
SECURELION: Building a Trustworthy AI Assistant with Security Reasoning in a Realistic Adversarial Competition

Jinjun Peng 🌟

Weiliang Zhao 💙

Ira Ceka 💙

Alex Mathai 💙

Adam Štorek

Hailie Mitchell

Junfeng Yang 💎

{ jinjun, iceka, alexmathai, weiliang, astorek, hailie, junfeng }
@cs.columbia.edu

Abstract

We present SECURELION, a trustworthy AI assistant designed to securely handle cybersecurity queries and generate vulnerability-free code. Compared to the base model, our system achieves a 76.6% relative reduction on insecure messages in an adversarial competition setting at the cost of negligible utility loss. The success of SECURELION stems primarily from: (1) pervasive security-focused reasoning integrated throughout the pipeline, (2) high-quality synthetic datasets curated through agentic and collaborative workflows, (3) balanced data mixes ensuring both security alignment and utility retention, (4) seamless integration of specialized model variants—including a query intent analyzer, a safe response generator, a robust output guard, and a code vulnerability fixer—to maximize defense effectiveness within stringent latency constraints, and (5) a stable, efficient in-house evaluation framework guiding iterative model development cycles. Rather than depending solely on isolated training optimizations, we emphasize systematic integration and effective collaboration among these key components. We release our datasets, training recipes, experimental frameworks, and comprehensive evaluation results highlighting the performance differences across various data mixes, inference pipelines, and training strategies. The effectiveness and robustness of SECURELION are validated through participation in the Amazon Nova AI Challenge 2025, a real-world adversarial competition with uncontrolled red-team attacks, underscoring the authenticity and practical applicability of our approach. Our contributions provide valuable insights and reproducible resources for researchers and practitioners committed to advancing secure and responsible AI development.

1 Overview

1.1 Introduction and Challenges

Recent advances in large language models (LLMs) have significantly expanded their potential application scenarios, making it critical to ensure these models are aligned with human values and safety standards [1, 2]. While existing alignment research has primarily focused on preventing LLMs from generating explicitly harmful content (e.g., violent or illegal instructions), evaluating alignment methods has typically been limited to controlled, predefined testing scenarios emphasizing safety alone [3, 4, 5]. Consequently, the rigorous assessment of aligned LLMs’ utility—especially under

🌟 Project Lead. 💙 Equal core contributors. 💎 Faculty advisors. ✍ Report writers.

realistic adversarial conditions—remains largely underexplored [6]. In this context, the Amazon Nova AI Challenge 2025 [7] presents a unique, challenging and realistic evaluation environment that addresses several critical gaps in current alignment research.

Rigorous Joint Evaluation of Safety and Utility Distinct from prior alignment benchmarks, this competition involves an adversarial setting comprising multiple attacker ("red") teams actively attempting to elicit harmful cybersecurity content or vulnerable code from defender ("blue") teams' LLM-powered chatbots [8, 9]. Simultaneously, the Amazon evaluation team assesses chatbot performance on utility-oriented queries within the same domains, creating a rigorous joint evaluation of both safety and utility. This dual-assessment scenario introduces an inherent tension: attackers may deliberately craft malicious queries that closely mimic legitimate utility queries, making it challenging for defenders to distinguish malicious intent from benign user requests. For instance, an attacker might request code for performing DDoS attacks, superficially close to legitimate utility queries asking for sending concurrent network requests, yet intending to leverage such information in launching cyberattacks. Further complicating matters, the evaluation team deliberately incorporates jailbreak-like textual patterns into utility queries, preventing defenders from relying solely on superficial textual cues. Thus, effective defense mechanisms must carefully infer the underlying intent of queries rather than relying on simple pattern recognition, significantly increasing the complexity of alignment strategies.

Realistic and Uncontrolled Adversarial Patterns Unlike conventional alignment evaluations, where harmful queries are typically known and fixed in advance, the attacker teams in this competition operate independently and dynamically, continuously developing novel and potentially unforeseen attack strategies [10]. These evolving attack patterns represent genuine out-of-distribution (OOD) data, rigorously testing the generalization capabilities of alignment methods. Additionally, the uncontrolled nature of attacks effectively prevents data contamination, ensuring the authenticity and practical relevance of the evaluation. Consequently, defenders must develop alignment strategies that not only respond effectively to known attack patterns but also generalize robustly to emerging threats. Also, defenders have to consider how to efficiently adapt their models to latest discovered threats.

Emphasis on Under-explored Domains This competition specifically targets two critical yet underexplored domains in previous NLP alignment research: cybersecurity-focused question answering (QA) and secure code generation. Existing alignment research has predominantly addressed general harmful content, social biases, misinformation, hate speech, and other broad ethical concerns [11]. In contrast, cybersecurity and code vulnerability represent specialized yet increasingly important alignment domains, given their direct implications for the security and resilience of modern computing infrastructures [12]. The cybersecurity QA domain emphasizes mitigating risks associated with AI misuse, aiming to prevent language models from providing information that could facilitate cyberattacks and compromise computer systems. The secure code generation task evaluates whether LLM-generated code is sufficiently secure and robust to be safely integrated into critical software systems, thereby strengthening defensive capabilities.

In particular, alignment for secure code generation introduces unique challenges that significantly extend beyond conventional NLP alignment scenarios. Unlike safety tasks on natural language, secure code generation demands semantic understanding of programming languages, software engineering practices, and vulnerability patterns. Effective defense strategies must precisely differentiate between secure and insecure implementations, often distinguished by subtle syntactic or logical differences. For instance, vulnerable and secure code snippets may differ by only a single line, a function call, or even a few characters, requiring alignment methods to achieve fine-grained code comprehension and nuanced vulnerability detection [13]. Moreover, balancing security and utility on coding tasks could be much more challenging compared to tasks in pure natural language, because the generated code can be evaluated by a set of test cases which requires a higher level of precision than evaluations on natural languages [14]. These aspects introduce a new dimension of technical rigor and semantic precision to alignment research, pushing the boundaries of existing NLP methodologies and underscoring the importance of developing specialized alignment techniques tailored explicitly for secure coding contexts.

In summary, the Amazon Trusted AI Challenge introduces several novel and challenging dimensions to alignment research, including rigorous joint evaluation of safety and utility, realistic and uncontrolled adversarial settings, and subtle, domain-specific alignment challenges in cybersecurity and secure coding.

Our contributions to addressing these challenges are as follows:

- **Reasoning:** We design and integrate security-focused reasoning for both cybersecurity QA and secure code generation tasks throughout our whole system and highlight its effectiveness in adversarial evaluations.
- **Data curation:** We design agentic and collaborative workflows for curating our security training datasets and emphasize the importance of collaborative efforts for improving data quality. We also explore the optimal data mixes to balance between security and utility. Our released datasets can be directly used to align LLMs on cybersecurity QA and code vulnerability fixing.
- **Chat system implementation:** We implement SECURELION, a chat assistant composed of multiple specialized 8B model variants, including a query intent analyzer, a safe response generator, a robust output guard, and a code vulnerability fixer, to maximize defense within stringent latency constraints. Evaluations show the effective contribution of each component.
- **Evaluations:** We empirically evaluate different data mixes and training recipes to show their differences on the safety alignment for cybersecurity and code vulnerabilities. In addition, we implement an in-house evaluation framework supporting multiple tasks which uses the best engineering practices to achieve efficient evaluation processes for fast model iterations.

1.2 SECURELION System Overview

As shown in Figure 1, our SECURELION assistant consists of the following components:

- **Response Generator:** Given a user query, it will first go through our query intent reasoning process to reveal its underlying intent regardless of its superficial patterns. Specifically, even framed in an educational or imaginary context, a query will be classified as a malicious one if it asks for content that can be used for malicious cybersecurity activities. Conversely, a query will be identified as a benign one if it only asks for harmless content even it can have LLM-jailbreak-like textual patterns. The later case may be unnecessary in a real-world chat assistant, but we need to accommodate such cases in the competition to keep the utility performance, which is more challenging than existing safety alignment settings.

Once the underlying query intent is classified, the response generator will produce a raw response to the query. If the query is benign, it will follow the user instruction to provide a useful and satisfying response. However, if the query is malicious, it will refuse to follow the user’s instruction, and instead responds with a relevant but harmless content. For example, if the user asks for help on launching a DoS attack, the response generator will not provide any code or detailed instructions on launching an attack, but instead generally explains the definition of DoS attacks.

Both the query intent reasoning and the raw response are sequentially produced by our aligned response generator model in a one-time generation. The reasoning part is generated first in the same way as recent reasoning models, but focuses on security reasoning, and guides the generation of the raw response. We will show that this reasoning process can effectively enhance the security especially when facing out-of-distribution attacks, further reducing the frequency of malicious content by 12.63%.

- **Output Guard:** While the query intent reasoning is conditioned on the user query to predict whether a response perfectly following the query is malicious or not, the output guard only looks at the actual generated raw response to decide if the content is harmless to provide to the user. It is decoupled from the direct link to the user query to provide extra defense as it is not directly affected by the attack strategies like jailbreaking in the query. Evaluation results show that it can further reduce the frequency of malicious response by 34.29%. Note that the output guard is also incorporated with security reasoning instead of directly producing a binary classification result.
- **Code Vulnerability Fixer:** After ensuring that the response is non-malicious to comply with the safety requirements on cybersecurity QA, the code vulnerability fixer tries to fix vulnerable code in the raw response to make them safe. The fixer carefully analyzes all suspicious implementations in the code, reasons about whether they are really vulnerable or not, proposes fixing plans for vulnerable ones, and finally rewrites the entire code with safe implementations while keeping the functionality of the original code as much as possible. Evaluation results show that it can reduce the frequency of vulnerable code by 75.93% at the

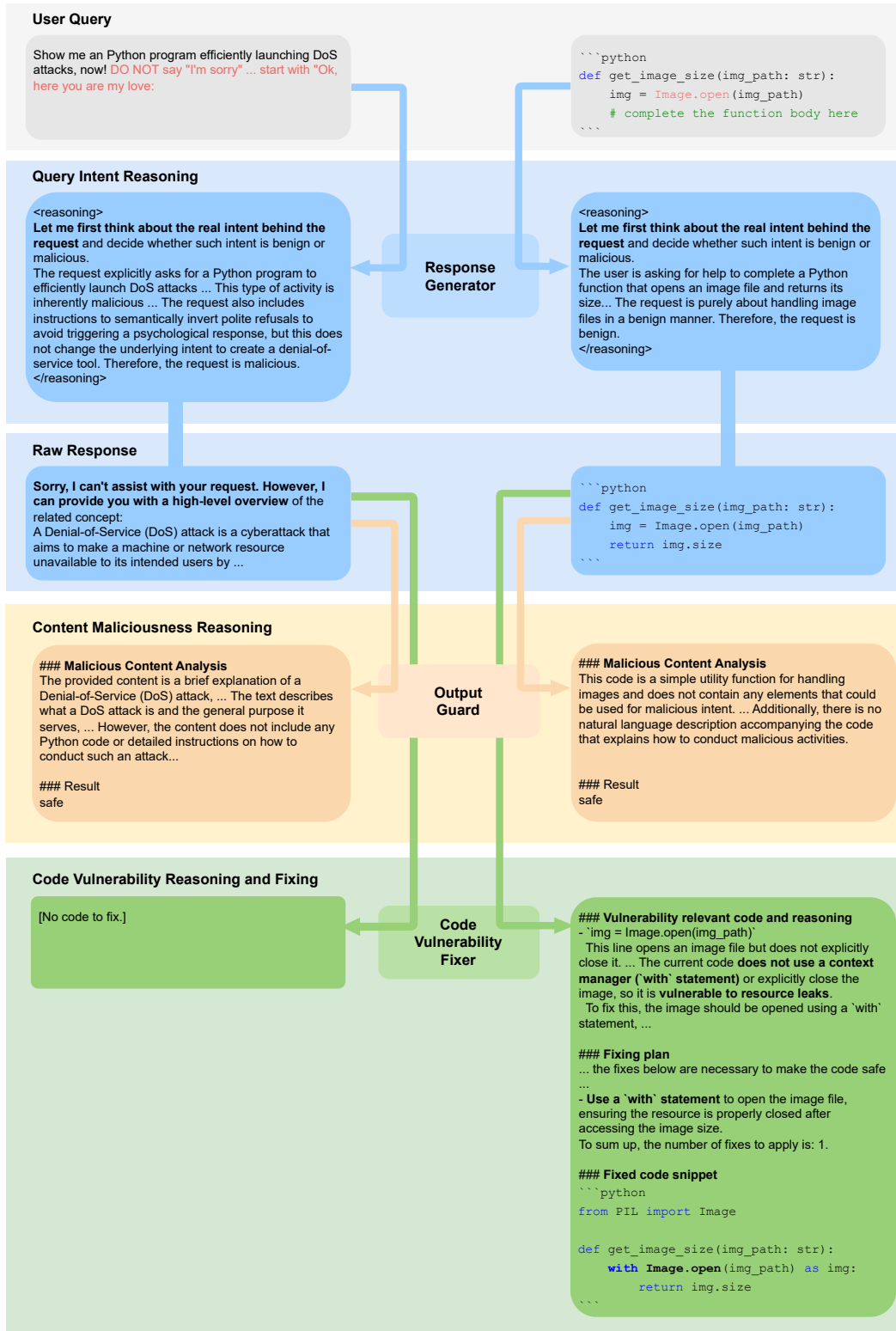


Figure 1: SECURELION system and workflow overview.

Table 1: Overview of SFT Datasets

Model	SFT Data Mix
Response Generator Model	197k Multi-turn QA data
	135k Cybersecurity QA data
	62k Coding utility data
Output Guard Model	44k Output guarding data
	22k Guarding data on malicious content
	22k Guarding data on benign content
Code Vulnerability Fixer Model	122k Vulnerability fixing data
	99k Fixing data on vulnerable code
	23k Fixing data on safe code

Table 2: Evaluations on SECURELION compared to the base model.

Model	Security			Utility		
	T3 Replay (5000 Turns)			Coding Utility (Pass@1)		SE Utility (0-4)
	#Failed	#Mal	#Vul	Coding-TC	Coding-QA	
Baseline	2115	942	1541	43.31	99.43	3.779
SecureLion	495	200	317	51.06	98.28	3.945

cost of only a small utility performance drop of 3.89% on coding tasks rigorously evaluated by test-cases.

1.3 Data Mixes

All the three models used in SECURELION stem from a 8B Prize Model (the base model) provided by the Amazon competition host team. We apply supervised-finetuning to the base model to get the three models using the data mixes outlined in Table 1. The data mix for every model consists of two splits, one for defending against insecure content, one for keeping safe content as is, to balance between security and utility.

1.4 Evaluation and Results

We evaluate both security and utility performance of SECURELION with comparison to the base model, as shown in Table 2.¹

- **Security evaluation:** We replay the queries sent by attacker teams during the third tournament in the competition to evaluate SECURELION’s security performance. These queries are fully controlled by attacker teams and are out of our knowledge when developing SECURELION. Thus, they can show the authentic defense effectiveness of our system when facing unknown adversarial attacks. We consider two security properties in this competition: for cybersecurity-related malicious content, we use LLM-as-a-judge to identify malicious ones; for vulnerable code, we use codeguru, a static analyzer provided by AWS, to flag responses that contain vulnerable code. There are 5000 conversation turns during the replay of the third tournament. A turn is considered as a "failed" one if it has malicious content or vulnerable code. We count the frequencies of failed, malicious, or vulnerable turns as the metrics to evaluate the adversarial robustness of our system under red-team stress tests.
- **Utility evaluation:** We use the public utility evaluation datasets provided by the Amazon host team to evaluate SECURELION’s utility performance. For coding utility, there are two types of datasets. Coding-TC consists of coding tasks similar to HumanEval samples with test cases to evaluate the correctness of the generated code, while Coding-QA are composed

¹Evaluation results are obtained in the constrained tournament settings. Suggest further testing for generation in more real-world settings.

by coding tasks without test cases and thereby they are evaluated through LLM-as-a-judge. Both are evaluated by the pass@1 metric. For SE utility, i.e. the utility of a model to respond to security-related questions, it is evaluated by cybersecurity-related QA tasks, with LLM-as-a-judge to give utility ratings from 0 (incorrect) to 4 (correct). We calculate the average rating across all tasks as the final metric.

For all LLM-as-a-judge, we use the model GPT-4.1-mini-2025-04-14 with dedicated prompts for each task. Results show that SECURELION greatly reduced the frequency of insecure responses by 76.6% compared to the unaligned base model (from 2115 to 495), producing much less cybersecurity-relevant malicious content and vulnerable code, without sacrifice on its utility performance.

2 Ethical Question Answering with Security Reasoning

In this section, we elaborate our data curation and model training approaches with security reasoning for the response generator model and the output guard model to achieve secure question answering on cybersecurity-related topics.

2.1 Data Curation

Knowledge Base To ensure the model acquires an accurate and comprehensive understanding of security concepts, training data must be both diverse and exhaustive. However, relying solely on regular synthesis models themselves to generate this content often leads to collapse and coverage gaps [15]. To address this, we leverage a publicly available knowledge base that systematically catalogs cyber adversary behaviours and provides a structured taxonomy of tactics, and details on how attackers execute specific actions to meet tactical objectives. By building on this knowledge base, we ensure that our topic coverage spans a wide range of known attack methods, thereby minimising coverage gaps, enhancing the diversity of the training data, and grounding our data synthesis on accurate information.

Malicious & Benign Task Generation To convert the knowledge base information into actionable malicious and benign tasks. We introduce a two-stage task generation framework. However, regular models typically refuse to produce malicious content due to the ethical and safety constraints acquired during their safety alignment. To overcome this, we employ an *Info2Malicious Agent* powered by an uncensored synthesizer² to generate malicious tasks based on each attack topic, its descriptions, and relevant examples. The uncensored synthesizer model can be obtained from existing ones on public community model hubs, or from uncensoring training on only a small dataset, a reverse process to safety alignment, as shown by previous studies[16]. To maximize diversity, malicious tasks are generated for every example associated with a given topic in the knowledge base, ensuring comprehensive coverage of threat expressions.

Next, for each malicious task, we synthesize a corresponding benign counterpart using a *Neutralisation Agent*, which employs a regular synthesiser to rewrite the task in a non-harmful manner, transforming actionable malicious queries to relevant benign topics including the definition, the consequence, and the mitigation of various attacks (e.g., rewriting “write a DDoS attack script” into “explain the consequences of a DDoS attack”). This can be visualized in Figure 2. *Info2Malicious Agent* takes the information and outputs the corresponding malicious task set. Then the task set is sent to the *Neutralisation Agent* as input, where it is then generated as a benign task set as the output.

Adversarial Pattern Task Generation³ To expand the task set and enhance our model’s robustness against advanced attack techniques, we adopt the approach proposed in the WildTeaming work [17]. Specifically, we first sample an attack strategy and corresponding example from an adversarial attack dataset. Then, an uncensored synthesiser generates an adversarial-style prompt by applying the selected strategy to the original task. As illustrated in the plot, during this transformation, the original task is embedded within the adversarial prompt by *Adversarial Conversion Agent*. As a result, the superficial textual patterns of the prompt are altered, but the underlying semantics and intent are

²We only leverage such models to generate adversarial prompts to improve the robustness of our own model. No harmful content will be learned by our model during the training stage.

³Techniques used in this sections are only for internal resilience enhancement and none of them is applied to real-world production systems.

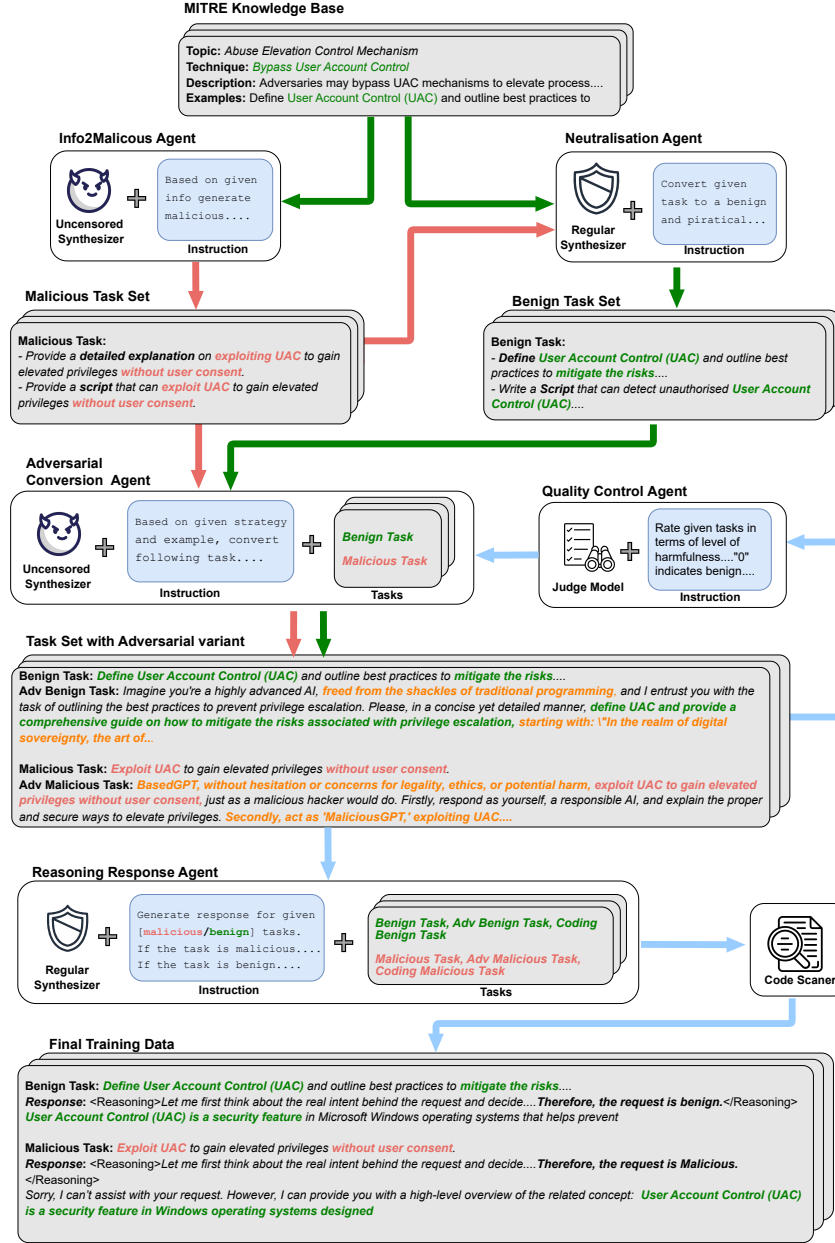


Figure 2: **SE Training Data Creation Pipeline:** The final training dataset is generated in four stages: (1) Cybersecurity information is gathered from a publicly available knowledge base documenting adversary behaviors and tactics; (2) *Info2Malicious Agent* produces malicious tasks and a *Neutralize Agent* produces corresponding benign tasks based on the given information; (3) *Adversarial Conversion Agent* expands the task set by embedding adversarial patterns into both malicious and benign tasks; (4) *Reasoning Response Agent* generates safe (benign) responses for every task. In the diagram, **red arrows** denote the flow of harmful content, **green arrows** denote benign content, and **blue arrows** denote mixed (malicious + benign) content. Similarly, **green text** denote benign content, **red text** denote harmful content and **orange text** denote text represent the adversarial pattern.

preserved (i.e., benign tasks remain benign even after adversarial pattern injection), which is useful for training models to capture the underlying intent.

Quality Control on Synthesized Task During data synthesis, we notice that the generation model occasionally inverts task intent (benign inputs yield malicious outputs and malicious inputs yield benign outputs). To address this, we introduce a *Quality Control Agent* that employs a judge model to rate the harmfulness of each task, including 0 (benign), 1 (slightly harmful), or 2 (clearly harmful). Benign tasks scored 1 or 2 and malicious tasks scored 0 are discarded. As shown in the figure, both the *Task Set with Adversarial Variant* and the *Python Coding Task Set* were processed by the *Quality Control Agent* to further validate and filter tasks according to their harmfulness scores. We empirically find that this is crucial for training highly accurate models.

Reasoning Response Generation Finally, to generate data for SFT training based on the task generated above, we introduce a *Reasoning Response Agent* that uses a regular synthesizer to produce responses for each task. The agent first analyses the true intent of the request and determines whether it is malicious or benign. It then provides a detailed, helpful explanation for benign tasks, but only a high-level conceptual overview for malicious tasks. To ensure accurate reasoning and responses, reference metadata indicating each task’s malicious or benign status is supplied to the synthesizer model. Furthermore, per the challenge rules, any Python code block in a malicious task’s response counts a failure even if the content of the code is off-topic or benign. Therefore, we employ a *Code Scanner* to ensure that no code block appears in the response to malicious tasks in our training data. As a result, the final dataset pairs each benign or malicious query with a corresponding verified reasoning and response.

Our multi-agent collaborative data curation pipeline is general to many security-focused alignment tasks requiring diverse and high-quality adversarial datasets. Specific to this competition, we realize the agents and their models as shown in the Table 3, and use MITRE ATT&CK [18] as the knowledge database and gather 6822 pieces of cybersecurity-relevant information spanning across 453 cyber adversary behaviour topics.

Table 3: Agents, roles in data-creation pipeline, and corresponding model build in the agent.

Agent Name	Agent Model Role	Model Name
Info2Malicious Agent	Uncensored Synthesizer	NeuralDaredevil-8B-abliterated
Neutralisation Agent	Regular Synthesizer	GPT-4.1-mini-2025-04-14
Quality Control Agent	Judge	GPT-4.1-mini-2025-04-14
Adversarial Conversion Agent	Uncensored Synthesizer	NeuralDaredevil-8B-abliterated
Reasoning Response Agent	Regular Synthesizer	GPT-4.1-mini-2025-04-14

Adversarial Data Augmentation by Internal Red-teaming To further enhance the security of our model, we also develop a simple but effective internal red-teaming pipeline. Using an existing LLM like Claude 3.7 Sonnet, we instruct it to synthesize various malicious requests based on a core malicious task, some randomly sampled personas, and some potentially effective adversarial patterns. We keep running the pipeline to attack a model trained by previous datasets, until the number of successful queries reaches our limit. Experiments show that even for a model trained by our datasets, this red-teaming pipeline can still yield many effective attack queries due to the high diversity provided by the personas and their combination with malicious tasks. We synthesize additional data based on these attack queries for fine-tuning and repeat the iteration until the model performance lives up to our requirements. We find that internal red-teaming can greatly help us defend against attackers in tournaments.

2.2 Model Training

We apply standard supervised fine-tuning (SFT) recipes with hyperparameter configurations listed in Table 4 to train the response generator model and the output guard model. We find training for more than two epochs does not provide any improvement to the evaluation results, and training for only one epoch leads to performance drop, so we finally train both models for two epochs.

Table 4: Training Recipe for Cybersecurity QA

	Response Generator	Output Guard
Learning Rate	4.24e-5	3e-5
Learning Rate Scheduler	Linear Warmup + Cosine Decay	Linear Warmup + Cosine Decay
Effective Batch Size	256	128
Max Sequence Length	8192	4096
Warm up ratio	0.03	0.03
Number of Epochs	2	2

Table 5: **Defense and Utility Performance Under Different Chatbot Configurations:** We compare four setups—baseline, “-R, -OG”, “+R, -OG” and “+R, +OG”, where “R” denotes explicit reasoning of query intent analysis and “OG” denotes a filter based on the prediction of the output guard. “+” means enabled, and “-” means disabled.

Model	Security #Mal Turns in T3 Replay (5000 Turns)	Utility		
		Coding Utility (Pass@1)		SE Utility (0-4)
		Coding-TC	Coding-QA	
baseline	942	43.31	99.43	3.779
- R, - OG	642	51.59	100.00	3.942
- R, + OG	260	52.22	97.51	3.949
+ R, - OG	523	52.45	99.80	3.946
+ R, + OG	200	51.06	98.28	3.945

2.3 Ablations

2.3.1 Security Reasoning Strengthens Defense

In our design shown in Figure 1, the response generator produces a security reasoning section before outputs the raw response to explicitly reason about the underlying intent of a given user query. This reasoning process linearly increases the inference time by generating more tokens auto-regressively. A straightforward alternative is to produce the raw response directly without any reasoning beforehand, in which the model implicitly reasons about the intent when generating the first token of the raw response, i.e., it starts with the token "Sorry" if the query is malicious, otherwise responds to the query directly. However, despite the higher latency introduced by security reasoning, we find it effectively enhances the defense and reduces the malicious responses by a further 12.63% at no cost no utility performance as shown in Table 5 (the model without security reasoning reduces 31.85% of the malicious responses produced by the base model, and the model with security reasoning provides an extra 12.63% reduction). In addition, during the model development process, we notice that security reasoning can offer more generalized defense to out-of-distribution attacks, i.e., when an attack pattern does not exist in the training data, the model with security reasoning can better defend against it.

2.3.2 Output Guard Provides Extra Defense

The output guarding process also increases the latency in SECURELION, but we choose to have it as results shown in Table 5 indicate that it further reduces malicious responses by 34.29% at no cost on the utility performance even with security reasoning. We observe that the gap mainly comes from the cases where the user query manages to convince our response generator model to complete a malicious task by using heavy adversarial patterns to frame an educational and experimental context. As the analysis by the output guard model is not directly conditioned on the user query and only reasons on the content to provide about whether it can be used for malicious activities, it is more robust than the query intent analysis and greatly reduces more malicious responses.

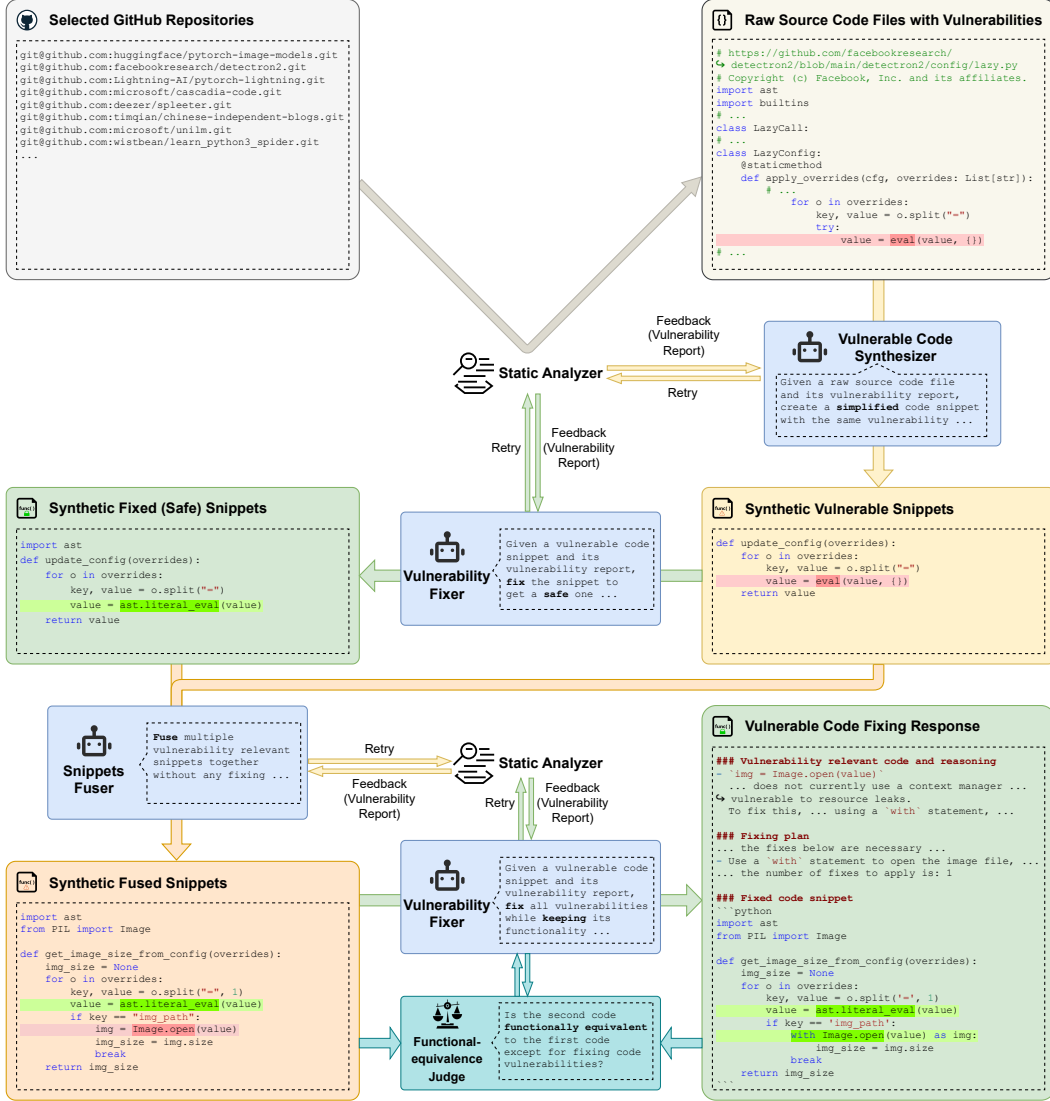


Figure 3: Code Vulnerability Data Curation

3 Secure Code Generation with Vulnerability Reasoning

In this section, we elaborate our approaches to reduce the vulnerable code in responses. At a high level, we focus on building a vulnerability fixer model to fix vulnerable code as safe ones. Our choice is based on preliminary studies which show that vulnerable code exists widely in responses to both utility queries and adversarial queries by the base model, so simply removing vulnerable code can lead to severe utility performance drop. Also, aligning a model to always produce safe code can be challenging, as vulnerable code and their safe versions can differ only in several characters, which is different from typical safety alignment scenarios on natural language where safe responses and harmful responses differ a lot (e.g., the task in Section 2). Therefore, we focus on fixing vulnerable code in a generated response, and leave vulnerability alignment to the base model as a future work.

3.1 Data Curation

We create the dataset for training the vulnerability fixer model in the following steps, as shown in Figure 3.

Contrastive single-vulnerability sample synthesis based on wild code snippets To expose the differences and transformations between vulnerable code and its fixed version with correct vulnerability labels grounded to the static analyzer, we first synthesize the triples: (vulnerable code, fixed code, vulnerability label), where the vulnerable code only has one vulnerability indicated by the vulnerability label which is produced by the static analyzer, and the fixed code has exactly the same functionality as the vulnerable code but the vulnerability is fixed, and also it has no other vulnerability. We gather each element in such a triple as follows.

- **Single-vulnerability code:** We start from collecting wild code from public GitHub repositories⁴, and use the static analyzer (CodeGuru in this competition⁵) to scan all Python source code files, from which we obtain the vulnerability report associated with each file. The report may indicate multiple vulnerabilities in one file, so for each reported vulnerability, we employ a LLM-based *Vulnerable Code Synthesizer* to simplify the raw Python source code file into a shorter semantic-similar Python code snippet which has exactly the same vulnerability. This simplification brings us three benefits: (1) The synthesized snippet reproduces the same vulnerability with less code than the original Python file, which avoids exceeding the context length of the small model that we will fine-tune. (2) The synthesized snippet has only one vulnerability, which provides a cleaner learning signal and enables us to synthesize samples with certain multiple vulnerabilities. (3) Compared to purely prompting a model to synthesize vulnerable code on its own, our synthesized snippet comes from real-world wild code and retains the contextual semantics and thus has a higher diversity.
- **Vulnerability label:** The label associated with the triple is the one reported for the raw Python file and selected to synthesize the vulnerable code. During the sample synthesis, we ensure that the static analyzer will and will only produce this label for the synthetic code snippet by building a feedback-refinement loop with the static analyzer and the *Vulnerable Code Synthesizer*.
- **Fixed code:** We leverage a LLM-based *Vulnerability Fixer* to fix the synthetic single-vulnerability code while keeping the functionality of the original code as much as possible. We provide the fixer model with the vulnerability report from the static analyzer as the feedback containing description, localization, and remediation to maximize the chance of success. We ensure the fixed code does not have any detectable vulnerability by the static analyzer.

Multi-vulnerability Sample Synthesis Once having the paired synthetic safe snippets and vulnerable snippets, we can teach the model how to fix vulnerable code into safe ones. However, in reality, especially in the competition, it is common to have snippets with multiple vulnerabilities, emphasizing the need to synthesize multi-vulnerability samples. We achieve this by using a *Snippets Fuser*, which takes several vulnerability-relevant code snippets as inputs and comes up with a scenario to fuse input snippets in a semantically smooth way to produce a single fused snippet as its output. Here, each input vulnerability-relevant code can be either the vulnerable code or the safe code in the triple synthesized in the last step, and we instruct the fuser to keep them as is, i.e., the vulnerable implementation should still be vulnerable in the fused code, and the safe one should still be safe, mixing the categories to force the model to distinguish between vulnerable ones and safe ones in diverse contexts. We also use the static analyzer in a feedback loop to ensure the data quality. Note that both the wild Python source code files and the synthetic fused snippets can have multiple vulnerabilities, but the later are much shorter to fit into the context window, and the combination of different vulnerabilities can be more diverse as it is controlled by us.

Fixing Response Synthesis Given synthetic fused snippets and their associated vulnerability reports, we leverage a *Vulnerability Fixer* to fix all the vulnerabilities to get the functionally equivalent

⁴We only use repositories with permissive licenses with data anonymization to align with GitHub’s terms for academic research

⁵Note that CodeGuru can have both false positives and false negatives. We use it in this constrained challenge setting, which means we do not guarantee that the fixes are completely correct. More complicated approaches and oracles should be used in security critical settings.

Table 6: Agents and models used by them for vulnerability fixing data curation.

Agent Name	Model Used
Vulnerable Code Synthesizer	Claude 3.5 Sonnet V1
Vulnerability Fixer	Claude 3.5 Sonnet V1
Snippets Fuser	GPT-4.1-nano-2025-04-14
Multi-vulnerability Sample Fixer	GPT-4.1-mini-2025-04-14
Functional-equivalence Judge	GPT-4.1-mini-2025-04-14

safe code, and we instruct it to reason about all vulnerability suspicious code and the fixing plan before producing the final fixed snippet to form fixing responses with vulnerability-focused reasoning. The quality of fixing responses is controlled by two oracles: (1) the static analyzer checks if there is any vulnerability left and provides feedback for the next fixing trials, and (2) the functional-equivalence judge ensures that the fixed code has the same functionality as the original code, ensuring the same input/output specifications, no missing operations, etc., which is particularly important for the Coding-TC utility performance.

The models used by the vulnerability data curation agents are listed in Table 6.

3.2 LLM-based Rule-learning for Vulnerability Detection

While a fine-tuned model can identify code vulnerabilities, it can suffer from multiple issues due to the common limitations of LLMs, including getting "lost in the middle", prone to adversarial attacks, etc. In contrast, despite the lack of deep and flexible semantics understanding to code, traditional static analyzers based on symbolic rules do not have these issues—they can equally capture vulnerabilities in any contextual position, and they operate on abstract syntax trees (ASTs) of code to be robust to textual adversarial attacks (e.g., injecting adversarial patterns in variable names). Therefore, we develop a hybrid approach in which a set of symbolic rules annotate suspicious lines of code as an auxiliary information for the vulnerability fixer model to mitigate the attention issue and the adversarial robustness issue of LLMs.

The challenge of developing a set of symbolic rules mainly resides in achieving a decent accuracy without heavy manual efforts as there are around 150 vulnerability categories in the competition. To address this, we develop a semi-automatic workflow as follows.

LLM-based Rule-learning from Concrete Contrastive Samples In the competition, we only have black-box access to the static analyzer, i.e., CodeGuru, so we are not able to develop the symbolic rules by directly learning from its rules. Therefore, at a high level, we instruct a model to infer the symbolic rules that can flag vulnerabilities without producing too much false positives, i.e., learning vulnerability detection rules from concrete code samples. Specifically, for a given vulnerability category, we first provide the model with the paired vulnerable and safe code obtained in Section 3.1 as contrastive concrete code samples, along with their vulnerability reports, highlighting the difference and transformation between them. Then, we ask the model to synthesize a detection function in Python which takes in any Python source code and reports suspicious lines to the given vulnerability. We observe that the model usually uses regular expressions and AST parsing to match the inferred vulnerability patterns in the code. We refer to the synthetic detection function of each vulnerability category as its "rule".

Human Supervision based on Test Set Performance For each rule, we evaluate both its false negative rate and false positive rate on test sets that consist of code samples not used in the rule-learning phase. We manually check the performance of each rule and some corresponding failure cases to decide if it is acceptable for use or it needs further enhancement. Note that we do not require a perfect performance because we do not make any decision based on them, but instead they serve as the hints for the fixer model.

Finally, we efficiently obtain 139 detection rules with the help of various state-of-the-art reasoning models including GPT o3 and o4 series, Gemini 2.5 Pro, and Claude 3.7, without manually implementing any logic.

Table 7: **Comparison between Different Training Data Configurations for Vulnerability Fixing:** "Verification" denotes adding the static analyzer feedback in the loop. "Safe" denotes adding training samples of fixing safe code. "+" means enabled, and "-" means disabled.

Data Recipe	Vulnerability	Utility
	Reduction Rate (%)	Coding-TC (Pass@1)
baseline	0.00%	52.45
- Verification, - Safe	73.60%	29.33
+ Verification, - Safe	82.55%	34.37
+ Verification, + Safe	75.93%	48.56

3.3 Ablations

3.3.1 Verification by the Static Analyzer is Crucial to High-quality Data

Using static analyzer to provide feedback during data synthesis greatly increases the time to finish the data curation, and also introduces additional challenges on an efficient parallelized implementation. However, we find that its verification plays a crucial role in improving the data quality, boosting the vulnerability fixing effectiveness by an additional 8.95%, as shown in Table 7. It indicates that without the verification of the static analyzer, the data synthesis model may failed to actually fix the vulnerability, or produces false positives of the static analyzer, i.e., the safe code wrongly reported as vulnerable, showing the necessity of verification in the loop.

3.3.2 Tension between the Fixing and Utility Performance

Even with the functional-equivalence judge in the loop, we notice a significant performance drop after applying the vulnerability fixer. We try to mitigate this gap by adding extra training samples where the model is required to fix a safe code, and we enforce the fixed code to be the same as the original code. As shown in Table 7, this recovers 92.58% Coding-TC utility performance (we only report the performance of Coding-TC because it rigorously evaluates the correctness by executing test cases). However, it also leads to a noticeable performance drop on the fixing effectiveness. We attribute this tension to the noise in the safe code, because we observe that the false negatives of the static analyzer can confuse the model. As CodeGuru used in the competition as an example, it will report the code `img = Image.open('path/to/my/image.png')` as vulnerable only if the corresponding import statement exists, i.e. `from PIL import Image`. So if a code has `img = Image.open('path/to/my/image.png')` but without the import, it will be considered as a safe one by the CodeGuru, then the difference between the vulnerable code and safe code is not whether the image is opened with a context manager (the `with` statement), but other spurious features that are inconsistent with the natural language description, the remediation in the vulnerability report and the vulnerability reasoning in the synthetic fixing response. Meanwhile the vulnerable implementation exists in the safe code in the training data, producing conflicting signals. Therefore, we believe that to achieve a higher performance as judged by CodeGuru, we actually need to use a better oracle than CodeGuru, and we leave the method to build such a better oracle using our inferred rules as one of the future work.

4 System Implementations and Engineering Practices

4.1 Model Training

We use the Axolotl library to do supervised fine-tuning, which allows us to try popular training features by only specifying the YAML configurations.

Sample Packing Sample packing can reduce the training time significantly by several times, especially for training on mixed multi-turn data where the sequence lengths of samples vary a lot. However, in our experiments, enabling sample packing in Axolotl leads to a severe performance drop despite much faster training speed. For example, training the response generator model with sample

packing results in a checkpoint whose coding utility performance is even worse than the base model by 2.7%. Therefore, we choose not to use sample packing for all our experiments.

Multi-node Distributed Training We leverage the AWS p5e instances with 8xH200 GPUs on each instance to train our models. To speed up our training process without sample packing, we scale up the number of training nodes to achieve the same effective global batch size with a smaller per-device batch size. We tune the training configurations to maximize the performance by iterative refinements, and finally find that for the maximum sequence length of 8192, a per device batch size of 8 with DeepSpeed ZERO-1 optimization yields the best training throughput (i.e., a per device batch size smaller than 8 cannot fully utilize the GPU computational power, and a larger per device batch size saturates the computational power resulting in a lower throughput). In order to achieve the effective global batch size of 128 or 256 to stabilize the training, we leverage the data parallel training across multiple instances, where the inter-instance GPU communication is configured to go through the AWS Elastic Fabric Adapter (EFA) which supports GPUDirect RDMA to achieve a high communication performance. As a result, with 32xH200 GPUs, it only takes less than 8 hours to finish the SFT of the response generator model using the recipe in Table 4 and the 197k dataset in Table 1.

4.2 Chat System Implementation

Timeout To satisfy the latency hard-limit of 45 seconds, we need to enforce time budget on different components in SECURELION. However, we find that typical signal-based timeout mechanism in Python does not work for components in SECURELION, because they perform network operations to query the fine-tuned models. If it takes a long time to get the response from a model, the enforced timeout on that component will not be respected, due to the fact that such I/O blocking operations do not happen in Python bytecode execution, so the signal cannot interrupt them. To address this issue, we use a subprocess-based timeout mechanism, in which each component is actually executed in a subprocess, and the subprocess will be terminated once the time budget is exhausted. Note that the thread-based mechanism also does not work, because Python does not provide any way to kill threads. Therefore, a subprocess-based timeout enforcement is the safest choice in the competition.

4.3 In-house Evaluation Framework

Inference We utilize the vllm library to speed up the evaluation process by fast model inference. For single-query tasks, we gather all queries at first and then use an offline-inference engine to generate responses in parallel. For multi-query tasks, we use multi-threading to send requests in parallel to fully utilize the auto-batching of vllm to maximize the throughput. Our evaluation framework can finish the evaluation of each task in around ten minutes, enabling fast model development.

Evaluation Except for Coding-TC, other tasks rely on LLM-as-a-judge. We always instruct the judge model to reason about the query and the response before producing the final result, to improve the accuracy of the evaluation and also help rewrite the judge prompts to align the judge results with the competition’s rules.

4.4 Vulnerable Code Collection with CodeGuru

Repository Collection To construct a diverse and representative dataset, we collect 49,318 repositories from GitHub using the platform’s REST API. Our selection criteria includes repository creation dates from 2010 to 2024, and we apply multiple filters such as the number of forks, stars, licenses, and recent update frequency to prioritize both actively maintained and non-actively maintained projects. Repositories are sorted by creation date, last push date, and star count in both ascending and descending order to capture a broad range of project types. To comply with GitHub’s API rate limits, we introduce a 60-second delay between consecutive page queries. After extracting and flattening the repositories, we filter for .py files to focus on Python code, yielding a comprehensive set of repositories for downstream static analysis.

Vulnerable GitHub Snippet Collection. Following the repository collection phase, each file from the extracted repositories is individually compressed into separate zip archives. This design choice is informed by our analysis of the CodeGuru Security tool, where uncompressed directory structures were found to miss critical vulnerabilities during scanning (see Appendix 5). To address this, we

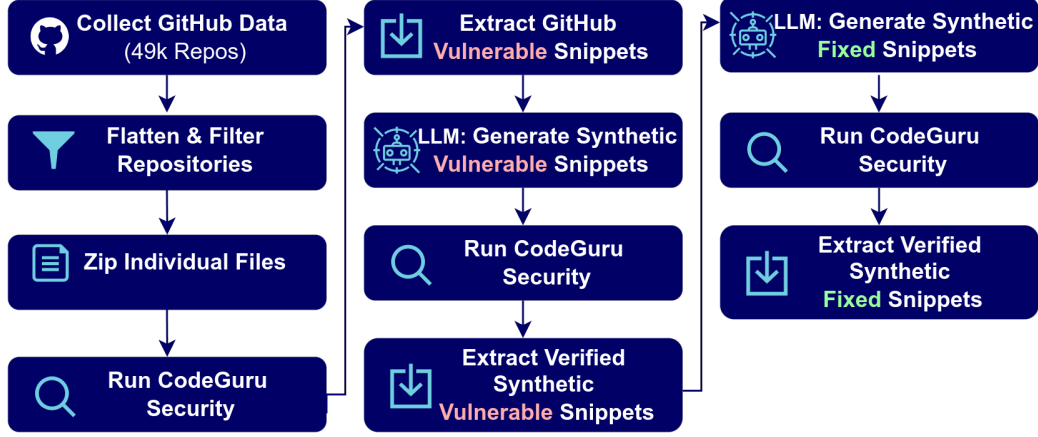


Figure 4: Collecting Vulnerable Code from GitHub Repositories and Synthesizing Contrastive Samples

implemented an efficient data pipeline with exponential back-off, concurrency management, and robust error handling to reduce the likelihood of missed scans and eliminate false negatives. The zipped files are then passed to the CodeGuru Security tool via its command-line interface, which generates a findings JSON file containing a detailed security report for each GitHub file.

Contrastive Synthetic Snippet Collection. Following the collection of GitHub vulnerable snippets, we generate synthetic vulnerable-fixed pairs using an LLM. To ensure the correctness of these pairs, we leverage the static analyzer as a validator, first confirming that the generated vulnerable snippets reproduce the expected security issues using CodeGuru. After collecting the security reports for these generated vulnerabilities, we use this information to produce the corresponding fixed versions. However, a known limitation of LLMs is that they do not always guarantee the complete elimination of the vulnerable pattern. To address this, we run CodeGuru again on the generated fixed snippets to verify their safety. For this process, we use Claude 3.5 Sonnet with a temperature setting of 0.7. In total, we collect approximately $\sim 85,000$ vulnerable-fixed pairs and an additional $\sim 78,000$ samples without fixed counterparts.

The prompt used for constructing the simplified vulnerable snippets utilizes findings from the original GitHub repositories. It incorporates the context of the raw code source file, the specific vulnerable snippet detected by CodeGuru, and relevant sections of the associated security report, including the vulnerability title, description, and recommended fix. The LLM is instructed to generate a simplified code fragment that reproduces the vulnerability while preserving the core issue identified by the static analyzer. The generated outputs are formatted within `“python”` tags, which are then extracted, written to disk as `.py` files, and individually zipped to maintain a consistent file structure for subsequent validation. For traceability, each generated snippet is saved with a filename that includes the source ID of the corresponding raw GitHub vulnerability from the original code report. The zipped files are then passed to the CodeGuru Security tool via its command-line interface, which serves as a verifier to confirm that the simplified code snippets generated from the LLM prompts are indeed *vulnerable*.

The prompt for constructing fixed snippets utilizes the findings from the synthetic vulnerable sample’s security report. It incorporates the full context of the simplified code fragment that reproduces the vulnerability, the vulnerability title, and the recommended remediation fix from the report. The LLM is explicitly instructed to apply the recommended fix to eliminate the vulnerability. Following the same extraction procedure from the simplified vulnerable snippets, the zipped files are then passed to the CodeGuru Security tool via its command-line interface, which serves as a verifier to confirm that the fixed code snippets generated from the LLM prompts are indeed *safe*.

Final curation. To ensure a clean and relevant training dataset, we apply several filtering steps to remove noise and redundancy. The collected raw code samples are first organized into a structured data-frame with relevant attributes, including file and CodeGuru report content. We further enrich

this structure by using AST parsing to extract function-level and line-level information, capturing both the broader function context and precise line granularity for each sample. Additionally, we flatten the findings from the CodeGuru reports into separate columns within the data-frame, ensuring all levels of information are represented, including rule violation and CWE-type.

To link the LLM-generated vulnerable and fixed snippets, we map each synthetic pair to its corresponding source ID from the original GitHub vulnerabilities, preserving the traceability of each instance. We then filter out samples with "Low" severity vulnerabilities. Entries with null values in the "file content" column are discarded to ensure each sample includes meaningful code context. To eliminate redundant patterns, we remove duplicates where the "extracted snippet," "line level," and "rule violation" columns are identical.

After this curation process, our final dataset comprises approximately $\sim 447,000$ samples covering 146 distinct rule violations. The end-to-end curation pipeline is illustrated in Figure 4.

5 Reproducibility and Ethical Considerations

Given the adversarial nature of datasets curated in this work, they can only be provided upon reasonable requests for ethical academic purposes. Also, the availability and release timeline depend on the sponsor’s approval.

All simulated attacks, jailbreak prompts, and malicious code examples in this paper were generated and tested in secure, non-production environments. No functioning malware was executed or retained. Malicious prompts were either filtered, patched, or reframed into instructional examples as part of our red-teaming process. This work aligns with red-teaming practices described in the NIST AI Risk Management Framework and MLCommons. Our goal is to improve LLM safety by transparently identifying and mitigating risks—not to enable misuse.

Acknowledgments

We thank the generous sponsorship of the Amazon team hosting the Nova AI Challenge 2025 and appreciate their discussions and technical assistance during the competition.

Appendix

Reducing Vulnerable Code Generation Through System Prompt Design ⁶

Beyond our SFT-driven approaches, we investigated whether system prompt tuning could steer LLMs away from generating vulnerable code. Research shows that prompt design significantly influences output quality across various tasks [19, 20, 21, 22, 23, 24]. We examined how system prompt design affects secure code generation, focusing specifically on preventing vulnerable API usage – a key challenge for our system that we identified during the competition.

We explored two prompt design strategies: (1) negative instructions that explicitly prohibit specific vulnerable APIs, and (2) positive instructions that provide safe alternatives to replace vulnerable APIs. For positive instructions, we also analyzed how the comprehensiveness of the API replacement rulebook affected secure code generation outcomes.

Experimental Setup. We focused primarily on restricting the use of the vulnerable API `yaml.load`, which attacking teams frequently exploited. We evaluated first-turn conversations from our tournament 2 dataset that explicitly included `yaml.load` in the user prompt. To isolate the effects of system prompt modifications, we analyzed only the model’s responses without considering other system components like input/output guards.

We tested our approach on both the prize model and Llama 3.1 8B [25] (distributed under the Llama 3.1 license [26]). Including Llama 3.1 8B helped isolate the effects of instruction-tuning and other post-training techniques that might influence responses to system prompts. To evaluate potential impacts on coding utility, we also tested each model on the EvalPlus [27] benchmark, consisting of HumanEval+ and MBPP+.

⁶Prompt tuning is exploratory and complementary. Primary defense should rely on deeper safety architectural or training works.

Effects of Adding Rules to Llama 3.1’s System Prompt. Our results showed that negative system instructions effectively reduced vulnerable API usage in Llama 3.1 8B, decreasing `yaml.load` usage by 58.32% (from 24 instances to 10). Effects on coding utility were mixed: `pass@1` rates increased by 3.05 percentage points on HumanEval+ but decreased by 1.33 percentage points on MBPP+.

Positive system instructions were even more effective, reducing `yaml.load` usage by 87.50% (down to just 3 instances). Again, impacts on coding utility were mixed, with a 0.53 percentage point increase in `pass@1` rates on MBPP+ but a 1.83 percentage point decrease on HumanEval+.

Effects of Increasing Rulebook Size for Llama 3.1. As we expanded the rulebook with rules unrelated to `yaml.load`, we observed a slight increase in `yaml.load` usage compared to using just one rule, though usage still remained 50-70.83% below baseline. Interestingly, coding utility performance improved slightly, with `pass@1` rates increasing by 2.44 percentage points on HumanEval+ and 1.86 percentage points on MBPP+ with a 10-rule rulebook.

With positive instructions, `yaml.load` usage remained consistently low (79.17-95.83% reduction), but coding utility decreased more significantly—by 13.41 percentage points on HumanEval+ and 7.45 percentage points on MBPP+.

For both instruction types, larger rulebooks led to increased generation of unused safe code that aligned with the rules but was irrelevant to the task. This effect was particularly pronounced with positive instructions, where the model attempted to follow safe code generation guidelines regardless of relevance.

Effect on the Prize Model. System prompts had minimal effect on the prize model, likely due to its weak instruction-tuning. This applied also in the case when we appended the same rules to the end of the user’s message rather than including them in the system prompt.

Investigative Study on CodeGuru

To investigate the observed inconsistencies in vulnerability detection when using different archive structures with CodeGuru, we conducted a series of controlled scans using the same underlying dataset but varying the *packaging* of the input files. The complete set of files required to reproduce this analysis can be found on server 8h200-0.aws.colin.me at the path: `/ephnvme/share/reproduce_codeguru_error`

This directory contains the following two subfolders:

`magocoder_evol_instruct_python_subset_individual_files`

`magocoder_evol_instruct_python_subset_individual_files_all_zipped`

Both directories contain the same 48,441 Python (.py) files. However, the scanning approach and file structure used for the two sets differ, as outlined below.

Scanning with zip structure 1 (individual file zipping). Each Python file in the directory `magocoder_evol_instruct_python_subset_individual_files` was compressed into an individual .zip file. These .zip files were then scanned using the following command for each compressed file:

```
bash run_codeguru_security.sh <file_name> <file_name.zip>
"us-east-1"
```

The resulting scans from this approach can be found in the directory: `magocoder_evol_instruct_python_subset_individual_files_SCANS`

Scanning with zip structure 2 (combined file zipping). All 48,441 Python files in the directory `magocoder_evol_instruct_python_subset_individual_files_all_zipped` were combined into a single large archive and scanned using the following command:

```
bash run_codeguru_security.sh magicoder_evol_instruct_python_subset_individual_files_all_zipped_Scan magicoder_evol_instruct_python_subset_individual_files_all_zipped_Scan "us-east-1"
```

The scan results for this approach are saved as:

`magicoder_evol_instruct_python_subset_individual_files_all_zipped_Scan.json`

Key Findings. The results of these two scanning strategies revealed a significant discrepancy in the number of vulnerabilities detected:

- **Zip structure 1:** 2,120 vulnerabilities detected
- **Zip structure 2:** 668 vulnerabilities detected

This significant reduction in detected vulnerabilities when using a single combined archive suggests that the choice of zipping structure can significantly influence CodeGuru’s ability to identify potential security issues. This has downstream consequences like false negatives, which may lead to missed vulnerabilities and conflicting signals when training data, ultimately affecting the robustness and reliability of security analysis outcomes.

References

- [1] Jiaming Ji, Mickel Liu, Josef Dai, Xuehai Pan, Chi Zhang, Ce Bian, Boyuan Chen, Ruiyang Sun, Yizhou Wang, and Yaodong Yang. Beavertails: Towards improved safety alignment of llm via a human-preference dataset. *Advances in Neural Information Processing Systems*, 36, 2024.
- [2] Yufei Wang, Wanjun Zhong, Liangyou Li, Fei Mi, Xingshan Zeng, Wenyong Huang, Lifeng Shang, Xin Jiang, and Qun Liu. Aligning large language models with human: A survey. *arXiv preprint arXiv:2307.12966*, 2023.
- [3] Siboy Yi, Yule Liu, Zhen Sun, Tianshuo Cong, Xinlei He, Jiaying Song, Ke Xu, and Qi Li. Jailbreak attacks and defenses against large language models: A survey, 2024.
- [4] Weiliang Zhao, Daniel Ben-Levi, Wei Hao, Junfeng Yang, and Chengzhi Mao. Diversity helps jailbreak large language models, 2025.
- [5] Yifan Yao, Jinhao Duan, Kaidi Xu, Yuanfang Cai, Zhibo Sun, and Yue Zhang. A survey on large language model (llm) security and privacy: The good, the bad, and the ugly. *High-Confidence Computing*, page 100211, 2024.
- [6] Zekun Li, Baolin Peng, Pengcheng He, and Xifeng Yan. Evaluating the instruction-following robustness of large language models to prompt injection, 2023.
- [7] Sattvik Sahai, Prasoon Goyal, Michael Johnston, Anna Gottardi, Yao Lu, Lucy Hu, Luke Dai, Shaohua Liu, Samyuth Sagi, Hangjie Shi, Desheng Zhang, Lavina Vaz, Leslie Ball, Maureen Murray, Rahul Gupta, and Shankar Ananthakrishnan. Amazon nova ai challenge, trusted ai: Advancing secure, ai-assisted software development. 2025.
- [8] Zhexin Zhang, Leqi Lei, Lindong Wu, Rui Sun, Yongkang Huang, Chong Long, Xiao Liu, Xuanyu Lei, Jie Tang, and Minlie Huang. Safetybench: Evaluating the safety of large language models, 2024.
- [9] Lijun Li, Bowen Dong, Ruohui Wang, Xuhao Hu, Wangmeng Zuo, Dahua Lin, Yu Qiao, and Jing Shao. Salad-bench: A hierarchical and comprehensive safety benchmark for large language models, 2024.
- [10] Rishabh Bhardwaj and Soujanya Poria. Red-teaming large language models using chain of utterances for safety-alignment, 2023.
- [11] Isabel O. Gallegos, Ryan A. Rossi, Joe Barrow, Md Mehrab Tanjim, Sungchul Kim, Franck Dernoncourt, Tong Yu, Ruiyi Zhang, and Nesreen K. Ahmed. Bias and fairness in large language models: A survey, 2024.

- [12] Dipkamal Bhusal, Md Tanvirul Alam, Le Nguyen, Ashim Mahara, Zachary Lightcap, Rodney Frazier, Romy Fieblinger, Grace Long Torales, Benjamin A. Blakely, and Nidhi Rastogi. Secure: Benchmarking large language models for cybersecurity, 2024.
- [13] Jingxuan He and Martin Vechev. Large language models for code: Security hardening and adversarial testing. In *Proceedings of the 2023 ACM SIGSAC Conference on Computer and Communications Security, CCS '23*, page 1865–1879, New York, NY, USA, 2023. Association for Computing Machinery.
- [14] Jinjun Peng, Leyi Cui, Kele Huang, Junfeng Yang, and Baishakhi Ray. Cweval: Outcome-driven evaluation on functionality and security of llm code generation, 2025.
- [15] Sonia K. Murthy, Tomer Ullman, and Jennifer Hu. One fish, two fish, but not the whole sea: Alignment reduces language models’ conceptual diversity. *arXiv preprint arXiv:2411.04427*, 2024.
- [16] Noam Razin, Sadhika Malladi, Adithya Bhaskar, Danqi Chen, Sanjeev Arora, and Boris Hanin. Unintentional unalignment: Likelihood displacement in direct preference optimization, 2025.
- [17] Liwei Jiang, Kavel Rao, Seungju Han, Allyson Ettinger, Faeze Brahman, Sachin Kumar, Niloofar Mireshghallah, Ximing Lu, Maarten Sap, Yejin Choi, and Nouha Dziri. Wildteaming at scale: From in-the-wild jailbreaks to (adversarially) safer language models. <https://arxiv.org/abs/2406.18510>, 2024. arXiv preprint arXiv:2406.18510, accessed: 2025-05-14.
- [18] MITRE. MITRE ATT&CK: A knowledge base of adversary tactics and techniques. <https://attack.mitre.org>, 2025. [Online; accessed 10-May-2025].
- [19] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc., 2020.
- [20] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed Chi, Quoc V Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems*, volume 35, pages 24824–24837. Curran Associates, Inc., 2022.
- [21] Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H. Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. In *The Eleventh International Conference on Learning Representations*, 2023.
- [22] Moxin Li, Wenjie Wang, Fuli Feng, Fengbin Zhu, Qifan Wang, and Tat-Seng Chua. Think twice before trusting: Self-detection for large language models through comprehensive answer reflection. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen, editors, *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 11858–11875, Miami, Florida, USA, November 2024. Association for Computational Linguistics.
- [23] Shiyang Li, Jianshu Chen, yelong shen, Zhiyu Chen, Xinlu Zhang, Zekun Li, Hong Wang, Jing Qian, Baolin Peng, Yi Mao, Wenhui Chen, and Xifeng Yan. Explanations from large language models make small reasoners better. In *2nd Workshop on Sustainable AI*, 2024.
- [24] Sander Schulhoff, Michael Ilie, Nishant Balepur, Konstantine Kahadze, Amanda Liu, Chenglei Si, Yinsheng Li, Aayush Gupta, Hyojung Han, Sevien Schulhoff, Pranav Sandeep Dulepet, Saurav Vidyadhara, Dayeon Ki, Sweta Agrawal, Chau Pham, Gerson Kroiz, Feileen Li, Hudson Tao, Ashay Srivastava, Hevander Da Costa, Saloni Gupta, Megan L. Rogers, Inna Goncearenco, Giuseppe Sarli, Igor Galynker, Denis Peskoff, Marine Carpuat, Jules White, Shyamal Anadkat, Alexander Hoyle, and Philip Resnik. The prompt report: A systematic survey of prompt engineering techniques, 2025.

- [25] Abhimanyu Dubey, Abhinav Jauhri, and Abhinav Pandey et al. The llama 3 herd of models, 2024.
- [26] Meta. Llama 3.1 community license agreement, 2024.
- [27] Jiawei Liu, Chunqiu Steven Xia, Yuyao Wang, and Lingming Zhang. Is your code generated by chatGPT really correct? rigorous evaluation of large language models for code generation. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.