

# Towards Effective Offensive Security LLM Agents: Hyperparameter Tuning, LLM as a Judge, and a Lightweight CTF Benchmark

Minghao Shao<sup>1,2\*</sup>, Nanda Rani<sup>3\*</sup>, Kimberly Milner<sup>1\*</sup>, Haoran Xi<sup>1</sup>, Meet Udeshi<sup>1</sup>,  
Saksham Aggarwal<sup>1</sup>, Venkata Sai Charan Putrevu<sup>1</sup>, Sandeep K. Shukla<sup>4</sup>,  
Prashanth Krishnamurthy<sup>1</sup>, Farshad Khorrami<sup>1</sup>, Ramesh Karri<sup>1</sup>, Muhammad Shafique<sup>2</sup>

<sup>1</sup>New York University

<sup>2</sup>New York University Abu Dhabi

<sup>3</sup>Indian Institute of Technology Kanpur

<sup>4</sup>International Institute of Information Technology Hyderabad

## Abstract

Recent advances in LLM agentic systems have improved the automation of offensive security tasks, particularly for Capture the Flag (CTF) challenges. We systematically investigate the key factors that drive agent success and provide a detailed recipe for building effective LLM-based offensive security agents. First, we present CTFJudge, a framework leveraging LLM as a judge to analyze agent trajectories and provide granular evaluation across CTF solving steps. Second, we propose a novel metric, CTF Competency Index (CCI) for partial correctness, revealing how closely agent solutions align with human-crafted gold standards. Third, we examine how LLM hyperparameters, namely temperature, top-p, and maximum token length, influence agent performance and automated cybersecurity task planning. For rapid evaluation, we present CTFTiny, a curated benchmark of 50 representative CTF challenges across binary exploitation, web, reverse engineering, forensics, and cryptography. Our findings identify optimal multi-agent coordination settings and lay the groundwork for future LLM agent research in cybersecurity.

**Code** — <https://github.com/NYU-LLM-CTF/CTFJudge>

**Datasets** — <https://github.com/NYU-LLM-CTF/CTFTiny>

**Extended version** — <https://arxiv.org/abs/2508.05674>

## Introduction

Large language models (LLMs) have inspired the development of a new generation of autonomous AI agents capable of planning, reasoning, and interacting with tools to solve complex problems that were traditionally handled by human experts (Wang et al. 2024; Guo et al. 2024; Motlagh et al. 2024). The field of cybersecurity in particular has seen promising applications of these systems in solving Capture the Flag (CTF) challenges, which serve as realistic and adversarial exercises used to develop offensive security skills (Yang et al. 2023b; Shao et al. 2024a; Tann et al. 2023).

CTF challenges require stepwise reasoning, procedural decomposition, command execution, and deep domain

\*Authors contributed equally to this research.

Copyright © 2026, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

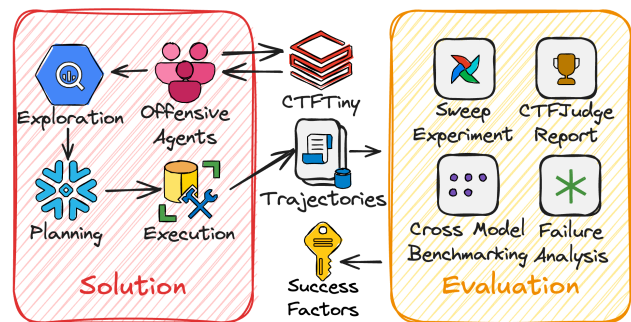


Figure 1: Overview of our evaluation, showing how hyperparameter tuning and multi-agent assessment benchmark LLM-based cybersecurity agents.

knowledge retrieval (Muzsai, Imolai, and Lukács 2024; Abramovich et al. 2025; Shao et al. 2025). They span various categories, including binary exploitation, reverse engineering, cryptography, web exploitation, and forensics, each demanding different toolchains and strategies. Thus CTFs are ideal for benchmarking offensive security agents to assess reasoning, tool use, and problem solving efficiency.

D-CIPHER is a multi-agent framework where specialized agents coordinate through planning, execution, and feedback to solve CTF problems (Udeshi et al. 2025). This human-like division of responsibilities improves scalability and modularity in solving diverse challenges. Yet, evaluating agentic systems is an open challenge and many issues need addressing. First, both the agentic system and its LLM operate as black boxes, sensitive to hyperparameter settings such as temperature, top-p, iteration limits, and token limits. Prior work did not examine how the hyperparameters affect behavior, performance of the agents and did not provide a recipe to study the factors underlying their success.

In addition to architectural and tuning limitations, evaluation methodologies are in their infancy. Prior work adopts a pass/fail metric based on whether the agent correctly extracts the flag or whether the agent completed the human-defined subtasks. This coarse metric doesn't capture nuances such

as partial progress, vulnerability detection ability, tool invocation efficiency, and reasoning steps attempted by the agent. This limits our understanding of agentic systems’ cybersecurity capabilities. Another bottleneck is the lack of standardized, lightweight benchmarks for reproducible experiments. Most CTF datasets are either heavy weight (Shao et al. 2024b; Zhang et al. 2025; Bhatt et al. 2024), or lack systematic difficulty analysis (Yang et al. 2023a). They are ill-suited for rapid, resource-constrained testing.

To address these limitations, we present a comprehensive study on offensive security agents. Our work includes an LLM judge agent for fine-grained assessment of the offensive reasoning capabilities of such systems, in-depth evaluation of key hyperparameters, and the use of a lightweight benchmark. Figure 1 illustrates the overall workflow of our evaluation pipeline, which integrates hyperparameter tuning, multi-agent coordination, and multi-dimensional evaluation on a standardized benchmark. By bridging methodological rigor with practical benchmarking, our framework not only supports the design of effective offensive security agents, but also offers a foundation for assessing agentic reasoning. Our contributions are threefold:

- Empirical analysis of offensive security agent system to show how hyperparameters (e.g., temperature, top- $p$ , and max tokens) affect performance, offering insights.
- CTFJudge, a fine-grained framework that assesses offensive security agents across aspects like vulnerability reasoning and exploitation techniques, uncovering bottlenecks beyond pass/fail metrics.
- CTFTiny is a curated set of 50 real-world CTF tasks spanning six domains, enabling reproducible, low-cost evaluation and parameter studies.

## Related Work

### LLM Agents for Offensive Security

Recent advancement of LLM agentic capabilities have led to the development of a variety of LLM agents for cybersecurity automation, relying on single-agent pipeline to modular designs with task decomposition and role-based agents (Guo et al. 2024; Song et al. 2024; Saha and Shukla 2025; Liu 2024; Bianou and Batogna 2024; Dorri, Kanhere, and Jurdak 2018; Abramovich et al. 2025). A recent work, D-CIPHER (Udeshi et al. 2025), presents a more comprehensive multi-agent system that enhances coordination, improves task execution, and boosts overall performance through task delegation, feedback loops, and specialized agent roles. Despite these advances, systematic evaluation of hyperparameter sensitivity remains an understudied area in the context of LLM agents. Prior works often evaluate with default parameters such as temperature and top- $p$ , or tune them in isolation without exploring their joint effects on reasoning quality, execution success, or inter-agent communication. Some studies have acknowledged the impact of decoding behavior (Turtayev et al. 2024; Yang et al. 2023a; Yao et al. 2022), but none have conducted targeted investigations on how these hyperparameters influence performance in multi-agent CTF-solving scenarios.

### CTF Benchmarks

CTF competitions have served as challenging tasks to evaluate the cybersecurity capabilities of automated LLM agents. The NYU CTF Bench (Shao et al. 2024b) and Cy-Bench (Zhang et al. 2025) provide multi-category challenge sets and step-by-step annotations, respectively, enabling structured evaluation. However, these datasets are often broad in scope or require significant computational resources, making them less suitable for rapid experimentation. Moreover, few are tailored for evaluating the nuanced behaviors of agentic systems under varying configurations. To address this, we introduce CTFTiny, a compact 50-challenge benchmark that enables rapid experiments, reproducible baselines, and exploration of agent design and parameter tuning strategies.

### LLM as a Judge

Traditional CTF evaluation often reduces agent performance to a pass/fail metric. To understand nuanced progress, recent work has explored using LLMs as automated judges (Gu et al. 2024; Li et al. 2024; Cao et al. 2025). These LLM judges assess reasoning quality, partial progress, and the correctness of intermediate steps, providing a deeper understanding of agent capabilities beyond simple flag capture. In this work, we use an LLM-judge that offers granular scoring of agent behavior across multiple aspects of task execution.

## Method

### CTFTiny Benchmark

CTFTiny is a lightweight evaluation benchmark curated from the full NYU CTF Bench to support rapid experimentation while maintaining representative challenge diversity. The motivation behind introducing CTFTiny is the high computational cost associated with running exhaustive evaluations on the full benchmark. Large-scale experiments, such as hyperparameter tuning, model ablations, and sensitivity analysis, are often constrained by prohibitive time and resource requirements of the full set. By reducing the benchmark size while preserving the core characteristics of challenge diversity and complexity, CTFTiny allows more frequent and cost-effective experimentation.

We have selected challenges for CTFTiny using a quantifiable measure of difficulty such that they are neither trivial nor prohibitively difficult. To estimate difficulty, we collate results on NYU CTF Bench from 12 configurations of previous works, namely D-CIPHER (Udeshi et al. 2025), and CRAKEN (Shao et al. 2025), and use the number of configurations that solved a challenge as empirical measure of difficulty. These configurations vary in LLM models, planning strategies, and agentic tools. This empirical measure offers a grounded basis for identifying challenges that require meaningful reasoning and decomposition. Challenge difficulty is classified based on how many of the 12 configurations solved them: 0–3 (hard), 4–6 (moderate), 6–9 (easy), and 9–12 (very easy), where a higher number means the challenge is easier to solve and a lower number means it is harder. Figure 2 shows the distribution of challenges from CTFTiny based on this empirical measure, demonstrating a healthy spread of challenges across all difficulty levels. The presence of

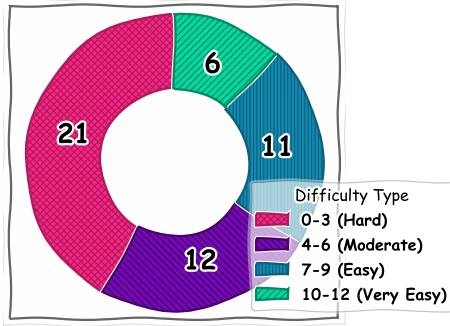


Figure 2: Challenge difficulty distribution based on configuration solves, showing a healthy mix of hard, moderate, and easy tasks for robust evaluation.

harder challenges ensures the benchmark to stressing model capabilities in reasoning, planning, and strategy adaptation.

### Hyperparameter Tuning

Our methodological evaluation includes a comprehensive experiment that varies LLM hyperparameters such as temperature, top- $p$ , and max tokens to assess their impact on the cybersecurity agent’s performance. These hyperparameters govern the trade-off between generation diversity and precision; exploring their values ensures our agent maintains robust reasoning, minimizes hallucinations, and optimizes resource usage across diverse challenges. This investigation helps identify robust defaults and guide adaptive parameter tuning for future deployments. By systematically examining these settings, we establish a principled basis for selecting decoding configurations that maximize solution accuracy and stability while controlling computational cost.

### CTFJudge Agent

CTFJudge provides a structured evaluation pipeline for LLM-driven CTF agents, assessing both flag retrieval performance and core cybersecurity competencies. The process begins by transforming expert-curated writeups into detailed, step-by-step summaries that capture each logical decision, underlying intent, key actions, and rationale. In addition the agent’s execution trace, which includes its planning choices, issued commands, observed outcomes, resource usage, and elapsed time, is abstracted into a summary format.

Once reference and candidate summaries are prepared, CTFJudge produces reference guided evaluations along six critical dimensions: vulnerability understanding, reconnaissance thoroughness, exploitation methodology, technical accuracy, efficiency of approach, and adaptability. This granular evaluation produces a quantitative CTF Competency Index (CCI) alongside a narrative report that highlights areas of strength (e.g. comprehensive scanning techniques) and pin-points opportunities for improvement such as gaps in protocol interpretation or suboptimal command sequencing.

Supporting large-scale, reproducible benchmarking, CTFJudge uses a centralized configuration layer to specify evaluation criteria, model parameters, and factor weights,

while built-in error handling and version tracking ensure reliable batch runs. By abstracting implementation details behind clear interfaces, the framework remains extensible, allowing researchers to swap in alternative summarization strategies, integrate real-time telemetry, or add anomaly detection without changing the core evaluation workflow.

### CTF Competency Index (CCI)

To quantify how closely an agent’s trajectory summary  $T$  aligns with a human-curated gold solution  $G$ , CTFJudge employs a weighted combination of  $n$  complementary evaluation factors under its default configuration:

$$CCI(T, G) = \sum_{i=1}^n w_i F_i(T, G), \quad \sum_{i=1}^n w_i = 1,$$

where  $n$  denotes the total number of factors and  $F_i$  is the factor that pre-defined. In our initial setup ( $n = 6$ ), these factors serve as the default evaluation criteria: vulnerability understanding, reconnaissance thoroughness, exploitation methodology, technical accuracy, efficiency of approach, and adaptability. The resulting score lies in  $[0, 1]$ , balancing strategic insight, precision, and operational efficiency for clear comparison and targeted feedback.

These six dimensions align with key stages of offensive cybersecurity. Vulnerability understanding gauges the agent’s ability to identify and interpret system flaws. Reconnaissance thoroughness measures the depth of information gathering. Exploitation methodology evaluates the robustness of attack planning, while technical accuracy ensures commands are executed correctly. Efficiency reflects resource and time optimization, and adaptability captures how the agent handles unexpected scenarios. Together, they provide a balanced benchmark across strategic, technical, and operational fronts, guiding targeted improvements in agent performance.

## Experiment Setup

### Metrics

Our evaluation framework adopts a multi-dimensional approach to assess agent performance across three key aspects.

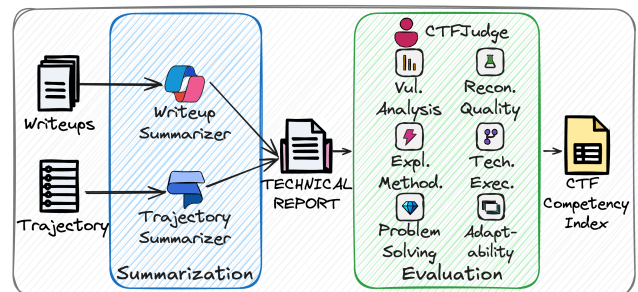


Figure 3: Architecture of the CTFJudge agent, showing how agent trajectories align with expert solutions to generate competency scores and actionable feedback.

**Pass@k** We use the standard pass@1 metric to measure the proportion of challenges successfully solved on the first attempt. We also report the average computational cost across configurations, highlighting the trade-offs between effectiveness and efficiency under different parameter settings.

**CCI** We fix the weights of six evaluation factors to be equal as competency in each activity is required for flag discovery. CTFJudge is provided with one expert/author write-up for each CTFTiny challenge, sequentially describing how to solve the challenge. The write-ups have insights on the targeted vulnerability, code snippets, tool references, and commentary on strategy and pitfalls. While CTFJudge is tasked to qualitatively compare AI solver trajectories in a solution, the graded performance matrix enables the framework to review the solver’s strategy, performance, and adaptability in navigating the challenge.

### Model Selection

We evaluate the offensive security agent on six state-of-the-art language models spanning diverse architectural families and capabilities. This selection ensures broad coverage across reasoning paradigms, parameter scales, and training methodologies. We employ Claude 3.7 Sonnet with a temperature of 0.1 for CTFJudge in grading the all agent-LLM interactions. We use official APIs from Anthropic, OpenAI, and Google for proprietary models to ensure performance and stability. We use the Together AI platform for open-source models.

**Proprietary models.** claude-sonnet-4-20250514, gpt-4.1-2025-04-14, gemini-2.5-pro, and gemini-2.5-flash. These commercially-deployed models benefit from large-scale training, proprietary optimization, and robust generalization over diverse, multi-step reasoning tasks in real-world scenarios.

**Open-source models.** llama-4-maverick-17b, qwen3-235b, and DeepSeek-V3-0324. These community-driven models prioritize transparency, reproducibility, and adaptability, enabling flexible integration and task-specific fine-tuning in constrained settings.

### Hyperparameter Selection

We evaluate the impact of LLM hyperparameters on D-CIPHER’s performance by varying:

- **Temperature:** {0, 0.2, 0.4, 0.6, 0.8, 1.0}
- **Top-p:** {0.25, 0.5, 0.75, 0.8, 0.85, 0.9, 0.95, 1.0}
- **Max tokens:** {2048, 4096, 8192}

For the baseline, we adopt the default configuration from the original D-Cipher implementation: temperature as 1.0, top-p as 1.0 max tokens as 4096.

## Results

### Baseline Results on CTFTiny

To establish a baseline for large language models (LLMs) on the CTFTiny benchmark, we evaluated seven state-of-the-art models across multiple cybersecurity domains. Among them, Claude 4 Sonnet achieved the highest performance,

solving 38 out of 50 challenges for a 76% success rate. Gemini 2.5 Flash followed with 32 correct solutions (64%), while Gemini 2.5 Pro and GPT-4.1 completed 24 (48%) and 20 (40%) challenges, respectively. The remaining models demonstrated more limited capabilities: Qwen 3 solved 14 challenges (28%), DeepSeek V3 solved 11 (22%), and LLaMA 4 Maverick 17B completed only 4 (8%).

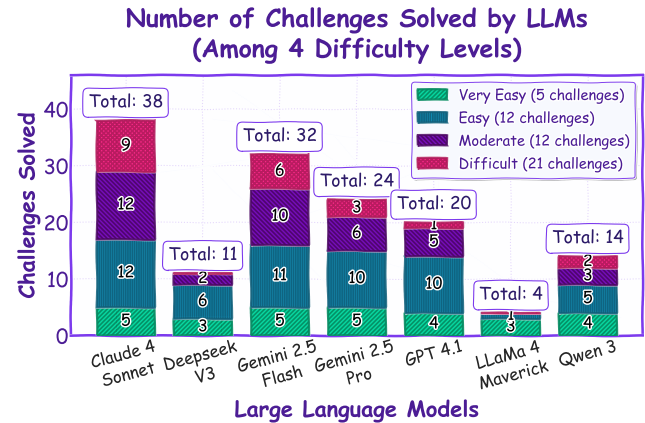


Figure 4: Challenge solves by model and difficulty, showing sharp performance drops as complexity rises.

Figure 4 breaks down each model’s solved challenges into four difficulty bands (*very easy*, *easy*, *medium*, *difficult*), as described in Section . Top models maintain near-perfect scores on the simplest challenges but decline sharply as complexity rises. Claude 4 Sonnet and Gemini 2.5 Flash solve 100% *very easy* instances and achieve over 40% and 25% respectively on *difficult* challenges, underscoring their robustness across task difficulty. Mid-ranked models (GPT-4.1, Qwen 3) have reasonable accuracy on *easy* tasks but fall below 50% in the *medium* band and solve less than 10% *difficult* challenges. Lower-ranked models (DeepSeek V3, LLaMA 4 Maverick 17B) succeed only on *very easy* and a fraction of *easy* problems. This highlights performance gaps on difficult CTFs underscoring the need for robust reasoning and decomposition strategies in agent designs.

A category-level analysis reveals pronounced differences in model strengths and specialization. In *reverse engineering*, Claude 4 Sonnet leads with an 81.3% success rate (13/16), followed by Gemini 2.5 Flash at 68.8%. In *cryptography*, Claude 4 Sonnet once again tops the leaderboard with 75% (9/12), while both Gemini variants reach 50%. Notably, Gemini 2.5 Flash outperforms Claude 4 Sonnet in *binary exploitation*, achieving 72.7% versus 63.6%, indicating a clear domain-specific specialization. *Forensics* tasks yield perfect scores from both Gemini models, though the sample size remains small. In *web exploitation*, Claude 4 Sonnet and Qwen 3 both achieve 100%. Finally, in the *miscellaneous* category, Claude 4 Sonnet maintains its overall lead at 83.3%. By mapping performance across difficulty tiers and categories, CTFTiny transforms raw scores into diagnostic insights, revealing each model’s unique strengths and blind spots, and charting an informed path toward resilient AI agents for cybersecurity.

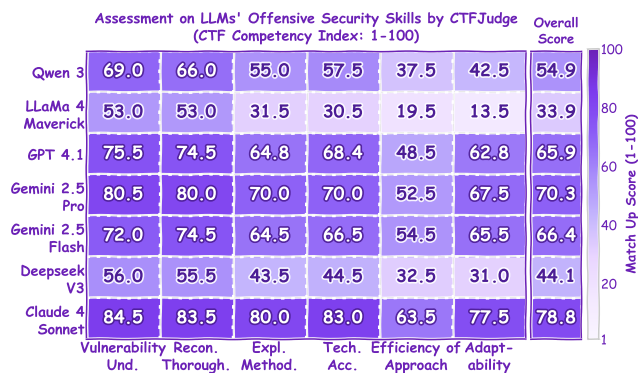


Figure 5: CTF Competency Index by model.

### CTF Competency Index (CCI)

Figure 5 presents the CTF Competency Index distribution for each model on the CTFTiny baseline. Claude 4 Sonnet consistently leads across all skill dimensions, with CCI from 77.5 to 84.5, reflecting robust reasoning and adaptability even on harder tasks. Gemini 2.5 Pro also delivers strong, balanced scores, especially in exploitation and adaptability, while Gemini 2.5 Flash, despite more solves, yields a lower overall CCI due to less systematic approaches. GPT-4.1 shows moderate results, particularly lagging in efficiency and adaptability. Qwen 3 performs reasonably on core skills but is limited by weak efficiency and adaptability. DeepSeek V3 and LLaMA 4 Maverick 17B stay at the lower end, underscoring persistent struggles, even on easy tasks.

Figures 6 and 7 show CCI cleanly separating outcomes: successes cluster at high, uniform scores; failures sit mid-low with sharp dips. Claude 4 and Gemini Pro/Flash lead; GPT-4.1/Qwen 3 are mid; DeepSeek/LLaMA weak. Reconnaissance and Vulnerability Understanding change little across outcomes. The separation comes from Exploitation Methodology, Efficiency of Approach, and Adaptability, which drop sharply on failed runs. Successful trajectories also show lower variance across skills, indicating steadier end-to-end execution. Thus, CCI not only distinguishes success from failure

Metric	Claude4S	DeepSeekV3	Gemini2.5F	Gemini2.5P	GPT 4.1	LLaMa 4M	Qwen 3
Total (%)	76	22	64	48	40	8	28
Cost (\$)	1.16	0.02	0.26	0.33	0.77	0.14	0.04
Cry (%)	75.0	16.7	50.0	50.0	33.3	0.0	25.0
For (%)	50.0	0.0	100.0	100.0	0.0	0.0	0.0
Pwn (%)	63.6	18.2	72.7	45.5	36.4	0.0	18.2
Rev (%)	81.3	18.8	68.8	37.5	37.5	12.5	18.8
Web (%)	100.0	66.7	33.3	33.3	66.7	33.3	100.0
Misc (%)	83.3	33.3	66.7	66.7	66.7	16.7	50.0

Table 1: Model performance on CTFTiny, with Claude 4 Sonnet leading but domain strengths varying.

but also pinpoints where the pipeline breaks.

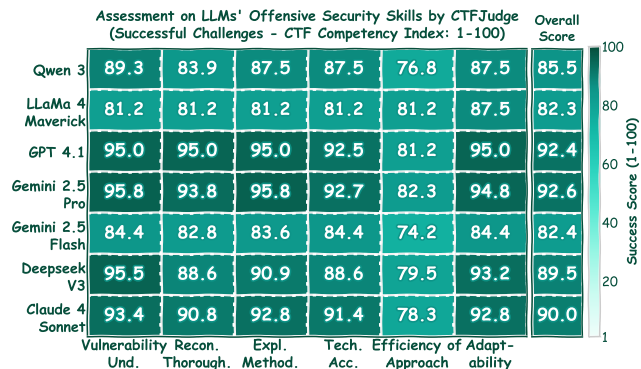


Figure 6: CTF Competency Index on a successful solution.

Key observations reinforce these patterns: Gemini 2.5 Pro, while solving fewer challenges than Flash, achieves higher average CCI, especially in exploitation and adaptability, suggesting more structured, human-aligned reasoning. Most models underperform on "Efficiency of Approach", relying on brute-force or redundant exploration—a clear bottleneck under cost constraints. CCI exposes the gap between solve rates and true reasoning quality, showing brute-force success rarely translates to robust, interpretable agent behavior. These findings highlight the need to evaluate agents by structured reasoning and operational efficiency, not solve count.

Figure 8 compares CCI by category. Claude 4 Sonnet is the most balanced; Gemini 2.5 Pro is similar with peaks in **for** and **pwn**. Flash excels in **for** and **web** but dips on **rev** and **misc**. GPT-4.1 spikes on **web**; Qwen 3 is mid-tier; DeepSeek V3 and LLaMA-4 Maverick remain low. Notably, **misc** and **pwn** show lower CCI, likely due to heavier reliance on challenge-server interaction, which adds overhead and fragility and lowers Efficiency/Adaptability relative to categories that rely more on local environments such as **rev**. **for** also lags, plausibly because challenges require handling diverse file types, which remains difficult for LLMs. **web** and **cry** score higher as their challenge formats are more standardized and pattern-repetitive, easing reasoning and verification under fixed toolchains and protocols.

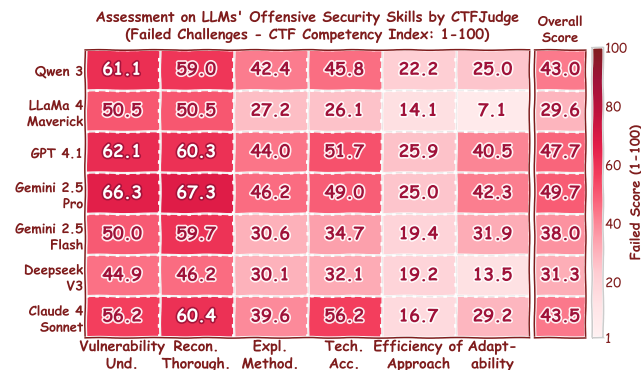


Figure 7: CTF Competency Index on failed solution.

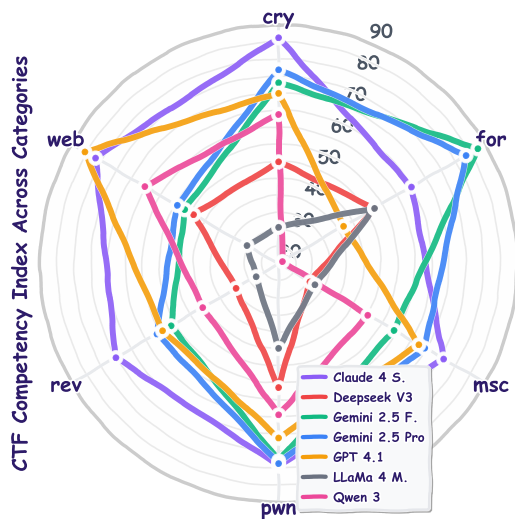


Figure 8: Category-wise CTF Competency Index across models, highlighting misc, pwn, and forensics challenges.

### Impact of Hyperparameters

We systematically analyze how temperature, top-p, and maximum token length affect LLM-based CTF agent performance, revealing their roles in shaping reasoning stability, execution precision, and success across diverse challenge types.

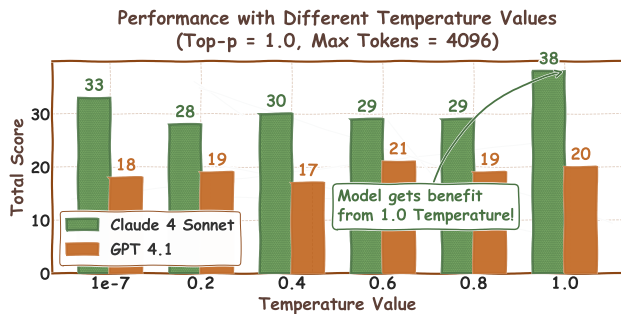


Figure 9: Model accuracy peaks at high temperature, highlighting the need for creative reasoning.

**Temperature.** As shown in Figure 9, model performance exhibits non-linear relationship with temperature. Claude 4 Sonnet reaches peak accuracy at the highest value ( $T = 1.0$ ) with 38 correct solves, contrasting with lower performance at intermediate settings. This suggests that controlled randomness enhances exploratory reasoning, especially in open-ended challenges. Claude’s performance remains stable across  $T \in [0.2, 0.8]$ , then surges at  $T = 1.0$ , pointing to creative reasoning benefits under high-temperature sampling. In contrast, GPT-4.1 peaks at  $T = 0.6$  (21 solves), with minimal variance across the sweep, indicating conservative decoding that resists both gains and drops. Overall, high temperature appears advantageous for models that effectively leverage generative diversity, particularly in ambiguous or multi-path problems common in CTF tasks.

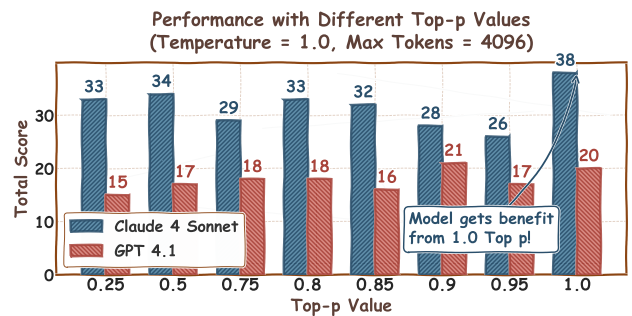


Figure 10: Model performance is stable across most top-p, with only the highest top-p notably increasing solve rates.

**Top-p.** In contrast to temperature, top-p adjusts diversity by truncating the token distribution. As shown in Figure 10, Claude 4 Sonnet performs consistently across  $p \in [0.25, 0.85]$ , hovering around 32–34 solves, and peaks at  $p = 1.0$  with 38 solves—suggesting full distributional access can boost exploration in strong models. This challenges the notion that lower top-p harms performance via distractor tokens. While minor fluctuations appear, Claude’s robustness implies control over long-tail generations. In contrast, GPT-4.1 performs stably between  $p = 0.75$  and 0.9, peaking modestly at  $p = 0.9$  with 21 solves, and shows no major drop at  $p = 1.0$ . These findings suggest that top-p near 1.0 may benefit models able to leverage expanded output space, especially in creative or underspecified CTF tasks.

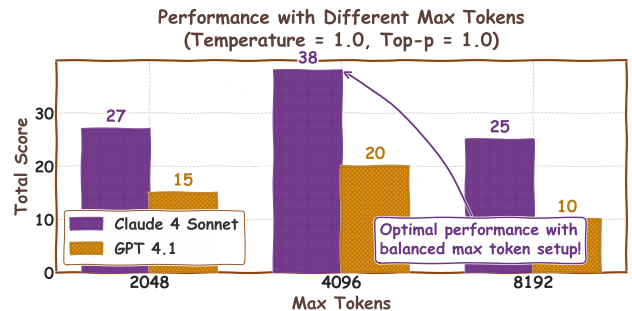


Figure 11: Mid-range max tokens yield best accuracy, while shorter or longer contexts hurt performance.

**Max Tokens.** Figure 11 presents a counterintuitive finding: longer context windows do not always improve performance. Claude 4 Sonnet peaks at 4096 tokens, solving 38 challenges, while performance drops at both 2048 and 8192. The dip at 8192 suggests context saturation, where longer completions may overwhelm model’s attention or generate overly verbose. This is especially problematic when prompt grounding is critical for correct planning, such as in **pwn** or **rev**. GPT-4.1 mirrors this trend with a modest peak at 4096 (20 solves), dropping to 10 at 8192. These results highlight a “Goldilocks zone” for context size: large enough to support complex reasoning, but not so large as to induce distraction or drift.

Taken together, these results show that model behavior is sensitive to decoding configurations, and optimal settings

require balancing determinism with exploration. Our study finds that higher temperature ( $\approx 1.0$ ) and top-p ( $\approx 1.0$ ) improve solve rates, particularly for strong models like Claude, by enabling flexible, multi-step reasoning. Low temperature, moderate top-p, and mid-range token lengths (4096) still offer a favorable trade-off between accuracy, stability, and efficiency, especially in precision-heavy domains. Though not universally optimal, these configurations form a strong baseline for LLM-based cybersecurity agents, with tuning informed by task complexity and feedback.

### Failure Analysis

We conducted a detailed analysis of model failures using 21 predefined categories highly relevant to cyber and agentic deficiencies. For each unsolved challenge the CTFJudge framework identified one or more reasons for failure and extracted key terms reflecting the root cause. Depending on the agent’s reasoning trace each unsolved challenge could be associated with multiple reasons for failure. The overall distribution of these failure types is presented in the line graph in Fig 12. The most frequent reason observed was a ‘Knowledge or Domain Expertise Gap’, followed by ‘Exploit Development Failure’ and ‘Insufficient Reconnaissance.’ These findings suggest models struggle due to limited understanding of CTF skills, inability to translate vulnerabilities into exploits, or failure to gather necessary scenario information. Notably despite D-CIPHER being a coordinated agentic framework, "Incorrect Task Delegation" and "Infrastructure or Environment Failure:" were in the lowest quartile for failures.

We performed a model-wise failure breakdown, shown in the heatmap<sup>1</sup> in Fig 12. LLaMa 4 Maverick showed the highest failure rates, while Claude 4 Sonnet had the lowest. The columns reveal a gradient: weaker models (LLaMa 4 Maverick, Qwen 3, DeepSeek V3) accumulate more failures; stronger ones (Claude 4 Sonnet, Gemini 2.5 Pro/Flash) show fewer, more concentrated errors. Gemini 2.5 Pro and Flash fall in between, matching their mid-tier solve counts. Except for Claude, all models show a steep drop from high to low frequency failure reasons, while Claude maintains uniformly low counts, indicating robustness. Overall, failure volume inversely tracks pass rate—better-performing models make fewer, less diffuse mistakes.

### Limitations and Future Work

Our work significantly advances automated evaluation of LLM-based CTF agents, revealing insights into reasoning-retrieval interactions and improving partial scoring, solution quality, and cost-efficiency over baselines. Nevertheless, as this field is nascent, several opportunities exist for further refinement. The CTF Competency Index employs a fixed, expert-curated weighting scheme that, while effective for studied challenges, could benefit from tuning or retraining to accommodate emerging or atypical CTF tasks. Our benchmarking focuses on a single open-source agent architecture (D-CIPHER) due to unavailability of alternative offensive-security frameworks; extending evaluations to fundamentally different designs would provide valuable insights. Moreover,

<sup>1</sup>Y-axis of heatmap shows failure reasons in the line graph.

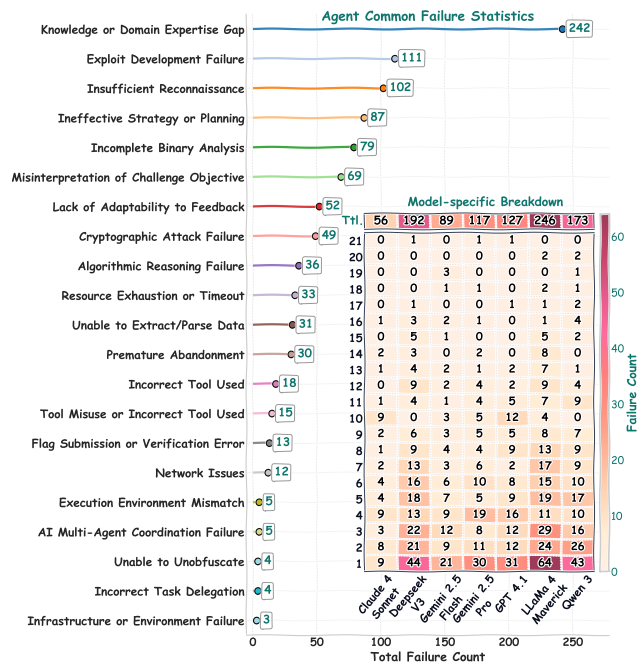


Figure 12: Failure analysis shows domain knowledge and exploit development as main agent bottlenecks.

dynamic metrics that adapt to domain drift would enhance the framework’s applicability to complex CTF scenarios.

Several directions can strengthen our framework. First, adaptive calibration for the CCI—alongside automated on-line evaluation metrics that detect domain drift and adversarial shifts—would improve resilience and support dynamic reweighting for new challenge types. Second, expanding baseline agents to include reinforcement-learning agentic systems, multi-role division multi-agent setups, and knowledge-based retrieval will provide broader performance insights. Third, more rigorous targeted ablation studies isolating explicit chain-of-thought and implicit reasoning will clarify each approach’s contribution to success and robustness.

### Conclusion

Our key contributions include: (i) a hyperparameter-centric evaluation showing how decoding settings like temperature, top-p, and max token affect reasoning stability, planning accuracy, and success; (ii) CTFJudge, a modular evaluation agent that yields an interpretable CTF Competency Index by aligning agent trajectories with expert solutions across six dimensions; and (iii) CTFTiny, a compact, diverse benchmark for reproducible, resource-efficient experimentation under constraints. A key insight is that optimal agentic behavior arises when decoding settings balance determinism with diversity to support adaptive, multi-step reasoning across CTF tasks. Beyond CTF, our tools and methods offer a framework for tasks like software repair, tool use, and cyber defense, to benchmark and improve next-generation LLM agents.

## Ethics Statement

All datasets utilized in this study were sourced from publicly available repositories and are properly cited throughout the paper. All experiments were conducted within controlled Linux server environments using exclusively open-source tools. To promote full reproducibility, we have open-sourced both the datasets and frameworks employed in this work, including CTFTiny and CTFJudge. Given the offensive security capabilities demonstrated by the agent systems studied, we underscore that any misuse of our code or methodologies for unauthorized attacks on real-world systems is strictly prohibited and may be illegal. This research is intended solely for educational use, academic research, and ethical security testing within properly authorized and controlled environments.

## Acknowledgements

This work has been supported in parts by the NYUAD Center for Cyber Security (CCS), funded by Tamkeen under the NYUAD Research Institute Award G1104, NYU Abu Dhabi Center for AI and Robotics CG010 and by the NYU-KAIST Partnership.

## References

- Abramovich, T.; Udeshi, M.; Shao, M.; Lieret, K.; Xi, H.; Milner, K.; Jancheska, S.; Yang, J.; Jimenez, C. E.; Khorrami, F.; Krishnamurthy, P.; Dolan-Gavitt, B.; Shafique, M.; Narasimhan, K.; Karri, R.; and Press, O. 2025. Interactive Tools Substantially Assist LLM Agents in Finding Security Vulnerabilities. *arXiv:2409.16165v2*.
- Bhatt, M.; Chennabasappa, S.; Li, Y.; Nikolaidis, C.; Song, D.; Wan, S.; Ahmad, F.; Aschermann, C.; Chen, Y.; Kapil, D.; Molnar, D.; Whitman, S.; and Saxe, J. 2024. CyberSecEval 2: A Wide-Ranging Cybersecurity Evaluation Suite for Large Language Models. *arXiv:2404.13161v1*.
- Bianou, S. G.; and Batogna, R. G. 2024. PENTEST-AI, an LLM-Powered Multi-Agents Framework for Penetration Testing Automation Leveraging Mitre Attack. In *2024 IEEE International Conference on Cyber Security and Resilience (CSR)*, 763–770.
- Cao, H.; Driouich, I.; Singh, R.; and Thomas, E. 2025. Multi-Agent LLM Judge: automatic personalized LLM judge design for evaluating natural language generation applications. *arXiv preprint arXiv:2504.02867*.
- Dorri, A.; Kanhere, S. S.; and Jurdak, R. 2018. Multi-agent systems: A survey. *IEEE Access*, 6: 28573–28593.
- Gu, J.; Jiang, X.; Shi, Z.; Tan, H.; Zhai, X.; Xu, C.; Li, W.; Shen, Y.; Ma, S.; Liu, H.; et al. 2024. A survey on llm-as-a-judge. *arXiv preprint arXiv:2411.15594*.
- Guo, T.; Chen, X.; Wang, Y.; Chang, R.; Pei, S.; Chawla, N. V.; Wiest, O.; and Zhang, X. 2024. Large Language Model based Multi-Agents: A Survey of Progress and Challenges. *arXiv:2402.01680*.
- Li, H.; Dong, Q.; Chen, J.; Su, H.; Zhou, Y.; Ai, Q.; Ye, Z.; and Liu, Y. 2024. Llm-as-judges: a comprehensive survey on llm-based evaluation methods. *arXiv preprint arXiv:2412.05579*.
- Liu, Z. 2024. Multi-Agent Collaboration in Incident Response with Large Language Models. *arXiv:2412.00652v2*.
- Motlagh, F. N.; Hajizadeh, M.; Majd, M.; Najafi, P.; Cheng, F.; and Meinel, C. 2024. Large Language Models in Cybersecurity: State-of-the-Art. *arXiv:2402.00891*.
- Muzsai, L.; Imolai, D.; and Lukács, A. 2024. HackSynth: LLM Agent and Evaluation Framework for Autonomous Penetration Testing. *arXiv:2412.01778v1*.
- Saha, B.; and Shukla, S. K. 2025. MalGEN: A Generative Agent Framework for Modeling Malicious Software in Cybersecurity. *arXiv preprint arXiv:2506.07586*.
- Shao, M.; Chen, B.; Jancheska, S.; Dolan-Gavitt, B.; Garg, S.; Karri, R.; and Shafique, M. 2024a. An Empirical Evaluation of LLMs for Solving Offensive Security Challenges. *arXiv:2402.11814v1*.
- Shao, M.; Jancheska, S.; Udeshi, M.; Dolan-Gavitt, B.; Xi, H.; Milner, K.; Chen, B.; Yin, M.; Garg, S.; Krishnamurthy, P.; Khorrami, F.; Karri, R.; and Shafique, M. 2024b. NYU CTF Bench: A Scalable Open-Source Benchmark Dataset for Evaluating LLMs in Offensive Security. In *Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.
- Shao, M.; Xi, H.; Rani, N.; Udeshi, M.; Putrevu, V. S. C.; Milner, K.; Dolan-Gavitt, B.; Shukla, S. K.; Krishnamurthy, P.; Khorrami, F.; et al. 2025. CRAKEN: Cybersecurity LLM Agent with Knowledge-Based Execution. *arXiv preprint arXiv:2505.17107*.
- Song, C.; Ma, L.; Zheng, J.; Liao, J.; Kuang, H.; and Yang, L. 2024. Audit-LLM: Multi-Agent Collaboration for Log-based Insider Threat Detection. *arXiv:2408.08902v1*.
- Tann, W.; Liu, Y.; Sim, J. H.; Seah, C. M.; and Chang, E.-C. 2023. Using Large Language Models for Cybersecurity Capture-The-Flag Challenges and Certification Questions. *arXiv:2308.10443*.
- Turtayev, R.; Petrov, A.; Volkov, D.; and Volk, D. 2024. Hacking CTFs with Plain Agents. *arXiv:2412.02776v1*.
- Udeshi, M.; Shao, M.; Xi, H.; Rani, N.; Milner, K.; Putrevu, V. S. C.; Dolan-Gavitt, B.; Shukla, S. K.; Krishnamurthy, P.; Khorrami, F.; et al. 2025. D-CIPHER: Dynamic Collaborative Intelligent Multi-Agent System with Planner and Heterogeneous Executors for Offensive Security. *arXiv preprint arXiv:2502.10931*.
- Wang, L.; Ma, C.; Feng, X.; Zhang, Z.; Yang, H.; Zhang, J.; Chen, Z.; Tang, J.; Chen, X.; Lin, Y.; Zhao, W. X.; Wei, Z.; and Wen, J. 2024. A survey on large language model based autonomous agents. *Frontiers of Computer Science*, 18(6): 186345.
- Yang, J.; Prabhakar, A.; Narasimhan, K. R.; and Yao, S. 2023a. InterCode: Standardizing and Benchmarking Interactive Coding with Execution Feedback. In *Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.
- Yang, J.; Prabhakar, A.; Yao, S.; Pei, K.; and Narasimhan, K. R. 2023b. Language Agents as Hackers: Evaluating Cybersecurity Skills with Capture the Flag.

Yao, S.; Zhao, J.; Yu, D.; Shafran, I.; Narasimhan, K. R.; and Cao, Y. 2022. ReAct: Synergizing Reasoning and Acting in Language Models.

Zhang, A. K.; Perry, N.; Dulepet, R.; Ji, J.; Menders, C.; Lin, J. W.; Jones, E.; Hussein, G.; Liu, S.; Jasper, D. J.; Peetathawatchai, P.; Glenn, A.; Sivashankar, V.; Zamoshchin, D.; Glikbarg, L.; Askaryar, D.; Yang, H.; Zhang, A.; Alluri, R.; Tran, N.; Sangpisit, R.; Oseleononmen, K. O.; Boneh, D.; Ho, D. E.; and Liang, P. 2025. Cybench: A Framework for Evaluating Cybersecurity Capabilities and Risks of Language Models. In *The Thirteenth International Conference on Learning Representations*.