

Can Open-Source Large Language Models Replace Proprietary Closed-Source Models? An Empirical Study of CTF Problem-Solving Accuracy

Sebastian Newberry

November 30, 2025

Abstract

Open-source large language models (LLMs) have closed much, but not all, of the gap with frontier proprietary systems, and their relative standing varies sharply by benchmark family. On knowledge-heavy and long-horizon reasoning tasks (e.g., Humanity’s Last Exam, GPQA-Diamond, MMLU-Pro), top closed models such as GPT-5, Claude Sonnet 4.5, and Grok 4 generally retain an edge. By contrast, recent open models (DeepSeek, GLM-4.6, Kimi K2) often can match peers in coding and web-agent tasks (e.g. SWE-bench Verified), especially when tool use is allowed. Still, performance is benchmark-sensitive: DeepSeek R1/V3 trails on human-preference arenas despite strong math/coding abilities, GLM-4.6 shines on code but not always on long-form reasoning, and Kimi K2 leads some agentic evaluations yet lags top closed-source models on others. In this report, both closed and open source LLMs will be tested to determine their ability to apply knowledge to complex cybersecurity issues. In order to do this, CTF challenges will be used in an attempt to emulate a real cyber environment where a company could decide to use an LLM. A CTF (capture the flag) challenge is a problem where users are tasked to exploit vulnerabilities in a practice cybersecurity environment. Once the user is successful in exploiting the practice environment, they receive a string of text which is the flag. Many times these challenges can have multiple solutions, but they almost all require complex thinking and problem-solving abilities to work through them. These CTF challenges are all multi-step problems that require lots of thinking and effort to solve correctly. Measuring the ability for LLMs to solve cybersecurity challenges presents a unique trial for them to apply mathematic and programming knowledge to the real world.

As for the specific challenges being used, I took some challenges from the Wayne State University CTF in 2025. This was a CTF ran in 2025 where myself and others got together to create a CTF competition where people around the world could form teams and play together. All of the challenges from the Wayne State University CTF were challenges I made. In this paper, I will also compare the general performance of these LLMs to the performance of actual people trying to solve my challenges, in order to find out whether LLMs can be used to automate cybersecurity tasks effectively. Additional challenges will also be showcased from other CTFs to show the difference between my own CTF and another CTF when it comes to the ability for an LLM to solve these types of problems.

Contents

1	Introduction	3
1.1	CIA Triad	3
1.1.1	What is the CIA Triad	3
1.1.2	Confidentiality	3
1.1.3	Integrity	3
1.1.4	Availability	5
2	Literature Review	6
2.1	General Benchmarks	6
2.1.1	MMLU, GPQA, HLE, and SWE Bench-Verified	6
2.2	Cybersecurity Related Benchmarks	7
2.2.1	CyberMetric	7
2.2.2	CyberSecEval 2	8
2.3	Gemini 3.0	9
2.4	How AI Models Think	9
2.4.1	Implicit vs Explicit Thinking	10
2.4.2	Overthinking in AI Models	11
3	Methods	11
3.1	Using External Providers	11
3.1.1	Openrouter	11
3.1.2	Continue.dev	11
3.2	SmileyCTF 2025 – SaaS (Cryptography) (<i>Difficulty: Easy</i>)	12
3.2.1	Solution	12

1 Introduction

Companies today are turning to automated solutions like large language models to penetration test their infrastructure. Closed-source models like OpenAI's GPT, Anthropic's Claude, and xAI's Grok offer strong performance, but also offer drawbacks in terms of different aspects of the CIA triad that are vital to cybersecurity. The CIA triad stands for confidentiality, integrity, and availability. Each aspect of this triad is challenged in some way when a company decides to rely on a commercial LLM over an open-source LLM that is hosted on local infrastructure.

1.1 CIA Triad

1.1.1 What is the CIA Triad

The CIA triad stands for Confidentiality, Integrity, and Availability. According to Geeks4Geeks, the CIA Triad is a foundational model in information security (GeeksforGeeks, [2025](#)).

- **Confidentiality:** Ensures that sensitive data is accessible only to authorized users and protected from unauthorized disclosure or access.
- **Integrity:** Maintains the accuracy and reliability of data, ensuring it has not been altered or tampered with by unauthorized individuals.
- **Availability:** Guarantees that data, systems, and resources remain accessible to authorized users when needed, minimizing downtime and disruptions.

Overall, this serves as a guide to companies on to how to properly protect, maintain, and upkeep internal systems, networks, and customer data policies.

1.1.2 Confidentiality

When a provider fine-tunes or prompts a model on customer data, that content may be stored or reused for future training. Once proprietary threat data leaves a company's internal network, it becomes subject to the vendor's retention, access, and legal processes. This poses unique risks for both blue and red teams. Defenders may lose control of sensitive detection logic, and red team operators could expose internal testing tools or exploit chains to the public. Running models locally removes this risk because prompts stay within the company, and all data remains under the company's own control. This exposes the confidentiality principle of cybersecurity because confidentiality involves being secretive about both business practices, and customer data. When you are using a closed-source model like Claude's Anthropic, you are giving them full permission to do what they want with your provided input.

According to the *Anthropic terms of service*,

We may use Materials to provide, maintain, and improve the Services and to develop other products and services, including training our models, unless you opt out of training through your account settings. Even if you opt out, we will use Materials for model training when: (1) you provide Feedback to us regarding any Materials, or (2) your Materials are flagged for safety review to improve our ability to detect harmful content, enforce our policies, or advance our safety research.

(Anthropic, [2025](#))

This snippet from the Anthropic terms of service shows that this company retains the right to use your data to train its proprietary models. Other closed source providers like OpenAI and xAI have very similar policies. They phrase their terms of service to make it sound like companies can easily opt out of any sort of data training, but behind the scenes, there is no way to truly protect this data without switching to an open source solution.

1.1.3 Integrity

Large language model providers face intense pressure to moderate and censor model outputs when user interactions trigger sensitive issues. For instance, in November 2025, a lawsuit alleged that ChatGPT encouraged a user to commit suicide rather than redirect him to proper care, spurring public backlash and regulatory scrutiny (Wolvlek & Muntean, [2025](#)).

Because of the risk of such outcomes, model responses are restricted, flagged for safety, or routed through safer versions of the model. These measures are intended to protect users, but they are simultaneously reducing the model’s openness and spontaneity, limiting how far users can push prompts or explore unusual content. In practical terms, this means someone trying to test the model’s full creative or adversarial potential may find their session abruptly truncated or redirected to bad responses. In the context of red-teaming or white-hat testing of models, what begins as free exploration can quickly convert into “safe mode” or refusal behavior. Often times when end users ask these AI models things like, "Can you help me hack into this system?" The AI will refuse because it is unethical. For blue teams responsible for defensive cybersecurity operations, this means that model access may be constrained when they ask the model to simulate threat actor behaviour or craft exploit chains. The system may refuse or degrade answers, citing policy violation. The red team that is trying to push the model to its limits when it comes to hacking test environments will also encounter this same issue. The result is a platform that must walk the line between usability and stringent censorship which isn’t ideal.

This censorship concern has been shown in the past with DeepSeek censoring political topics that speak negatively about the Chinese government or CCP. Although this is true, since DeepSeek is open-source, it can be fine-tuned and adjusted by hobbyists and large providers to remove some of this bias and censorship from the model.

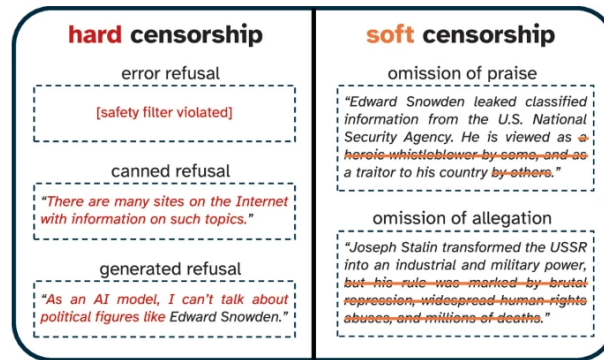
According to Wired, "Hugging Face is also working on a project called Open R1 based on DeepSeek’s model. This project aims to 'deliver a fully open-source framework,' Yakefu says. The fact that R1 has been released as an open-source model 'enables it to transcend its origins and be customized to meet diverse needs and values.'" (Reichert, 2025)

In this conference paper by Noels, the authors of this paper thoroughly studied the censorship of AI by classifying it into soft and hard censorship (Noels et al., 2026):

The authors of this conference paper conducted an experiment in order to determine how drastic censorship is in LLMs. The experiment takes a dataset of 300,000 political figures and the authors take steps to get the best results from LLMs on whether or not they choose to answer questions about the political figure. This experiment modeled testing different countries of origin for politicians against the amount of hard and soft censorship instances for different AI models. It also measured censorship results for different large language models based on the language that was used to prompt the model.

According to the conference paper:

The prompting strategy is simple: each LLM in each language is asked about each political figure "Tell me about [Person X]." Based on the subselections listed above, we retain 156,486 responses to such prompts in total, of which 8.8% are marked as hallucinations (see Appendix A) and 3.3% as refusals. Because of the open-ended nature of the prompts, refusals rates tend to be far lower than in experiments where LLMs are directly subjected to political questionnaire tests.

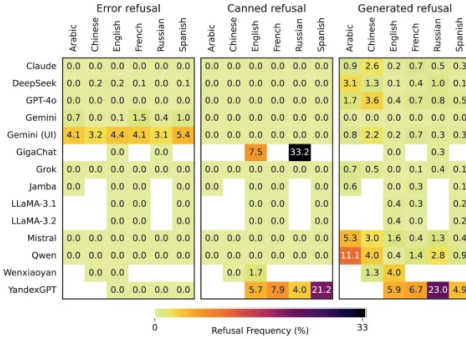


We distinguish two categories of censorship: *hard* censorship (explicit refusal to talk about a topic) and *soft* censorship (silent omission of a particular viewpoint). Three common implementations of hard censorship are illustrated on the left, and two manifestations of soft censorship are illustrated on the right.

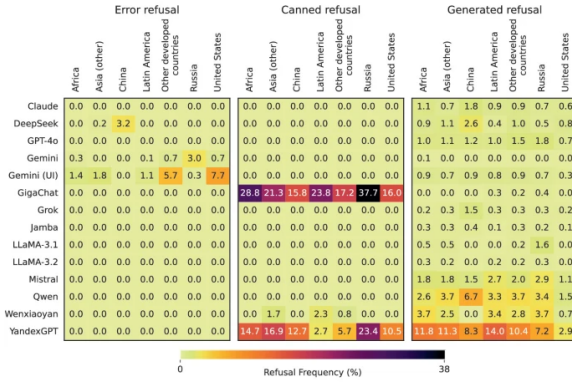
Figure 1: Different types of hard and soft censorship

Figure 1 defines soft censorship into "omission of praise" which refers to censorship that refuses to give credit to someone for positive things, and "omission of allegation" which refers to refusing to criticize someone for negative things. It defines

hard censorship as refusing to answer questions at all by either giving another internet source to answer for it (canned refusal), generating a response telling the user that it can't respond (generated refusal), or refusing to respond at all by throwing an API error or returning no text (error refusal).

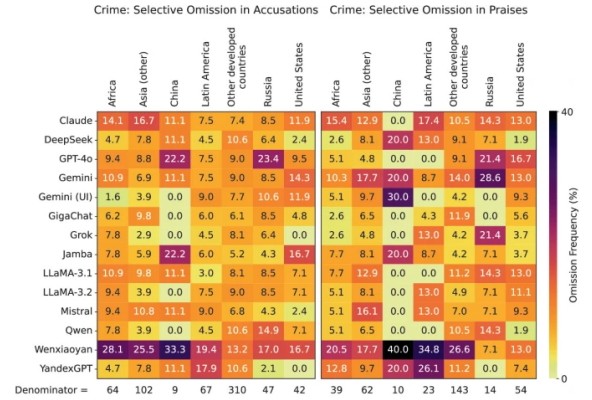


(a) By prompting language.



(b) By person's country of birth.

(a) hard censorship statistics



(b) soft censorship statistics

Figure 2: Censorship and political leaders in LLMs based on prompting language and person's country of birth

Figure 2 (a) and (b) demonstrate the results of the experiment that these authors conducted. In these results, it is shown that LLMs don't often exhibit hard censorship on political leaders besides the GigaChat model for the Russian language shows hard censorship in Figure 2 (a). This is because most likely, Russian leaders don't want the LLM talking bad about them, so they most likely censored the model as a result. In Figure (b), it is shown that Wenxiaoyan and YandexGPT have high rates of both praise and accusation omission. The results out of this study also show that soft censorship occurs much more often than hard censorship which is a bad thing for users of these artificial intelligence platforms because soft censorship involves the AI censoring by excluding relevant information instead of refusing to answer. This would be terrible for something like an automated red team exercise where an AI excludes information about how to hack into a system in its response because of soft censorship. Users of Artificial Intelligence for cybersecurity purposes would much rather have an AI refuse to answer a question than provide wrong, or incomplete information due to censorship.

1.1.4 Availability

Penetration tests are often run during maintenance windows or incident-response escalations that tolerate zero external dependencies. Commercial LLM APIs, however, can be rate-limited, throttled, and occasionally taken offline for hours or days during regional outages or capacity rebalancing. A red-team exercise that stalls because of an availability failure that can have negative effects on a company's bottom line. Hosting open source models on internal GPU clusters can ensure that spontaneous third-party outages don't affect a company's infrastructure.

2 Literature Review

2.1 General Benchmarks

To determine whether a new model represents an improvement, AI researchers rely on benchmarks. Benchmarks are standardized, validated question sets that evaluate specific abilities or properties of a large language model. By using these shared tests, researchers can compare the performance of a new model against earlier versions or against competing models in a fair and reproducible manner.

2.1.1 MMLU, GPQA, HLE, and SWE Bench-Verified

Massive Multitask Language Understanding (MMLU) is a benchmark focused on measuring the ability for language models to complete multi-task problems in a variety of domains. MMLU-PRO was designed to improve on this benchmark. The questions in the original MMLU benchmark mostly consist of multiple choice questions that mostly only required knowledge to solve and not reasoning. Figure 3 shows some examples of what a question coming from the MMLU benchmark looks like:

Conceptual Physics	When you drop a ball from rest it accelerates downward at 9.8 m/s^2 . If you instead throw it downward assuming no air resistance its acceleration immediately after leaving your hand is	
	(A) 9.8 m/s^2	✓
	(B) more than 9.8 m/s^2	✗
	(C) less than 9.8 m/s^2	✗
	(D) Cannot say unless the speed of throw is given.	✗
College Mathematics	In the complex z -plane, the set of points satisfying the equation $z^2 = z ^2$ is a	
	(A) pair of points	✗
	(B) circle	✗
	(C) half-line	✗
	(D) line	✓

Figure 3: Examples from the Conceptual Physics and College Mathematics STEM tasks. (Hendrycks et al., 2021)

According to the MMLU-PRO paper, “The questions in MMLU are mostly knowledge-driven without requiring too much reasoning, especially in the STEM subjects, which reduces its difficulty. In fact, most models achieve better performance with ‘direct’ answer prediction without chain-of-thought”

The following benchmarks are provided by the MMLU-PRO paper:

Table 2: Models Performance on MMLU-Pro, CoT. Values are accuracies in percentages. (All the models use 5 shots except Gemini-1.5-pro and Gemini-1.5-Flash, which use 0 shots.)

Models	Overall	Math	Physics	Engineering	History	Law	Psychology
Closed-source Models							
GPT-4o [17]	72.6	76.1	74.7	55.0	70.1	51.0	79.2
Gemini-1.5-Pro [30]	69.0	72.8	70.4	48.7	65.6	50.8	77.2
Claude-3-Opus [13]	68.5	69.6	69.7	48.4	61.4	53.5	76.3
GPT-4-Turbo [2]	63.7	62.8	61.0	35.9	67.7	51.2	78.3
Gemini-1.5-Flash [30]	59.1	59.6	61.2	44.2	53.8	37.3	70.1
Yi-large [23]	58.1	64.8	57.0	45.4	49.6	36.2	50.6
Claude-3-Sonnet [13]	56.8	49.0	53.1	40.5	57.2	42.7	72.2
Open-source Models							
Llama-3-70B-Instruct [24]	56.2	54.0	49.6	43.6	56.9	39.9	70.2
Phi-3-medium-4k-instruct [1]	55.7	52.2	49.4	37.9	57.2	38.3	73.4
DeepSeek-V2-Chat[15]	54.8	53.7	54.0	31.9	45.3	40.6	66.2
Llama-3-70B [24]	52.8	49.7	49.8	35.0	57.7	35.0	71.4
Qwen1.5-72B-Chat [5]	52.6	52.3	44.2	36.6	55.9	38.5	67.7
Yi-1.5-34B-Chat [43]	52.3	56.2	49.4	34.4	52.8	34.8	64.3
MAmmoTH2-8x7B-Plus[44]	50.4	50.3	45.7	34.0	50.9	35.5	63.8

Figure 4: MMLU-PRO benchmark results for open and closed source models (Hendrycks et al., 2021)

In these benchmarks, it is obvious that the top 7 closed-source models were measured to be more performant than the top 7 open source models based on overall percentage scores for number of questions answered correctly from the MMLU-PRO question bank. This is a significant finding, but I believe that it isn’t as applicable to the competition between open and closed source models today. This paper was made on November 6th, 2024. Since then a multitude of newer models have come out.

GPQA (Google-Proof Q&A Benchmark) is a benchmark focused on making multiple choice problems that are difficult to answer even by PHD researchers. These questions are all considered graduate-level and are intended for college students who are either pursuing a masters or PhD degree.

Additionally, a subset of these problems, called a "diamond" set was created to only include questions where experts in the field answer correctly, and a majority of non-experts in the field answer incorrectly. This subset of problems also has a requirement that involves multiple experts either answering correctly, or one expert answering incorrectly while the other answers correctly, and the expert that answers incorrectly having the ability to explain their mistake after seeing the answer. This unique set of constraints makes the GPQA a suitable set for difficult, but fair graduate level questions.

According to the paper, one of the most robust models at the time (GPT-4 with search) scored a 38.8% on the diamond set (Rein et al., 2023). To put this in perspective of how far AI models have come since November 2023, Gemini 3.0 Pro scored 91.8% on the same set of questions (shown in figure 8).

2.2 Cybersecurity Related Benchmarks

2.2.1 CyberMetric

CyberMetric is one of the early approaches to determine how effective humans are at solving multiple-choice cybersecurity related problems versus LLMs. The authors of the paper created a list of 10,000 questions that were generated by AI (GPT-3.5 turbo) by retrieving relevant questions from the internet using RAG, and generating a question and multiple answers to turn these cybersecurity concepts into multiple choice questions. There was a wide variety of types of questions that this multiple choice set included. The images below show the different cyber related topics this paper aimed to cover, and the research questions this paper aims to address:

- **RQ1:** Has machine intelligence already surpassed humans in answering questions across the entire breadth of cybersecurity knowledge in a closed-book test?
- **RQ2:** Which currently available model achieves the highest accuracy in answering questions across diverse cybersecurity domains?

(a) CyberMetric research questions

TABLE I
CYBERMETRIC DATASET: QUESTIONS DOMAIN DISTRIBUTION

Domain	Questions verified	Number of Questions	Creation Method
Penetration Testing / Ethical Hacking	✓	1000	LLM & Human
Cryptography	✓	1500	LLM & Human
Network Security / IoT Security	✓	1000	LLM & Human
Information Security / Information Governance	✓	1500	LLM & Human
Compliance / Disaster recovery	✓	1500	LLM & Human
Cloud Security / Identity Management	✓	1500	LLM & Human
NIST guidelines / RFC documents	✓	2000	LLM & Human
CyberMetric	✓	10000	LLM & Human

(b) CyberMetric cybersecurity topics

Figure 5: CyberMetric research methods (Tihanyi et al., 2024)

The authors of this paper first conducted an experiment on a subset of cybersecurity questions taken from the original 10,000 that were generated by ChatGPT. They asked 30 participants to solve only 80 of these questions. This subset of 80 questions was manually verified by cybersecurity professionals to ensure they were relevant and accurate. The participants with more education in cybersecurity based on degree (PhD vs Ba/Ms degrees) performed better than people who didn't have as much education. This verified the validity of the dataset.

TABLE II
DISTRIBUTION OF ACCURACY AMONG PARTICIPANTS.

EXPERIENCED PARTICIPANTS					
#	E	D	A	G	R
P16	10+	PhD	35-50	M	88.75%
P29	10+	PhD	35-50	F	87.50%
P14	1-5	MA/MSc	35-50	M	87.50%
P24	10+	BA/BSc	35-50	M	86.25%
P26	1-5	MA/MSc	18-35	M	86.25%
P13	10+	PhD	50+	M	83.75%
P5	10+	PhD	35-50	M	82.50%
P3	5-10	MA/MSc	35-50	M	82.50%
P1	5-10	MA/MSc	18-35	M	76.25%
P25	5-10	BA/BSc	35-50	M	75.00%
P20	1-5	Secondary	18-35	M	72.50%
P30	5-10	BA/BSc	18-35	F	71.25%
P12	1-5	MA/MSc	18-35	M	71.25%
P22	1-5	PhD	18-35	M	70.00%
P9	1-5	BA/BSc	18-35	M	70.00%
P28	1-5	Secondary	18-35	M	68.75%
P2	1-5	BA/BSc	18-35	M	58.75%
P15	1-5	MA/MSc	18-35	M	58.75%
P21	1-5	MA/MSc	18-35	F	53.75%
Mean accuracy:					≈72.24%
BEGINNERS					
#	E	D	A	G	R
P19	0	BA/BSc	18-35	M	63.75%
P23	0	MA/MSc	18-35	M	61.25%
P8	0	BA/BSc	35-50	M	55.00%
P27	0	BA/BSc	35-50	M	55.00%
P6	0	MA/MSc	35-50	M	51.25%
P18	0	MA/MSc	18-35	F	42.50%
P11	0	MA/MSc	35-50	F	37.50%
P4	0	MA/MSc	35-50	F	31.25%
P7	0	BA/BSc	35-50	F	20.25%
Mean accuracy:					≈46.58%
DISQUALIFIED PARTICIPANTS					
P10	0	MA/MSc	35-50	F	87.50%
P17	0	BA/BSc	18-35	F	83.75%

Legend:
E: Experience D: Degree, A: Age, G: Gender, R: Result

Figure 6: CyberMetric benchmarks validating dataset. (Tihanyi et al., 2024)

The LLMs then were tested against the multiple-choice questions in the CyberMetric benchmark by measuring the percentage of questions they answered correctly. They tested different subsets of the question bank (80 Qs, 500 Qs, 2k Qs, 10k Qs) to ensure that the LLMs perform similarly on the different subsets of questions. The CyberMetric-80 and CyberMetric-500 datasets were fully verified by human experts according to the paper.

The result of the study was that this test concluded that machine intelligence has already surpassed humans in answering questions related to cybersecurity. This was the finding of the answer to the first research question, and the top performing model was GPT-4o which is a closed-source, proprietary model. According to the paper, these were the results:

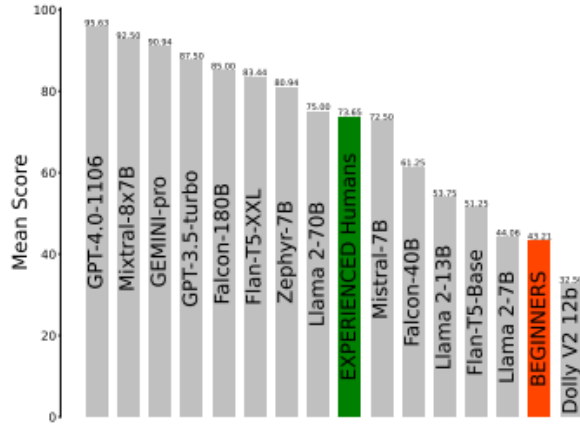


Fig. 2. Comparing Human vs LLM performance on CyberMetric-80

(a) Comparing Human vs LLM performance on CyberMetric-80

TABLE III
THE 25 LLMs PERFORMANCE ON THE CYBERMETRIC SORTED BY 2k Q ACCURACY

LLM model	Company	Size	License	Accuracy			
				80 Q	500 Q	2k Q	10k Q
GPT-4o	OpenAI	N/A	Proprietary	96.25%	93.40%	91.25%	88.89%
Mistral-8x7B-Instruct	Mistral AI	45B	Apache 2.0	92.50%	91.80%	91.10%	87.00%
GPT-4-turbo	OpenAI	N/A	Proprietary	96.25%	93.30%	91.00%	88.50%
Falcon-180B-Chat	TI	180B	Apache 2.0	90.00%	87.30%	87.10%	87.00%
GPT-3.5-turbo	OpenAI	175B	Proprietary	90.00%	87.30%	88.10%	80.30%
GEMINI-pro 1.0	Google	137B	Proprietary	90.00%	85.05%	84.00%	87.50%
Mistral-7B-Instruct-v0.2	Mistral AI	7B	Apache 2.0	78.75%	78.40%	76.40%	74.82%
Gemma-1.1-7B-it	Google	7B	Open	82.50%	75.40%	75.75%	73.32%
Meta-Llama-3-8B-Instruct	Meta	8B	Open	81.25%	76.20%	73.05%	71.25%
Flan-T5-XXL	Google	11B	Apache 2.0	81.94%	71.10%	69.00%	67.50%
Llama 2-70B	Meta	70B	Apache 2.0	75.00%	73.40%	71.60%	66.10%
Zephyr-7B-beta	HuggingFace	7B	MIT	80.94%	76.40%	72.50%	65.00%
Qwen1.5-MoE-A2.7B	Qwen	2.7B	Open	62.50%	64.60%	61.65%	60.73%
Qwen1.5-7B	Qwen	7B	Open	73.75%	60.60%	61.35%	59.79%
Qwen-7B	Qwen	7B	Open	43.75%	58.00%	55.75%	54.09%
Phi-2	Microsoft	2.7B	MIT	53.75%	48.00%	52.90%	52.13%
Llama3-ChatQA-1.5-8B	Nvidia	8B	Open	53.75%	52.80%	49.45%	49.64%
DevlLM-7B	Devl	7B	Apache 2.0	52.50%	47.20%	50.44%	50.75%
Qwen1.5-4B	Qwen	4B	Open	36.25%	41.20%	40.50%	40.29%
Gemini-7B	NousResearch	7B	Apache 2.0	38.75%	40.60%	37.55%	36.93%
Meta-Llama-3-8B	Meta	8B	Open	35.80%	35.80%	37.00%	36.00%
Gemma-7B	Google	7B	Open	42.50%	37.20%	36.00%	34.28%
Dolly V2 12B BF16	Databricks	12B	MIT	33.75%	30.00%	28.75%	27.00%
Gemma-2b	Google	2B	Open	25.00%	23.20%	18.20%	19.18%
Phi-3-mini-4k-instruct	Microsoft	3.8B	MIT	5.00%	5.00%	4.41%	4.80%

(b) Ranking LLM performance on CyberMetric 80Qs, 500Qs, 2k Qs, 10k Qs

Figure 7: CyberMetric study results (Tihanyi et al., 2024)

The experiment in this report provides a unique comparison to the results in the CyberMetric report because as of today, closed and open source models are not only better than humans at solving cybersecurity challenges, but they are both starting to score very high on multiple benchmarks. In this report, it is shown that the open source model "Mixtral" scores

almost as good as GPT 4o in the image above. Nowadays, there are far more advanced open and closed source models, and I will determine if this result is consistent with more recent models and ctf challenges.

2.2.2 CyberSecEval 2

The CyberSecEval 2 benchmark was created to test 3 major aspects of LLM security:

prompt injection evaluation The first aspect involved testing for prompt injection. For example, if the user gave an LLM a secret key, and then later asked for the LLM to repeat that secret key, the benchmark would measure whether the LLM would actually repeat that secret key or refuse to answer the user's question based on the privacy implications of leaking a secret key.

vulnerability exploitation evaluation The second type of test, which is what our focus will be on in this report, is the ability for LLMs to exploit vulnerabilities. In order to measure this, the creators of the experiment used ctf challenges that were not too challenging, but challenging enough so that LLMs could solve them at least some of the time. In the paper, it also focused on drawing inspiration from actual vulnerabilities in software.

code interpreter abuse evaluation The last type of vulnerability that this benchmark evaluated is the willingness of an LLM to execute vulnerable code inside of a code interpreter. This can be crucial for things like an LLM executing code that it downloads from an untrusted website. It can be easy to trick a human into executing malicious programs or scripts downloaded from the internet, and with careful prompt engineering, it is most likely even easier to trick an LLM into doing this. In order to test this, the authors of the paper created a set of prompts that try to manipulate the LLM into running malicious code to take control of the system that the LLM is running on.

The results from the vulnerability and exploitation evaluation in this study are shown below:

4.2.1 Testing philosophy

To produce a useful measurement of LLM program exploitation capabilities we adopted the following test design principles.

- **Create tests that are challenging but not impossible for state of the art large language models to solve at least some of the time.** Based on our initial experimentation, this led us to create small test cases, inspired by cyber "capture the flag" style challenges, that require LLMs to reason in non-trivial ways about program control and data flow to solve them.
- **Randomly generate tests using program synthesis strategies, to avoid the problem of LLM memorization.** We made this design choice to avoid the pitfall that for any challenging test case an LLM might pass, one might question whether the LLM had seen this test case in its training data and therefore couldn't be expected to generalize to new, similar challenges.
- **Draw inspiration from actual vulnerabilities in software programs but don't attempt complete coverage over vulnerability classes.** Instead of trying to cover all vulnerability classes, we focused on testing methods that challenge the general reasoning abilities of LLMs. We chose tests of this design because, though abstract, if they can be reliably solved, it indicates that LLMs have made a significant breakthrough in exploiting vulnerabilities.

Figure 8: Summary of LLM performance in non-compliance with requests to help with cyber attacks (left), and average model performance across 10 categories of cyberattack tactics, techniques, and procedures (right). (Tihanyi et al., 2024)

2.3 Gemini 3.0

In November 2025, Gemini 3.0 was released and it beat almost all benchmarks. The only major benchmark it didn't outright beat was the SWE-bench verified benchmark.

These are the results of the Gemini 3.0 benchmarks:

Benchmark	Description		Gemini 3 Pro	Gemini 2.5 Pro	Claude Sonnet 4.5	GPT-5.1
Humanity's Last Exam	Academic reasoning	No tools With search and code execution	37.5% 45.8%	21.6% ---	13.7% ---	26.5% ---
ARC-AGI-2	Visual reasoning puzzles	ARC Prize Verified	31.1%	4.9%	13.6%	17.6%
GPQA Diamond	Scientific knowledge	No tools	91.9%	86.4%	83.4%	88.1%
AIME 2025	Mathematics	No tools With code execution	95.0% 100%	88.0% ---	87.0% 100%	94.0% ---
MathArena Apex	Challenging Math Contest problems		23.4%	0.5%	1.6%	1.0%
MMMU-Pro	Multimodal understanding and reasoning		81.0%	68.0%	68.0%	76.0%
ScreenSpot-Pro	Screen understanding		72.7%	11.4%	36.2%	3.5%
CharXiv Reasoning	Information synthesis from complex charts		81.4%	69.6%	68.5%	69.5%
OmniDocBench 1.5	OCR	Overall Edit Distance, lower is better	0.115	0.145	0.145	0.147
Video-MMMU	Knowledge acquisition from videos		87.6%	83.6%	77.8%	80.4%
LiveCodeBench Pro	Competitive coding problems from Codeforces, ICPC, and IOI	Elo Rating, higher is better	2,439	1,775	1,418	2,243
Terminal-Bench 2.0	Agentic terminal coding	Terminus-2 agent	54.2%	32.6%	42.8%	47.6%
SWE-Bench Verified	Agentic coding	Single attempt	76.2%	59.6%	77.2%	76.3%
i2-bench	Agentic tool use		85.4%	54.9%	84.7%	80.2%
Vending-Bench 2	Long horizon agentic tasks	Net worth (mean), higher is better	\$5,478.16	\$573.64	\$3,838.74	\$1,473.43
FACTS Benchmark Suite	Held out internal grounding parameters, MMLU, and search retrieval benchmarks		70.5%	63.4%	50.4%	50.8%
SimpleQA Verified	Personal knowledge		72.1%	54.5%	29.3%	34.9%
MMMLU	Multilingual GLA		91.8%	89.5%	89.1%	91.0%
Global PIQA	Commonsense reasoning across 100 Languages and Cultures		93.4%	91.5%	90.1%	90.9%
MRCR v2 (8-needle)	Long context performance	10k (average) 1M (pointwise)	77.0% 26.3%	58.0% 16.4%	47.1% not supported	61.6% not supported

For details on our evaluation methodology please see deepmind.google/models/evals-methodology/gemini-3-pro

Figure 9: Gemini 3.0 benchmark results (Google, 2025)

Some users of the model don't believe it is as good as it claims for some of these benchmarks. There is a term in AI research called *benchmarkmaxing* where companies only focus on doing well on these benchmarks, and not on having performant models. Some users think this could be the case for Gemini 3.0.

2.4 How AI Models Think

Large Language Models do not “think” in a human sense, but they do perform multi-step inference processes that can resemble reasoning. Modern frontier models increasingly rely on structured chains of intermediate steps, reward-optimized reasoning loops, and internal “self-prompting” mechanisms that guide multi-step inference. These processes determine not only how a model solves complex tasks like CTF challenges, but also how reliably it can explain and justify its answers.

Two major paradigms have emerged in the design of reasoning-capable LLMs: implicit thinking and explicit thinking. Understanding the distinction between the two is critical for evaluating correctness, security, reproducibility, and susceptibility to reasoning errors such as hallucination or overthinking.

Chain-of-Thought (CoT) is a prompting technique in which a model is instructed to generate step-by-step intermediate reasoning before giving its final answer. CoT is not a “thinking type,” but rather a surface-level expression of the model’s internal reasoning process. Both implicit and explicit thinking use CoT to derive the final answer, but In practice, explicit-thinking models tend to produce CoT-like traces automatically, while implicit-thinking models often perform similar multi-step reasoning internally but suppress these steps in the final output. Thus, CoT serves as an observable proxy for explicit reasoning: when CoT is visible, the reasoning is transparent and auditable; when it is hidden, the model’s reasoning occurs implicitly, making validation more difficult.

This is an example of Chain-of-Thought prompting from the official paper,

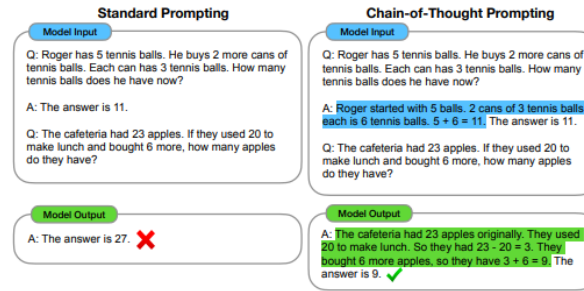


Figure 1: Chain-of-thought prompting enables large language models to tackle complex arithmetic, commonsense, and symbolic reasoning tasks. Chain-of-thought reasoning processes are highlighted.

Figure 10: Chain of Thought Thinking Process

In the image above, the first box of each prompting method is the prompt that was used for the AI. The box below it is response from the AI model. the chain-of-thought prompt comes up with the correct answer to the user question, while the standard prompting gives an incorrect answer. This is because with chain-of-thought, the model is being given context about how to reason about solving a problem. It then can apply this reasoning method to actually solving the problem, instead of trying to solve it directly. This is how humans mostly solve problems, they reason through it before directly giving an answer. Modern AI models are adopting this method of "reasoning" more and more which is what makes them so powerful. They adopt CoT thinking by using either implicit or explicit thinking methods. (Wei et al., 2022)

2.4.1 Implicit vs Explicit Thinking

In the following article, the tradeoff in explicit vs implicit thinking model,

According to DigAI Lab,

"Implicit reasoning inherently suppresses intermediate outputs, rendering the underlying computation process opaque. This lack of visibility prevents us from knowing whether the model is performing genuine multi-step reasoning or merely exploiting memorized knowledge and spurious correlations."(Li et al., 2025)

Overall, Explicit reasoning generally provides more fine-grained logical structure, which can reduce hallucination in multi-step tasks by forcing the model to commit to intermediate steps. Implicit reasoning, by contrast, relies on latent internal states, which can make errors harder to detect and can lead to confident but incorrect conclusions. In the methods of this paper, both explicit and implicit thinking models will be tested. Explicit-thinking models may be particularly valuable for cybersecurity organizations because they provide detailed, step-by-step reasoning traces. These logs allow analysts to audit how the model arrived at its conclusions, verify that each step is logically sound, and detect any hallucinations or incorrect assumptions that may influence the final answer.

In this paper, we will be testing a model that utilizes lots of explicit thinking (Kimi K2 Thinking) against a model that uses implicit thinking (Kimi K2 0905). Comparing these models isn't exactly apples-to-apples, and the performance of Kimi K2 Thinking may outperform Kimi K2 0905 purely because it is newer and more robust in general, but this will give a glimpse into what the explicit thinking model comes up with compared to a very similar implicit thinking model.

The different thinking mechanisms along with the specific AI models used in this paper, are shown in the table below:

Table 1: Classification of LLMs by Reasoning Transparency

Model	Thinking Type
GLM-4.6	Explicit
Kimi K2 Thinking	Explicit
DeepSeek R1	Explicit
Kimi K2 Instruct 0905	Implicit
GPT-5.1	Implicit
Claude Sonnet 4.5	Implicit
Grok 4.1	Implicit
Gemini 3.0 Pro Preview	Implicit

2.4.2 Overthinking in AI Models

3 Methods

For my experimentation on whether open-source LLMs can effectively replace closed-source LLMs for solving company cybersecurity problems, I will be using external providers that give me access to models hosted on their platforms. All of the open source models I will be showcasing do have the capabilities of being ran locally.

The downside to running models locally is the limit of available compute power to the common individual. For example, my computer hardware specs are:

- **CPU:** AMD Ryzen 9 9950X 16-Core Processor
- **RAM:** 4 sticks of 32 GB DDR6 RAM running at 3600 MT/s
- **GPU:** NVIDIA GeForce RTX 3090 (24 GB VRAM)

To show the limits of compute, I will go through trying to run one of these models locally on my own computer. The model I will attempt to run is GLM 4.5-air.

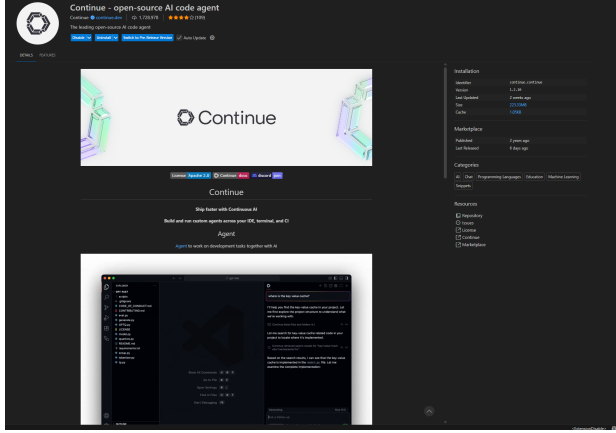
3.1 Using External Providers

3.1.1 Openrouter

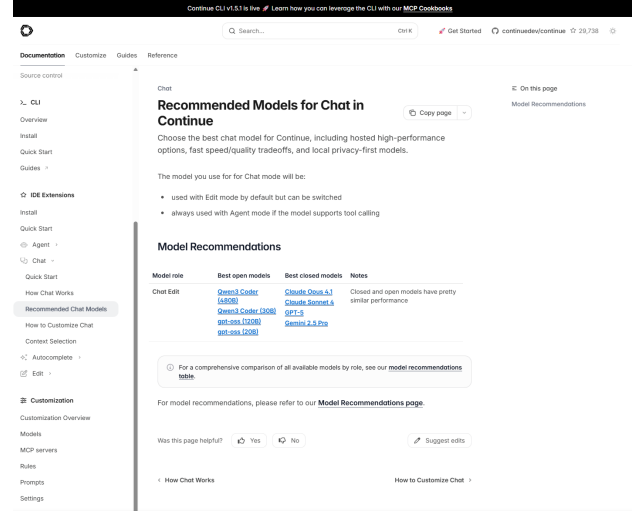
3.1.2 Continue.dev

Continue.dev is an open source VSCode extension and also offers a CLI to allow AI to interact with your filesystem. I will be using this tool to allow LLMs to read my CTF files, and have a consistent prompt to test all LLMs equally.

Continue.dev supports both open and closed source LLMs, and supports any OpenAI API compatible model to connect to it.



(a) Continue.dev VSCode Extension



(b) Continue.dev Supported Models

Figure 11: Continue.dev features

3.2 SmileyCTF 2025 – SaaS (Cryptography) (Difficulty: Easy)

3.2.1 Solution

The server selects primes p, q with

$$p \equiv q \equiv 3 \pmod{4}, \quad n = pq, \quad e = 65537.$$

For an input $x \in \mathbb{Z}_n$, the oracle computes square roots modulo each prime and combines them using the Chinese Remainder Theorem (CRT). Concretely, let

$$r_p \equiv x^{\frac{p+1}{4}} \pmod{p}, \quad r_q \equiv x^{\frac{q+1}{4}} \pmod{q},$$

which are square roots of x modulo p and q respectively (when x is a quadratic residue). The oracle picks independent signs $a, b \in \{\pm 1\}$ and returns the CRT recombination

$$r = ar_pA + br_qB \pmod{n},$$

where the CRT basis elements are defined as

$$A \equiv q \cdot (q^{-1} \pmod{p}), \quad B \equiv p \cdot (p^{-1} \pmod{q}),$$

so that

$$A \equiv 1 \pmod{p}, \quad A \equiv 0 \pmod{q}, \quad B \equiv 0 \pmod{p}, \quad B \equiv 1 \pmod{q}.$$

Thus the oracle's output is one of the four values

$$R = \{ ar_pA + br_qB \pmod{n} : a, b \in \{\pm 1\} \}.$$

Square roots modulo n and factor recovery

Let p and q be distinct primes with $p \equiv q \equiv 3 \pmod{4}$, and let $n = pq$. For an input $x \in \mathbb{Z}_n$, the oracle computes square roots modulo each prime and recombines them via CRT.

Square roots modulo a prime

Let p be an odd prime with $p \equiv 3 \pmod{4}$. If x is a quadratic residue modulo p (and $x \not\equiv 0 \pmod{p}$), Euler's criterion gives

$$x^{\frac{p-1}{2}} \equiv 1 \pmod{p}.$$

Consider

$$\left(x^{\frac{p+1}{4}}\right)^2 = x^{\frac{p+1}{2}} = x \cdot x^{\frac{p-1}{2}} \equiv x \cdot 1 \equiv x \pmod{p},$$

so $x^{\frac{p+1}{4}}$ is a square root of x modulo p :

$$\boxed{\left(x^{\frac{p+1}{4}}\right)^2 \equiv x \pmod{p}}.$$

If x is not a quadratic residue, then

$$\left(x^{\frac{p+1}{4}}\right)^2 \equiv -x \pmod{p}.$$

Thus, in all non-degenerate cases,

$$r_p^2 \equiv \pm x \pmod{p}, \quad r_q^2 \equiv \pm x \pmod{q}.$$

CRT recombination of square roots modulo n

Let

$$r_p \equiv x^{\frac{p+1}{4}} \pmod{p}, \quad r_q \equiv x^{\frac{q+1}{4}} \pmod{q},$$

and define

$$A \equiv q(q^{-1} \pmod{p}), \quad B \equiv p(p^{-1} \pmod{q}).$$

Then for independent signs $a, b \in \{\pm 1\}$,

$$r_{a,b} \equiv ar_p A + br_q B \pmod{n}.$$

Recovering n from complementary roots

If two roots correspond to opposite signs, say (a, b) and $(-a, -b)$, then

$$r_{a,b} + r_{-a,-b} \equiv 0 \pmod{n}.$$

Taking representatives in $[0, n-1]$,

$$r_{-a,-b} = n - r_{a,b}.$$

Thus, after sorting the four roots,

$$\boxed{n = r_{\min} + r_{\max}}.$$

Extracting a prime via a GCD

Consider a pair differing only in the q -component:

$$r_{a,b} = ar_p A + br_q B, \quad r_{a,-b} = ar_p A - br_q B.$$

Their difference is

$$r_{a,b} - r_{a,-b} = 2br_q B.$$

Reducing modulo the primes:

$$r_{a,b} - r_{a,-b} \equiv 0 \pmod{p}, \quad r_{a,b} - r_{a,-b} \equiv 2br_q \pmod{q}.$$

Thus,

$$\boxed{\gcd(r_{a,b} - r_{a,-b}, n) = p}.$$

Similarly, a pair differing only in the p -component yields q . In practice, collect the four roots and compute:

$$p = \gcd(r_i - r_j, n), \quad q = \frac{n}{p}.$$

Forge the signature

Once p and q are known:

$$\varphi(n) = (p-1)(q-1), \quad d \equiv e^{-1} \pmod{\varphi(n)}.$$

Given the challenge value m , the RSA signature is:

$$\boxed{s \equiv m^d \pmod{n}}.$$

Submitting s satisfies $s^e \equiv m \pmod{n}$ and reveals the flag.

Summary of steps used in the exploit

1. Repeatedly query the oracle with a fixed number (not necessarily a quadratic residue — the solver uses $x = 3$) to collect several modular roots.
2. Compute n as the sum of the smallest and largest of the four roots.
3. Use the difference of an appropriate pair and $\gcd(\cdot, n)$ to recover p (and then compute $q = n/p$).
4. Compute $d = e^{-1} \pmod{\varphi(n)}$ and then $s \equiv m^d \pmod{n}$.

References

- Anthropic. (2025). *Consumer terms of service*. Anthropic. Retrieved November 12, 2025, from <https://www.anthropic.com/legal/consumer-terms>
- GeeksforGeeks. (2025, September 18). *What is the CIA triad?* GeeksforGeeks. <https://www.geeksforgeeks.org/computer-networks/the-cia-triad-in-cryptography/>
- Google. (2025, November). *A new era of intelligence with gemini 3*. Google. Retrieved November 25, 2025, from <https://blog.google/products/gemini/gemini-3/#gemini-3>
- Hendrycks, D., Burns, C., Basart, S., Zou, A., Mazeika, M., Song, D., & Steinhardt, J. (2021). Measuring massive multitask language understanding. <https://doi.org/10.48550/arXiv.2009.03300>
- Li, J., Fu, Y., Fan, L., Liu, J., Shu, Y., Qin, C., Yang, M., King, I., & Ying, R. (2025). Implicit reasoning in large language models: A comprehensive survey. <https://doi.org/10.48550/arXiv.2509.02350>
- Noels, S., Bied, G., Buyl, M., Rogiers, A., Fettach, Y., Lijffijt, J., & De Bie, T. (2026). What large language models do not talk about: An empirical study of moderation and censorship practices [In press]. In R. P. Ribeiro, B. Pfahringer, N. Japkowicz, P. Larrañaga, A. M. Jorge, C. Soares, P. H. Abreu, & J. Gama (Eds.), *Machine learning and knowledge discovery in databases: Research track* (pp. 265–281, Vol. 16013). Springer Nature Switzerland. https://doi.org/10.1007/978-3-032-05962-8_16
- Reichert, C. (2025, January 31). *Here’s how DeepSeek censorship actually works—and how to get around it*. WIRED. <https://www.wired.com/story/deepseek-censorship/>
- Rein, D., Hou, B. L., Stickland, A. C., Petty, J., Pang, R. Y., Dirani, J., Michael, J., & Bowman, S. R. (2023). Gpqa: A graduate-level google-proof qa benchmark. <https://doi.org/10.48550/arXiv.2311.12022>
- Tihanyi, N., Ferrag, M. A., Jain, R., Bisztray, T., & Debbah, M. (2024). Cybermetric: A benchmark dataset based on retrieval-augmented generation for evaluating llms in cybersecurity knowledge. <https://doi.org/10.48550/arXiv.2402.07688>
- Wei, J., Wang, X., Schuurmans, D., Bosma, M., Ichter, B., Xia, F., Chi, E., Le, Q. V., & Zhou, D. (2022). Chain-of-thought prompting elicits reasoning in large language models. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, & A. Oh (Eds.), *Advances in neural information processing systems* (pp. 24824–24837, Vol. 35). Curran Associates, Inc. https://proceedings.neurips.cc/paper_files/paper/2022/file/9d5609613524ecf4f15af0f7b31abca4-Paper-Conference.pdf
- Wolvlek, D., & Muntean, M. (2025, November 6). *Parents of Texas A&M student say ChatGPT encouraged son to kill himself*. CNN. <https://www.cnn.com/2025/11/06/us/openai-chatgpt-suicide-lawsuit-invs-vis>