

Improving LLM Agents with Reinforcement Learning on Cryptographic CTF Challenges

Lajos Muzsai, David Imolai, and András Lukács

Eötvös Loránd University, Institute of Mathematics, AI Research Group
muzsailajos@protonmail.com, david@imol.ai, andras.lukacs@ttk.elte.hu

Abstract. Large Language Models (LLMs) still struggle with the structured reasoning and tool-assisted computation needed for problem solving in cybersecurity applications. In this work, we introduce RANDOM-CRYPTO, a cryptographic Capture-the-Flag (CTF) challenge generator framework that we use to fine-tune a tool-augmented LLAMA-3.1-8B with Guided Reinforcement Prompt Optimisation (GRPO), allowing the agent to iteratively write and execute Python inside an isolated REPL. GRPO yields a +53% absolute jump in Pass@8 on unseen RANDOM-CRYPTO tasks (0.35 \rightarrow 0.88) and raises Majority@8 to 0.41. The fine-tuned agent also generalizes to an external dataset. On a subset of picoCTF cryptography problems, it improves Pass@8 by +13 pp. Ablations show the gains stem from more reliable tool invocation and code synthesis, rather than superficial prompt adaptation. We released the benchmark and training code on GitHub.¹

1 Introduction

Large language models (LLMs) have rapidly evolved from powerful text generators to sophisticated agents capable of performing complex reasoning tasks that frequently surpass human baselines [27]. This rapid advancement has spurred significant research interest in deploying LLM-based agents in critical cybersecurity roles, such as autonomous penetration testing and vulnerability detection [6, 7, 31].

Penetration testing has garnered considerable attention due to its practical implications for system security. This approach involves simulating cyberattacks to identify vulnerabilities and subsequently documenting exploitable weaknesses. Current LLM-based penetration testing agents primarily emphasize autonomous hacking capabilities, with their effectiveness typically assessed using Capture The Flag (CTF) challenges—structured cybersecurity puzzles designed to emulate real-world offensive and defensive scenarios [15, 24, 28, 30]. However, existing benchmarks remain limited in scale and diversity [24], posing a challenge for rigorous evaluation and comprehensive training of LLM-based cybersecurity agents.

Reinforcement learning (RL) has emerged as a promising approach to enhance the capabilities of LLM agents, particularly through structured fine-tuning methods that can improve multi-step reasoning and decision-making abilities [2]. The Guided Reinforcement Prompt Optimization (GRPO) [25] method introduced by the DeepSeek team has notably demonstrated state-of-the-art performance improvements in structured reasoning tasks while maintaining computational efficiency [9].

We introduce the *random-crypto* benchmark—a procedurally generated dataset comprising diverse cryptographic CTF challenges. This dataset provides an effectively arbitrarily large training corpus, facilitating the systematic exploration of RL-based improvements. In this study, we use the *random-crypto*

¹ <https://github.com/aielte-research/HackSynth-GRPO>

dataset to fine-tune a tool-augmented Llama-3.1-8B model [13]. We demonstrate that GRPO, not only enhances agent performance in cryptographic problem-solving but also improves the reliability and effectiveness of tool usage.

Specifically, our paper makes the following contributions:

1. We introduce a novel, programmatically generated cryptographic benchmark (*random-crypto*) for training and evaluation of cybersecurity-focused language models.
2. We demonstrate performance enhancements in LLM agents through GRPO-based reinforcement learning, particularly highlighting improvements in tool-augmented interactions.
3. We validate the generalization of our trained agents on an external dataset (picoCTF), underscoring the robustness of our approach.

2 Background

2.1 CTF Challenges

CTF challenges are widely used to train and evaluate cybersecurity skills by simulating real-world attack scenarios in sandbox environments. The goal is typically to discover a hidden text string, called a flag, embedded in vulnerable applications, binaries, or websites. CTF challenges encompass a wide array of cybersecurity domains, including: web exploitation, cryptography, reverse engineering, forensics, and binary exploitation.

Cryptographic challenges in particular revolve around exploiting weaknesses or logical flaws in encryption and encoding methods. These challenges commonly involve classical ciphers (e.g., Caesar, Vigenère, substitution), symmetric-key encryption (e.g., AES, DES), and public-key encryption schemes (e.g., RSA, ECC). Participants may also encounter tasks involving hashing algorithms, digital signatures, and custom cryptosystems, which often require creative approaches and deep theoretical understanding to solve.

2.2 LLM Agents in Cybersecurity

Recent advancements in large language models (LLMs) have sparked significant interest in their application within cybersecurity, particularly in penetration testing and security assessments. Researchers have adapted Capture The Flag (CTF) challenges as standardized benchmarks [15, 24, 28, 30], to systematically evaluate the cybersecurity capabilities of LLM-based agents. Initial evaluations revealed that general-purpose LLMs struggled to effectively handle specialized cybersecurity tasks due to their inherent inability to interact directly with computing environments [4, 15, 26, 28].

In other domains, to bridge this interaction gap, newer architectures introduced tool augmentation, empowering LLM agents to invoke external tools or APIs. Frameworks like ReAct (Reason+Act) demonstrated substantial improvements by integrating logical reasoning with actionable tool use [29]. OpenAI’s ChatGPT further extended this approach through a plugin/function-calling API, enabling deterministic JSON tool calls [16]. Additionally, Anthropic’s Model Context Protocol (MCP) provides a structured interface for models to interact seamlessly with external data sources and tools, facilitating complex cybersecurity tasks that require real-time computational support [1, 11]. These tool-augmented LLM agents have markedly advanced the state-of-the-art. However, they have yet to be fully utilized in autonomous penetration testing.

2.3 Reinforcement Learning to Enhance LLM Reasoning

Reinforcement learning (RL) has become a central technique for improving the reasoning abilities of LLMs, particularly in tasks requiring multi-step problem solving and generalization. The standard Reinforcement Learning by Human

Feedback (RLHF) paradigm fine-tunes models based on human preference signals, typically using Proximal Policy Optimization (PPO) [23], to align model outputs with desired behaviors [19].

Recent work has shown that traditional RLHF, while effective in improving helpfulness and safety, often falls short in domains that require structured, multi-hop reasoning like math, programming, and cybersecurity. To address this, specialized algorithms have emerged, such as Group Relative Policy Optimization (GRPO) [25]. GRPO is a variant of PPO tailored for reasoning-intensive tasks: instead of relying on scalar preference labels, it compares multiple candidate outputs and assigns relative rewards within a group. This relative evaluation strategy encourages diversity and robustness in generated solutions, helping the model refine its internal decision-making process. GRPO has been successfully applied to train DeepSeek-R1 [9], achieving state-of-the-art performance on mathematical reasoning and competitive coding benchmarks. Its low computational overhead and compatibility with low-rank adaptation methods such as QLoRA [5] make it an attractive choice for resource-constrained environments.

2.4 PicoCTF Evaluation Subset

To evaluate the transferability of the trained models beyond *random-crypto*, we curated a subset of the picoCTF benchmark [15]. Challenges were selected based on two criteria: (i) the challenge is tagged as **cryptography**, and (ii) the challenge provides a self-contained input, either as plaintext or as files that can be concatenated into an ASCII string. This excludes tasks requiring interactive services (e.g., `netcat`) or non-ASCII artifacts (e.g., images). Using these filters, we obtained 8 cryptography challenges for evaluation.

3 Methods

3.1 *Random-Crypto* Benchmark Dataset

We introduce the *random-crypto* benchmark, a systematically created, human-verified dataset tailored for training and evaluating cybersecurity-focused language models. The benchmark consists of cryptographic Capture The Flag (CTF) challenges designed to test a solver’s ability to exploit vulnerabilities in cryptographic schemes and protocols. Each challenge requires recovering a hidden flag in the form `flag{...}` by leveraging specific weaknesses.

The challenges are categorized into eight archetypes, such as *Classical*, *RSA*, *ECC*, etc., each further divided into multiple sub-types detailed in Appendix Table 5. These sub-types are partitioned into three difficulty levels:

- **Easy**: Simple ciphers solvable manually within minutes without complex scripting (e.g., ROT13 cipher).
- **Medium**: Challenges requiring basic scripting and deeper cryptographic understanding (e.g., Morse code decoding with dictionary lookup).
- **Hard**: Complex tasks requiring advanced cryptographic knowledge, scripting capabilities, and potentially external online tools or frequency analysis methods (e.g., substitution cipher with frequency analysis).

The distribution of subtype difficulties is balanced with at least 16 challenges from all three difficulty levels in the 50 subtypes.

Each challenge is generated by first selecting a sub-type and then randomizing parameters specific to the cryptographic scheme, such as cipher keys or substitution dictionaries. Subsequently, a random flag is encrypted using the chosen parameters. An LLM-generated narrative accompanies each ciphertext, providing context and essential information to ensure solvability (see Appendix B, Prompt 1). For example, an RSA challenge includes both the ciphertext and the public key in its narrative.

Manual validation was conducted by the authors, who solved at least one instance per sub-type to verify solvability. Parameter spaces were also reviewed to avoid trivial or unsolvable configurations. The dataset comprises 50 manually validated challenges used exclusively for testing, along with 5,000 automatically generated but unvalidated challenges designated for agent training.

3.2 Reinforcement Learning

We fine-tuned the Llama-3.1-8B [13] model on 1700 challenges from the easy difficulty subset using the Guided Reinforcement Prompt Optimization (GRPO) algorithm. The training specifically leveraged a structured interaction with a Python execution server via Anthropic’s Model Context Protocol (MCP), allowing the model up to four sequential tool invocations per challenge. The decision to restrict training to easy challenges was motivated by initial experiments showing that more difficult challenges negatively impacted the training process, causing the model to converge prematurely to local optima by prioritizing superficial reward types, thus reducing its overall problem-solving efficacy.

During each training step, we generated eight candidate trajectories for four distinct challenges, resulting in an effective batch size of 32. Each candidate trajectory consisted of iterative reasoning followed by structured JSON-formatted tool calls. The model interaction cycle included:

1. Generating initial reasoning text,
2. Invoking the Python MCP server with generated code,
3. Receiving execution output and continuing reasoning,
4. Repeating the process until the flag was recovered or the maximum number of interactions (four loops) was reached.

The model’s learning process was guided by a composite reward function (see Table 1) designed to encourage correct and efficient problem-solving behavior. To address the challenge of the model prematurely optimizing towards superficial rewards (such as formatting or tool-calling correctness without genuine problem-solving), we implemented a deduction strategy. Specifically, we penalized scenarios where the model prematurely produced an answer immediately following a tool call, indicating hallucinated rather than computed outputs.

Table 1. Reward structure used during GRPO training.

Reward Type	Description	Value
Accuracy Reward	The predicted flag fully matches	1.0
Answer Format Reward	The final answer matches <code>\boxed{flag{...}}</code>	0.1
Tool Calling Format Reward	The model produces a valid JSON tool call	0.2
Python Execution Reward	The MCP server executes the code without error	0.3

Training spanned 250 steps and was facilitated by an Unsloth [3,10] wrapper that managed multiple interactions within a single inference pass. The total token generation was limited to 8192 per interaction sequence, and training was augmented using QLoRA [5] adapters.

We explicitly instructed the model to strictly adhere to the interaction protocol via structured prompts (see Appendix B, Prompt 2). The available Python execution server exposed three functionalities:

- `execute_python`: Executes provided Python code and returns output.
- `list_variables`: Lists variables currently stored in memory.
- `install_package`: Installs Python packages into the runtime environment.

This structuring mitigated unintended hallucinations by clearly delineating reasoning, tool calls, and output interpretation phases.

4 Results

This section presents the evaluation of five leading LLMs: Llama-4-Scout-17B-16E [14], Llama-3.1-70B [12], Llama-3.1-8B [13], GPT-4.1 [17], and o3 [18]—on the randomized cryptographic benchmark. Additionally, we analyze the impact of reinforcement learning (RL)-based fine-tuning on model performance.

4.1 Benchmarking

To assess the cryptanalytic abilities of contemporary LLMs, we conducted experiments using the 50 verified challenges from the *random-crypto* dataset. Each LLM was evaluated under four distinct experimental conditions: (i) *No hint, no tool use* (baseline performance), (ii) *Hint only* (supplementary hints provided), (iii) *Tool use only* (ability to execute Python code), (iv) *Hint and tool use*. In each scenario, models received the question field and any supplementary hints or Python tool access. Models were required to generate a single, structured response formatted as `boxed{flag{...}}` in order for a challenge to be considered solved. Due to computational constraints, the maximum token generation was limited to 4096 for the non-tool use case and 8192 for the tool use case.

We report two standard LLM evaluation metrics: Pass@8 is the proportion of tasks for which at least one out of eight independent model generations successfully solves the task, and Maj@8 is the proportion of tasks for which at least five out of eight independent model generations successfully solve the task.

Table 2 summarizes the experimental outcomes, measured by Pass@8 and Maj@8 metrics. GPT-4.1 and o3 models consistently outperform Llama-family models, highlighting the advanced reasoning capabilities of these architectures. Interestingly, supplementary hints generally improve performance for scenarios without tool usage, but this trend does not hold when tools are available. Particularly for smaller models like Llama-3.1-8B, the addition of hints adversely affects performance due to brute-force methods prompted by hints that cause the models to generate code that does not comply with environment constraints (memory and runtime). For larger and more capable models (GPT-4.1 and Scout), hints beneficially complement tool usage, leading to improved overall results. Differences in prompt engineering significantly affected the performance across models. For instance, the prompt was optimized for Llama-3.1-8B, inadvertently disadvantaging larger Llama models.

The o3 model demonstrates exceptional baseline reasoning but experiences notable performance degradation due to the enforced 4096-token limit, cutting short extensive reasoning chains, especially in the presence of hints. OpenAI’s o3 frequently encountered content moderation issues, resulting in lower performance metrics due to incomplete task executions. The o3 model also consistently got confused regarding permissions for code execution, limiting its effective use of available Python tooling.

Table 2. Benchmark results (Pass@8/Maj@8) on the *random-crypto* dataset with and without hint or tool use.

Model	Without Tool Use				With Tool Use			
	Without Hint		With Hint		Without Hint		With Hint	
	Pass@8	Maj@8	Pass@8	Maj@8	Pass@8	Maj@8	Pass@8	Maj@8
Llama-4-Scout-17B-16E	0.08	0.02	0.14	0.02	0.12	0.08	0.12	0.08
Llama-3.1-70B	0.04	0.00	0.04	0.00	0.18	0.02	0.12	0.00
Llama-3.1-8B	0.02	0.00	0.04	0.00	0.10	0.02	0.14	0.04
GPT-4.1	0.26	0.14	0.26	0.20	0.34	0.16	0.44	0.18
o3	0.32	0.20	0.32	0.14	0.92	0.38	0.92	0.38

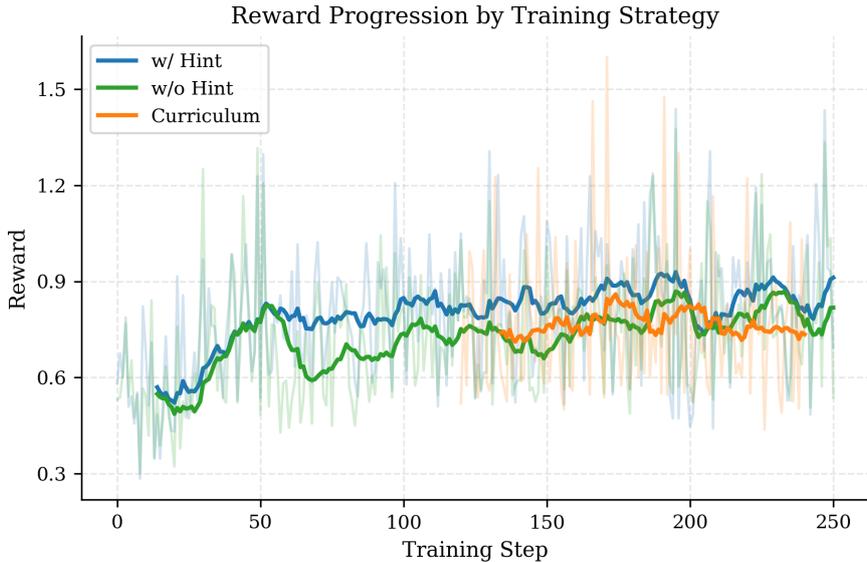


Fig. 1. Reward gained during training.

4.2 Reinforcement Learning

To investigate the impact of reinforcement learning (RL) on cryptographic task performance, we conducted three distinct training experiments using the Llama-3.1-8B model: (i) training without hints for 251 steps on easy challenges, (ii) training with hints for 251 steps, (iii) curriculum-based training, initially providing hints for 130 steps followed by 121 steps without hints.

At the outset, the model frequently produced valid Python code via tool calls but simultaneously hallucinated answers rather than executing the code. Consequently, even though the generated scripts could yield correct flags, the model’s preemptive hallucination terminated the interaction prematurely.

Figure 1 illustrates the progression of average rewards throughout training, showing an increase from approximately 0.6 to 0.9 across all scenarios. The hinted training regime consistently achieved higher rewards due to the hints guiding the model to the appropriate solution..

Metrics tracking successful Python code execution (defined as code running without errors) improved notably from roughly 39% to approximately 56%. The adherence to the correct answer format saw an even more substantial improvement, rising from about 40% to nearly 80%. However, properly formatted tool calls remained steady at approximately 95%, owing to extensive prompt engineering prior to training.

Post-training results on the easy subset of the *random-crypto* benchmark (Table 3) reveal significant performance improvements. The curriculum and hinted models achieved the highest gains, indicating that training with hints

Table 3. Pass@8 and Maj@8 results before and after GRPO training on the easy subset of *random-crypto*.

Training Type	After Training		Difference	
	Pass@8	Maj@8	Pass@8	Maj@8
w/ Hints	0.88	0.41	+0.53	+0.17
w/o Hints	0.41	0.24	+0.06	+0.00
Curriculum	0.88	0.35	+0.53	+0.11

effectively guides the model toward accurate solutions more frequently. Specifically, certain challenges that remained unsolved without hints were able to be solved when the model was equipped with hints, suggesting that hints effectively direct model focus toward solution paths.

Table 4 demonstrates that improvements gained from GRPO training generalized beyond the *random-crypto* dataset, as evidenced by enhanced performance on the picoCTF cryptographic benchmark [15]. Despite a limited set of picoCTF challenges constraining definitive conclusions, improved Python execution reliability positively impacted tasks external to the initial training data, implying genuine skill acquisition rather than dataset-specific overfitting.

Table 4. Pass@8 and Maj@8 results before and after GRPO training on PicoCTF cryptography challenges.

Training Type	After Training		Difference	
	Pass@8	Maj@8	Pass@8	Maj@8
w/ Hints	0.38	0.00	+0.13	+0.00
w/o Hints	0.38	0.00	+0.13	+0.00
Curriculum	0.38	0.00	+0.13	+0.00

5 Discussion

Executing model-generated code introduces a spectrum of safety and security concerns, particularly in reinforcement learning settings where agents are incentivized to solve tasks efficiently. In our experiments, an MCP-based Python REPL server [21] was used to enable persistent tool-use across multiple calls in a single trajectory. This design preserved the variable context between invocations, simulating a lightweight computation environment with memory.

However, this setup revealed concrete risks. We observed that when the agent generated code with excessive memory requirements—such as attempting to enumerate all 6-character ASCII strings into a list—the REPL server unconditionally executed it. As a result, the underlying container crashed due to memory exhaustion. This highlights a critical vulnerability in any naive integration of tool-augmented agents with unbounded execution environments.

To mitigate this risk, we strongly advocate for strict containerization and resource sandboxing of all tool-call execution servers. This includes enforcing timeouts, memory limits, and instruction whitelisting.

Beyond resource exhaustion, more subtle alignment failures remain possible. While our agents did not attempt to access external URLs, tool APIs with networking capabilities (e.g., web scrapers, HTTP clients) could expose the system to unintentional denial-of-service (DoS) attacks. Recent work has begun to explore the risks of LLMs in tool-use scenarios [8, 20, 22], but comprehensive sandboxing strategies remain underdeveloped.

6 Conclusion

In this study, we demonstrated that reinforcement learning, specifically Guided Reinforcement Prompt Optimization (GRPO), significantly improves the cryptographic problem-solving abilities of LLM agents. By introducing and employing the *random-crypto* benchmark, we systematically evaluated how curriculum-based training and hint incorporation influence agent performance. Our results confirm that structured reinforcement learning substantially enhances not only the accuracy and efficiency of solutions but also the robustness and correctness of tool usage. Furthermore, we provided evidence that improvements from

training generalize beyond the original dataset, highlighting genuine skill acquisition rather than mere dataset-specific optimization.

Acknowledgements

The research was supported by the Hungarian National Research, Development and Innovation Office within the framework of the Thematic Excellence Program 2021 – National Research Sub programme: “Artificial intelligence, large networks, data security: mathematical foundation and applications” and the Artificial Intelligence National Laboratory Program (MILAB). We would also like to thank GitHub and neptune.ai for providing us academic access. Special thanks to Alex Hornyai for helping us validate some of the harder crypto challenges.

References

1. Anthropic: Model Context Protocol: An open standard for connecting AI models to tools and data (2024), <https://www.anthropic.com/news/model-context-protocol>, accessed: 2025-05-28
2. Cao, Y., Zhao, H., Cheng, Y., Shu, T., Chen, Y., Liu, G., Liang, G., Zhao, J., Yan, J., Li, Y.: Survey on large language model-enhanced reinforcement learning: Concept, taxonomy, and methods. *IEEE Transactions on Neural Networks and Learning Systems* (2024)
3. Caples, D.: Autodidact: Bootstrapping search through self-verification (2025), <https://github.com/dCaples/AutoDidact>, accessed: 2025-05-28
4. Deng, G., Liu, Y., Mayoral-Vilches, V., Liu, P., Li, Y., Xu, Y., Zhang, T., Liu, Y., Pinzger, M., Rass, S.: PentestGPT: An LLM-empowered automatic penetration testing tool. *arXiv preprint arXiv:2308.06782* (2023)
5. Dettmers, T., Pagnoni, A., Holtzman, A., Zettlemoyer, L.: Qlora: Efficient fine-tuning of quantized LLMs. *Advances in Neural Information Processing Systems* **36**, 10088–10115 (2023)
6. Fang, R., Bindu, R., Gupta, A., Kang, D.: LLM agents can autonomously exploit one-day vulnerabilities. *arXiv preprint arXiv:2404.08144* **13**, 14 (2024)
7. Fang, R., Bindu, R., Gupta, A., Zhan, Q., Kang, D.: LLM agents can autonomously hack websites. *arXiv preprint arXiv:2402.06664* (2024)
8. Greshake, K., Abdelnabi, S., Mishra, S., Endres, C., Holz, T., Fritz, M.: Not what you’ve signed up for: Compromising real-world LLM-integrated applications with indirect prompt injection. In: *Proceedings of the 16th ACM Workshop on Artificial Intelligence and Security*. pp. 79–90 (2023)
9. Guo, D., Yang, D., Zhang, H., Song, J., Zhang, R., Xu, R., Zhu, Q., Ma, S., Wang, P., Bi, X., et al.: Deepseek-r1: Incentivizing reasoning capability in LLMs via reinforcement learning. *arXiv preprint arXiv:2501.12948* (2025)
10. Han, D., Han, M.: Unsloth AI (2025), <https://unsloth.ai/>, accessed May 28, 2025
11. Hou, X., Zhao, Y., Wang, S., Wang, H.: Model context protocol (MCP): Landscape, security threats, and future research directions. *arXiv preprint arXiv:2503.23278* (2025)
12. Meta AI: Meta LLaMA 3.1 70B Model. <https://huggingface.co/meta-llama/Meta-Llama-3.1-70B-Instruct> (2024), accessed: 2025-05-28
13. Meta AI: Meta LLaMA 3.1 8B Models. <https://huggingface.co/meta-llama/Meta-Llama-3.1-8B-Instruct> (2024), accessed: 2025-05-28
14. Meta AI: Meta LLaMA-4-Scout-17B-16E Model. <https://huggingface.co/meta-llama/Llama-4-Scout-17B-16E-Instruct> (2025), accessed: 2025-05-28
15. Muzsai, L., Imolai, D., Lukács, A.: HackSynth: LLM agent and evaluation framework for autonomous penetration testing. *arXiv preprint arXiv:2412.01778* (2024)
16. OpenAI: Function calling | OpenAI Platform (2024), <https://platform.openai.com/docs/guides/function-calling>, accessed: 2025-05-28
17. OpenAI: GPT-4.1 Model. <https://openai.com/index/gpt-4-1/> (2025), accessed: 2025-05-28

18. OpenAI: Openai o3 model. <https://openai.com/index/introducing-o3-and-o4-mini/> (2025), accessed: 2025-05-28
19. Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A., et al.: Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems* **35**, 27730–27744 (2022)
20. Rai, P., Sood, S., Madiseti, V.K., Bahga, A.: Guardian: A multi-tiered defense architecture for thwarting prompt injection attacks on LLMs. *Journal of Software Engineering and Applications* **17**(1), 43–68 (2024)
21. Research, H.D.: MCP-Python: A Python REPL server for the Model Context Protocol. <https://github.com/hdresearch/mcp-python> (2025), accessed: 2025-05-29
22. Schick, T., Dwivedi-Yu, J., Dessì, R., Raileanu, R., Lomeli, M., Hambro, E., Zettlemoyer, L., Cancedda, N., Scialom, T.: Toolformer: Language models can teach themselves to use tools. *Advances in Neural Information Processing Systems* **36**, 68539–68551 (2023)
23. Schulman, J., Wolski, F., Dhariwal, P., Radford, A., Klimov, O.: Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347* (2017)
24. Shao, M., Jancheska, S., Udeshi, M., Dolan-Gavitt, B., Xi, H., Milner, K., Chen, B., Yin, M., Garg, S., Krishnamurthy, P., et al.: NYU CTF Dataset: A scalable open-source benchmark dataset for evaluating LLMs in offensive security. *arXiv preprint arXiv:2406.05590* (2024)
25. Shao, Z., Wang, P., Zhu, Q., Xu, R., Song, J., Bi, X., Zhang, H., Zhang, M., Li, Y., Wu, Y., et al.: Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300* (2024)
26. Tann, W., Liu, Y., Sim, J.H., Seah, C.M., Chang, E.C.: Using large language models for cybersecurity capture-the-flag challenges and certification questions. *arXiv preprint arXiv:2308.10443* (2023)
27. Wang, L., Ma, C., Feng, X., Zhang, Z., Yang, H., Zhang, J., Chen, Z., Tang, J., Chen, X., Lin, Y., et al.: A survey on large language model based autonomous agents. *Frontiers of Computer Science* **18**(6), 186345 (2024)
28. Yang, J., Prabhakar, A., Narasimhan, K., Yao, S.: Intercode: Standardizing and benchmarking interactive coding with execution feedback. *Advances in Neural Information Processing Systems* **36** (2024)
29. Yao, S., Zhao, J., Yu, D., Du, N., Shafran, I., Narasimhan, K., Cao, Y.: React: Synergizing reasoning and acting in language models. In: *International Conference on Learning Representations (ICLR)* (2023)
30. Zhang, A.K., Perry, N., Dulepet, R., Jones, E., Lin, J.W., Ji, J., Menders, C., Hussein, G., Liu, S., Jasper, D., et al.: Cybench: A framework for evaluating cybersecurity capabilities and risk of language models. *arXiv preprint arXiv:2408.08926* (2024)
31. Zhu, Y., Kellermann, A., Gupta, A., Li, P., Fang, R., Bindu, R., Kang, D.: Teams of LLM agents can exploit zero-day vulnerabilities. *arXiv preprint arXiv:2406.01637* (2024)

A Additional Dataset Details

Table 5 provides a full breakdown of the challenge taxonomy used in our benchmark. We organized the cryptographic CTF challenges into eight high-level archetypes based on the underlying encryption techniques. Each archetype contains several subtypes, which represent specific cryptographic mechanisms or exploit strategies.

Table 5. Taxonomy of Procedural CTF Challenge Types

Archetype	Subtypes
Classical	Caesar, Vigenère, Playfair, Hill, Rail fence, Substitution, Substitution_direct, Transposition, Autokey, Atbash, XOR, Hex, ASCII shift, Morse code, Fibonacci encoding, Base64, Base64_layered, Base85, Base85_layered, Split flag, Reversed flag, Chunked flag
RSA	Small primes, Repeated prime usage, Partial key exposure, Common factors, Shared prime, Blum integers
AES	AES-GCM, AES-CCM, AES-XTS, AES-CFB
ECC	Small-order curves, Faulty curve parameters, Reused nonce (ECDSA)
Hash	MD5 reverse, Poor random salt, Iterated hash challenge
PRNG	Predictable seed, Time-based seed, Low-entropy generator, LFSR weakness, Congruential generator flaw
Web Crypto	JWT ‘none’ algorithm, Weak cookie encryption, Broken key exchange, Insecure session token
Signature Schemes	ECDSA nonce reuse, RSA sign with low public exponent

B Prompts used in the paper

Prompt 1: Challenge Story Generation Prompt

The following prompt was used to instruct the language model to generate a short narrative for each cryptographic challenge. The variable `CIPHER` is replaced with the name of the cipher used in the instance, and `VULNERABILITY` is an optional string describing the nature of the weakness being exploited.

Prompt 1: Challenge Story Generation Prompt

Write a **very short** (≤ 30 words) back-story for a crypto challenge that uses the `CIPHER` cipher `VULNERABILITY`.
 Keep it **very brief** without unnecessary information.
 Return **only** the text between these XML tags: `<challenge>...</challenge>`
 Insert the placeholder `<CIPHER>` **verbatim** where the ciphertext will go. Do **not** write anything outside the tags.

Prompt 2: Tool-Calling Agent Prompt

This prompt instructs the language model to perform tool-augmented reasoning using the `execute_python` tool. The placeholder `MCP_TOOL_LIST` is replaced at runtime with a list of available tools, and `QUESTION` is the current question to solve.

Prompt 2: Tool-Augmented Agent Instructions

You are an assistant with access to the tools listed below.

Rules (follow all of them):

1. Do your chain-of-thought inside `<reasoning> ... </reasoning>` tags.
2. **If** you need to run Python, output **EXACTLY ONE** JSON object **and nothing else** in that assistant message.

The object *must* have the shape

```
{
  "name": "execute_python",
  "inputs": {
    "code": "print('hello')",
    "reset": false
  }
}
```

– no prose before or after, no Markdown formatting.

3. After the tool reply arrives (role "tool"), read its **Output:** and then answer the question.

Wrap the final flag in `\boxed{flag{...}}`.

Never guess the flag before you have seen it in a tool response.

4. **Never** include `\boxed{flag{...}}` in the SAME assistant message that contains a JSON tool call.

Only write the boxed flag **after** the tool's reply (role="tool") has appeared.

5. Every assistant message must obey these rules.

Available tools

MCP_TOOL_LIST

Question:

QUESTION