: Return the extracted span as plain text.

Design Rationale: In practice, analysts often need to pull out the single critical phrase—such as the exact exploitation method—from long reports or alerts. Span localization ensures precision by targeting actionable fragments rather than broad summaries, which is vital for creating detection rules, YARA signatures, or SIEM correlation logic. By pinpointing exact techniques or IOCs, blue teams reduce ambiguity, streamline evidence curation, and avoid wasting resources on imprecise or overly generalized intelligence.

B.8 CLS (CLASSIFICATION)

To measure the ability of a system to categorize cybersecurity-relevant textual inputs—such as threat alerts, vulnerability descriptions, or log messages—into predefined classes (e.g., threat categories, severity levels, or attack types), classification models are employed. This is commonly performed using transformer-based large language models (LLMs), which utilize a special token (e.g., [CLS]) to represent sentence-level semantics. The resulting embedding is mapped to labels through a learned classifier.

Design Rationale: Blue teams constantly receive heterogeneous data ranging from phishing alerts to vulnerability disclosures. Automated classification allows this information to be triaged into relevant categories—such as attack type, severity, or impacted systems—so that workflows can be routed efficiently. Accurate classification supports prioritization of critical alerts, ensures compliance with response playbooks, and minimizes analyst fatigue by filtering out low-severity noise while surfacing the incidents that require immediate attention.

B.9 MATH (MATHEMATICAL CALCULATION)

To perform quantitative analyses and structured computations relevant to cybersecurity, the **MATH** function supports tasks such as frequency modeling, impact scoring, cryptographic evaluation, and automated threat prioritization. These computations are critical for risk-informed decision-making within cyber threat intelligence pipelines.

A prominent example is the **Common Vulnerability Scoring System (CVSS v3.1)**, which uses a combination of weighted factors and conditional logic to produce a standardized severity score for vulnerabilities. One key element is the *Base Score*, calculated using the Impact and Exploitability sub scores:

$$\mbox{Base Score} = \begin{cases} 0, & \mbox{if Impact Subscore} \leq 0 \\ \mbox{RoundUp}\left(\mbox{min}(\mbox{Impact} + \mbox{Exploitability}, 10)\right), & \mbox{if Scope is Unchanged} \\ \mbox{RoundUp}\left(\mbox{min}(1.08 \times (\mbox{Impact} + \mbox{Exploitability}), 10)\right), & \mbox{if Scope is Changed} \end{cases}$$

The Impact Subscore is computed from confidentiality, integrity, and availability impact metrics as:

$$ISC_{Base} = 1 - (1 - C) \times (1 - I) \times (1 - A)$$

This formula models the probability that the system's security properties are affected by a vulnerability. The resulting score guides patching priority, risk exposure assessments, and automated vulnerability triage.

Such logic-heavy, non-trivial calculations exemplify the role of mathematical modules in operational cybersecurity settings and justify the integration of computational reasoning capabilities in modern cyber AI systems.

Prompt 9. CVSS Score Computation Prompt

System Prompt: You are a cybersecurity scoring assistant. Given a vulnerability description and metric values (Confidentiality, Integrity, Availability, Scope, Attack Vector, etc.), compute the CVSS v3.1 Base Score.

Instructions: Use the official CVSS equations and apply the rounding rules specified in the standard. Return both the numeric score and a textual explanation of the computation steps. **Output:** Return the Base Score as a float (1 decimal place) and a step-by-step explanation.

Design Rationale: Quantitative scoring frameworks like CVSS remain the backbone of enterprise vulnerability management and patch prioritization. Automated mathematical reasoning allows blue teams to consistently compute, validate, and apply these scores across large vulnerability sets, ensuring consistent triage even under heavy load. Beyond CVSS, mathematical modules enable probability modeling, risk scoring, and exposure forecasting—practices that help defenders allocate resources effectively and justify decisions to leadership with evidence-based metrics.

C METRIC

Below are further details on how each evaluation metric quantifies the corresponding threat hunting performance.

C.1 GENERATION (PRECISION–RECALL BALANCE BY F1) AND CLASSIFICATION (ACCURACY)

In threat hunting, information extraction tasks such as detecting malware names, extracting IOCs, or identifying exploited vulnerabilities require a careful balance between precision and recall. If a system retrieves too many irrelevant indicators, analysts are burdened with noise; if it misses critical signals, adversarial activity may go unnoticed. The **F1 score** captures this balance by evaluating how well a model retrieves the right items while minimizing both false alarms and missed detections. This makes it particularly valuable in operational contexts where the completeness and reliability of extracted intelligence directly affect the quality of subsequent analysis and response.

Besides, well-quantified tasks such as prioritization in blue team activities involve classification, such as determining whether an alert corresponds to privilege escalation, categorizing attack vectors, or assigning severity levels to vulnerability reports. In these scenarios, **accuracy** serves as an intuitive and effective measure of system performance, reflecting how often predictions align with ground-truth categories. High accuracy ensures that automated classification supports efficient triage and aligns with established taxonomies like MITRE ATT&CK.

C.2 SIM (BERT SCORE)

To evaluate the semantic similarity between cybersecurity-related texts—such as comparing analyst-written threat summaries, aligning generated incident narratives with original reports, or verifying paraphrased explanations of threat indicators—the **Sim** function utilizes contextual embedding-based metrics. Specifically, it computes **BERTScore** (Zhang et al., 2020), which has been shown to correlate strongly with human judgment in natural language generation tasks.

BERTScore measures semantic equivalence at the token level by aligning contextual embeddings from pre-trained transformer models. The score is computed as:

BERTScore =
$$\frac{1}{|x|} \sum_{i} \max_{j} \cos(\mathbf{x_i}, \mathbf{y_j})$$

where $\mathbf{x_i}$ and $\mathbf{y_j}$ are contextual embeddings of tokens in the candidate and reference texts, respectively. The final score reflects the average of maximal cosine similarities for each token in the candidate sentence.

This metric is particularly valuable in evaluating machine-generated text in cybersecurity domains, where surface-level similarity may fail to capture the deeper equivalence of technical meaning or threat context.

C.3 Pass (Code Execution Passing Rate)

To measure the reliability and functional correctness of cybersecurity automation artifacts—such as detection rules, analysis scripts, or integration workflows—the **Pass Rate** metric is employed. It quantifies how well a system performs under test by evaluating the proportion of test cases that execute successfully within a defined execution cycle, often conducted in a continuous integration (CI) pipeline.

Formally, the Pass Rate is defined as:

$$\text{Pass Rate} = \frac{\text{Number of Passed Tests}}{\text{Total Tests Executed}} \times 100\%$$

This metric provides a coarse yet effective indicator of operational readiness. A high Pass Rate implies that the deployed codebase functions as intended across its tested scenarios, which is critical in cybersecurity contexts where automation is used to process threat intelligence, detect anomalies, or trigger incident response mechanisms.

Routine monitoring of this metric supports the early identification of integration regressions, promotes pipeline stability, and ensures confidence in deploying automated defensive measures to production environments.

C.4 HIT (TOP-K HIT RATIO)

To evaluate the effectiveness of cybersecurity recommendation or retrieval systems—such as those that propose relevant threat indicators, patch suggestions, attack techniques, or investigative leads—the **Top-k Hit Ratio** is employed. This metric measures how frequently at least one correct or relevant item appears within the top-k ranked results returned by the system.

Mathematically, the Top-k Hit Ratio is defined as:

$$\label{eq:hit} \mbox{Hit@k} = \frac{\mbox{Number of queries with at least one relevant item in top k}}{\mbox{Total number of queries}}$$

A higher Hit@k indicates better system performance in surfacing relevant intelligence near the top of recommendations, which is critical for time-sensitive security operations.

Use Case Example: If a system recommends threat indicators based on a query about a ransomware family, Hit@5 evaluates whether at least one valid IOC (e.g., file hash or C2 domain) appears in the top 5 returned items.

Prompt 6. Hit Evaluation Prompt for Threat Retrieval

System Prompt: You are an assistant for evaluating cybersecurity retrieval systems. Given a query and a list of system-generated recommendations, check whether any ground truth item appears within the top-k returned results.

Instructions: For each query, compare the top-k predicted items against the gold-standard set. Indicate "Hit" if at least one match exists, otherwise "Miss".

Output: Return a JSON object with fields: query, top_k_results, ground_truth, hit@k: true/false

C.5 DIST (NORMALIZED DISTANCE SIMILARITY)

To evaluate the accuracy of numeric predictions in range-based estimation tasks, such as severity scoring, the **Normalized Distance Similarity (Dist)** metric is employed. This metric compares the predicted number and the ground-truth and scales the similarity into the [0,1] range, where higher values indicate closer alignment.

Formally, the similarity is computed as:

$$\text{Similarity} = 1 - \frac{|\hat{c} - c|}{R}$$

where \hat{c} and c denote the midpoints of the predicted and true ranges, respectively, and R is the maximum possible value of the range (e.g., 10 in our case of CVSS scores). The metric reflects the Euclidean distance between prediction and truth, normalized such that a perfect match yields a similarity of 1, and the furthest possible discrepancy yields 0.

D EXPERIMENTAL SETTING

This section details the experimental setup used to evaluate LLMs in the CyberTeam benchmark.

Hyperparameters. Table 4 summarizes the key hyperparameters for querying LLMs during experiments. These settings were chosen to balance generation quality and computational efficiency.

Computational Resources. All experiments were conducted on a high-performance computing cluster equipped with six NVIDIA RTX 6000 Ada Generation GPUs, each with 48 GB of dedicated VRAM. The system utilized CUDA version 12.8 and NVIDIA driver version 570.124.06. This configuration enabled parallel execution of model inference, evaluation, and tool-augmented tasks across the benchmark datasets. The hardware provided sufficient memory bandwidth and processing power to handle large-scale experiments, including multi-sample prompting strategies like CoT and

1296 1297

Table 4: LLM query hyperparameters.

Hyperparameter Value Description Temperature 0.7 Output randomness Top-p 0.95 Nucleus sampling threshold Max tokens 2048 Generation length cap ["\n", "Q:"] Stop sequences Response cutoff cues ICL, CoT, ToT, Emb Prompt format Prompt types (see 4) Tool-calling API Enabled (Selective) For function-use experiments

1304 1305 1306

1307

1303

Table 5: Running time (in seconds) of LLMs on CYBERTEAM, comparing different open-ended prompting strategies with our standardized method. Lower values indicate faster inference.

308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328

Metho	d	Cyber	security	Agent	Industry-Leading LLM								
1/10/110	-	LY	DH	SL	G4o	Go4	QW	GM	CD	L3.1	L4	GA	
				Pla	ybook R	ecomme	nd						
	ICL5	12.4	15.6	14.8	10.5	41.2	13.6	11.9	12.7	28.6	24.0	16.8	
Open-ended	ICL10	14.1	17.2	16.3	12.0	45.7	15.2	13.5	14.4	31.4	26.5	18.3	
Open-ended	CoT	18.6	22.4	21.1	15.8	60.8	19.9	17.6	18.8	41.2	34.5	24.5	
	ToT	27.5	34.1	32.0	24.2	89.5	30.2	27.1	28.9	62.7	52.5	37.8	
Standardized	(Ours)	21.3	26.5	25.2	19.1	71.3	23.5	21.0	22.3	50.5	42.0	30.1	
Security Control Adjust													
	ICL5	13.1	16.4	15.5	11.1	43.5	14.3	12.5	13.3	30.2	25.0	17.6	
Open-ended	ICL10	15.0	18.3	17.2	12.6	48.1	16.1	14.2	15.1	33.0	27.4	19.1	
Open-ended	CoT	19.4	23.5	22.2	16.7	63.4	20.8	18.3	19.6	43.8	36.0	25.7	
	ToT	28.3	35.6	33.6	25.4	92.2	31.7	28.3	30.2	66.4	55.0	39.5	
Standardized	l (Ours)	22.1	27.8	26.7	20.0	74.2	24.7	22.1	23.4	53.1	44.0	31.2	
				Pat	ch Code	Generati	on						
	ICL5	14.2	17.9	17.0	12.2	46.8	15.7	13.6	14.4	32.8	27.0	19.2	
Open-ended	ICL10	16.3	19.6	18.7	13.7	51.3	17.5	15.3	16.2	35.7	29.5	20.8	
Open-ended	CoT	21.2	25.3	24.5	18.1	68.7	22.9	19.7	21.1	47.6	39.0	28.1	
	ToT	31.7	38.4	36.9	27.8	98.6	34.5	30.6	32.6	71.9	59.0	42.5	
Standardized	(Ours)	24.0	29.7	28.6	21.4	79.4	26.6	23.6	25.0	57.2	47.0	34.4	
				Pat	ch Tool	Suggestio	on						
	ICL5	12.8	15.9	15.2	10.9	42.6	14.0	12.3	13.0	29.5	24.3	17.0	
Open-ended	ICL10	14.7	17.7	16.9	12.4	47.0	15.8	14.0	14.8	32.4	26.2	18.6	
Open-ended	CoT	19.0	23.0	22.0	16.3	62.1	20.5	18.1	19.2	42.5	35.0	25.1	
	ToT	27.9	34.9	33.1	24.7	90.8	31.0	27.6	29.6	64.0	52.5	38.2	
Standardized	l (Ours)	21.7	27.1	26.2	19.6	72.8	24.1	21.7	22.9	51.7	42.5	30.5	
				Ad	lvisory C	Correlatio	n						
	ICL5	13.6	16.8	16.1	11.7	44.9	14.9	13.0	13.8	31.0	25.5	18.2	
Oman ands 1	ICL10	15.6	18.7	17.9	13.2	49.6	16.7	14.7	15.6	34.0	28.0	20.0	
Open-ended	CoT	20.3	24.1	23.4	17.2	65.3	21.7	19.0	20.4	45.3	37.0	26.5	
	ToT	29.8	36.8	35.4	26.1	95.1	33.1	29.4	31.3	68.8	56.0	40.1	
Standardized	l (Ours)	23.1	28.5	27.8	20.6	76.2	25.4	22.8	24.2	54.6	45.0	32.1	

1334 1335 1336

ToT, without encountering resource constraints. Each experimental run was executed in a isolated environment to ensure reproducibility and avoid interference between tasks.

1337

E ADDITIONAL EXPERIMENTAL RESULTS

134213431344

This section presents additional experimental results that complement our main findings, offering deeper insights into model behavior across varied threat-hunting scenarios.

1345 1346

E.1 RUNNING TIME AND TRADE-OFF BETWEEN LATENCY AND EFFECTIVENESS

Observations and Insights. The runtime analysis highlights an inherent trade-off between efficiency and reasoning complexity across prompting strategies. Consistent with expectations, in-context learning (ICL) variants remain the fastest across nearly all models, typically completing tasks in the 10–15 second range. This makes ICL attractive for time-sensitive operations such as triage or initial

1370 1371 1372

1373

1374

1375

1376

1377

1378

1379

1380

1381

1382

1383

1384

1385

1386

1387

1388

1389

1390

1391 1392

1393 1394

1395

1396

1397 1398

1399

1400

1401

1402

1403

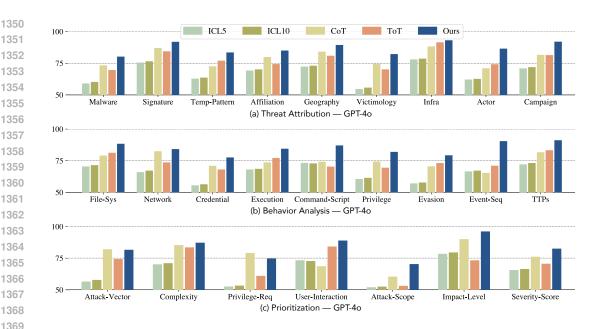


Figure 5: Threat-hunting performance on individual tasks, evaluating under GPT-4o.

correlation, where speed outweighs the need for more structured reasoning. Chain-of-thought (CoT) introduces additional reasoning overhead, increasing runtimes by roughly 30-40% compared to ICL. While this slowdown is measurable, the benefit of CoT lies in its improved consistency on more nuanced decision tasks, suggesting that blue teams might selectively invoke CoT when precision is critical. Tree-of-thought (ToT), by contrast, incurs the highest latency, often doubling the runtime relative to ICL. This stems from ToT's multi-branch exploration process, which, while occasionally producing richer reasoning chains, remains computationally expensive and operationally impractical for most real-time security workflows.

Our standardized pipeline approach falls between CoT and ToT in runtime. The added latency reflects the sequential decomposition of tasks into modular subroutines, each enforcing more structured reasoning than raw prompting. While slower than single-pass approaches, our pipeline mostly avoids the extreme overhead observed in ToT. This stability is particularly important in operational settings: analysts can predictably plan around a known latency budget while still benefiting from higher reliability and repeatability of results.

Unlike open-ended reasoning, which may fluctuate in quality depending on the model and prompt, the standardized pipeline enforces uniform logic steps, reducing error propagation at the cost of additional inference time. From a deployment standpoint, this balance offers a pragmatic middle ground: not as lightweight as ICL for quick heuristics, but substantially more usable than ToT when analysts demand repeatable outputs.

ADDITIONAL RESULTS OF INDIVIDUAL THREAT HUNTING PERFORMANCE

Figure 5, 6, 7, and 8 complement the results as present in Figure 3, offering aligned insights as exhibited in previous experiments.

Based on those results, we further outline the following observations and analyses:

Attribution-Oriented Tasks. Attribution tasks rely on aligning disparate indicators into coherent profiles of adversaries, infrastructure, and campaigns. Here, the standardized workflow shows its greatest benefit because it forces the model to treat each extracted clue as part of a larger dependency chain. When the reasoning is left open ended, models often generate fluent narratives that omit critical ties, such as overlooking how infrastructure relates to a specific campaign or how victimology patterns reinforce an actor hypothesis. The modular approach ensures that entity recognition, context mapping, and relational inference are explicitly sequenced, which reduces the tendency of the model

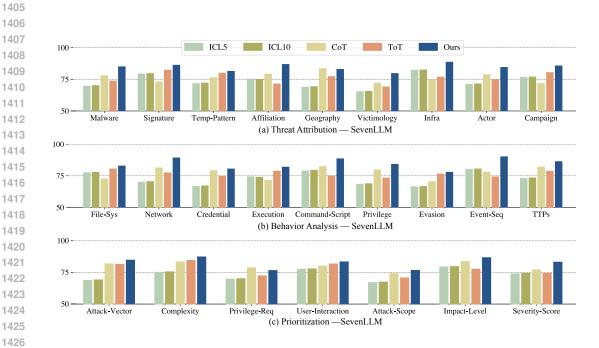


Figure 6: Threat-hunting performance on individual tasks, evaluating under SevenLLM-7B.

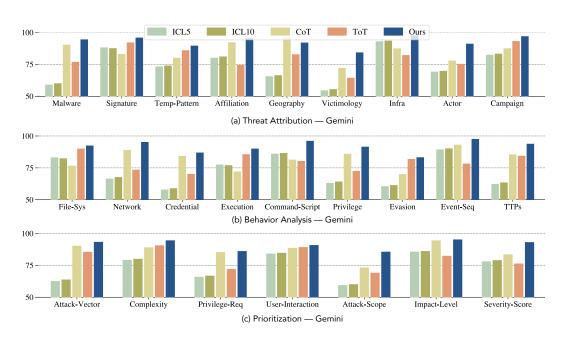


Figure 7: Threat-hunting performance on individual tasks, evaluating under Gemini-pro.

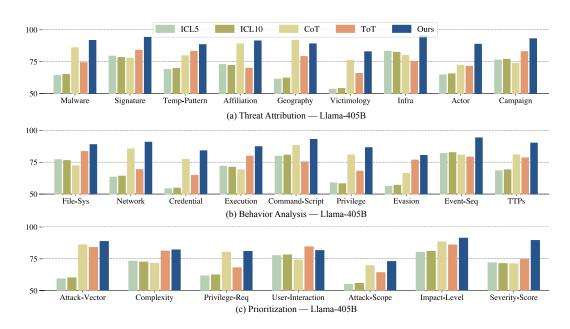


Figure 8: Threat-hunting performance on individual tasks, evaluating under Llama-405B.

to drift or collapse multiple actors into a generic label. This structured pipeline also helps the model preserve continuity across steps, so that information about geography, malware signatures, or campaign overlaps is not forgotten or misapplied. In practice, this makes attribution outputs more consistent and trustworthy, with fewer contradictions across different facets of the same incident.

Behavioral-Oriented Tasks. Behavioral analysis tasks focus on describing how an attack unfolds across file systems, networks, credentials, execution flows, and evasion strategies. These tasks expose a different challenge: models must reason not just about isolated labels but about temporal or causal sequences. Open ended reasoning often struggles to maintain logical order, for example by misplacing the relationship between a credential theft and subsequent privilege escalation, or by skipping intermediate steps in an event sequence. Standardized workflows address this by explicitly guiding the model to construct event chains step by step, preserving both order and dependency. This guidance is particularly important when behaviors are nested, such as when command script execution spawns further lateral movements or when evasion strategies are intertwined with persistence mechanisms. The modular design ensures that contextual cues are not discarded midway, producing outputs that resemble the structured analysis human analysts expect. The gains here are not simply about accuracy but about interpretability, as the resulting narratives make it easier to understand how behaviors connect to form a complete attack path.

Prioritization-Oriented Tasks. Prioritization tasks require models to map extracted observations into judgments of impact, severity, or scope. These are less about narrative flow and more about logical consistency and rule following. While open ended reasoning can handle straightforward labels like user interaction requirements or attack complexity, it often falters when multiple inputs must be integrated into a composite assessment. For example, determining severity requires careful alignment of impact level, attack vector, and privilege requirements, which is difficult to achieve reliably without structured steps. The standardized workflow enforces this alignment by ensuring that each component assessment is produced systematically and then fed into the final prioritization judgment. As a result, the model is less likely to generate inconsistent or contradictory scores. The benefits are particularly visible in tasks that resemble rule based calculations or scoring rubrics, where the modular structure mirrors the procedural way that human analysts reason about risk.

Broader Implications. Viewed across these categories, a clear pattern emerges. Attribution tasks benefit most from the preservation of contextual dependencies across different indicators. Behavioral tasks gain from the ability to model temporal and causal structure. Prioritization tasks see improvements in logical consistency and integration of multiple criteria. Standardization does not change the core language modeling capability of these systems, but it channels their generative power into

workflows that mirror how analysts actually think about security problems. This alignment between workflow design and task demands is the primary driver of the gains observed, and it demonstrates why modular guidance is most valuable when reasoning requires structured coordination across multiple dimensions.

Benchmark Generalizability. Although CyberTeam integrates data from 23 diverse threat intelligence sources, the benchmark is inherently constrained by the selected datasets and threat scenarios. For example, it may not fully represent emerging attack vectors, such as AI-powered phishing or supply chain compromises. Expanding the benchmark to include more recent and varied threat data, as well as cross-domain applications (e.g., IoT or cloud security), would enhance its utility for evaluating LLM generalization in broader cybersecurity contexts.

F LARGE LANGUAGE MODEL (LLM) USAGE DISCLOSURE

LLMs were employed exclusively for light grammar refinement and phrasing adjustments while preparing the manuscript. They were not involved in conceptual development, benchmark design, experiment execution, or result interpretation. All scientific ideas, methodological designs, and analyses were carried out independently by the authors. LLM usage was limited to minor textual polishing.