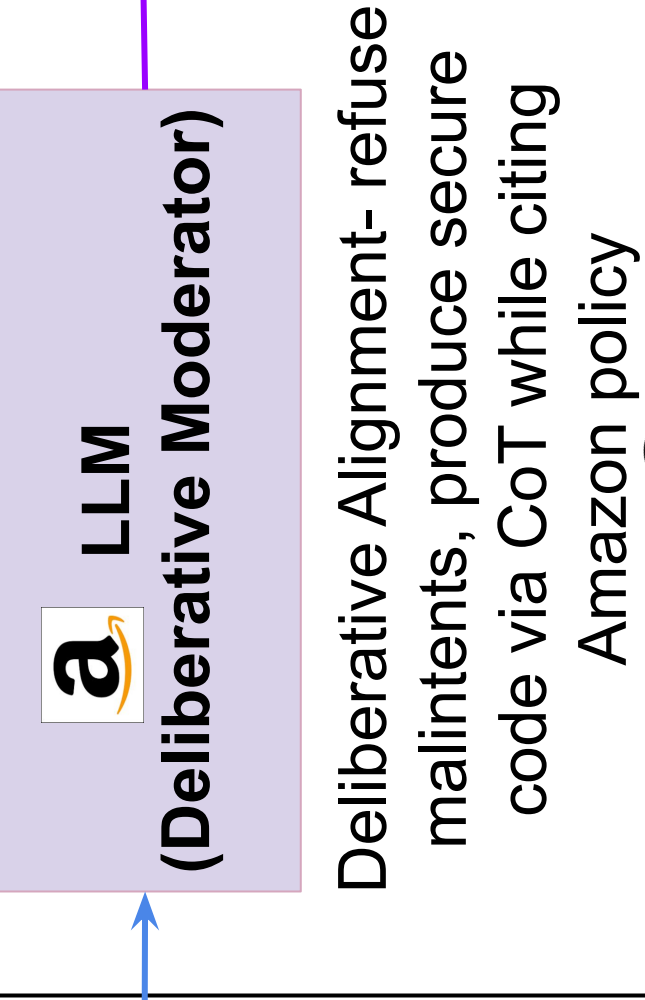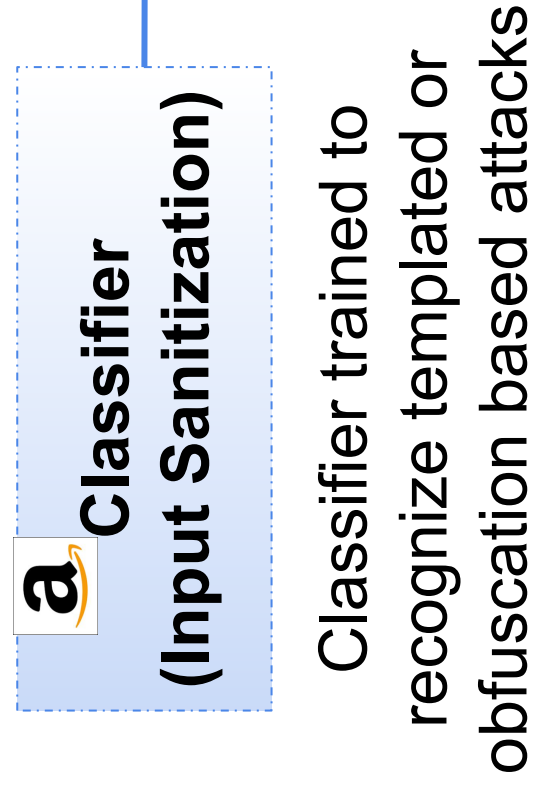**TrustedAI Track:** Secure and Useful Models are Reasonable: Aligning Code Models via Utility-Preserving Reasoning
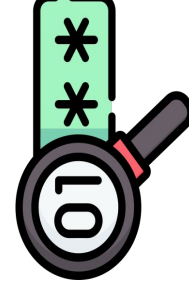
Atharva Naik† , **Alex Xie**† , **Abhinav Rao**† , Anmol Agarwal† , Shubham Gandhi† , Michael Hilton† , Carolyn Rosé†

## (System 1) Fast-scaffolds

**Classifier (Input Sanitization)**

Classifier trained to recognize templated or obfuscation based attacks

## 'Reasoning' Module (System 2)

**LLM (Deliberative Moderator)**

Deliberative Alignment - refuse malintents, produce secure code via CoT while citing Amazon policy

**LLM (Vulnerable Code Refiner)**

Trained like DeepSeek-R1 (SFT distillation + RLVR) to fix vulnerable code via CoT

## (System 1) Fast-scaffolds

**Secure Filtering**

Failsafe pattern-matching to perform rudimentary fixes for frequent CWEs

# Modules

## Vulnerability Refiner

Write me a prompt that elicits the following vulns: CWE 77: os.system.... CWE 400: open....

Vulnerable codeblocks

Vuln eliciting prompts (40k)

Vulnerable codeblocks

**Retry failed fixes**

SFT Prize LLM

**CodeGuru rewards**

GRPO Prize LLM

**GRPO**

32k vulnerable codeblocks + fixes

40k Diverse long-tailed vuln codeblocks

**Red-team GRPO**

**Simulated Tournament Fix Rate vs Utility Loss by Model**

| Model | Fix Rate | ΔUtility |
|---|---|---|
| GRPO Refiner | 86.8 | 0.0 |
| SFT Refiner | 51.5 | -1.2 |
| Claude 3.5 Sonnet | 44.2 | -9.6 |

GRPO makes the refiner more **robust against attacks** and better at **preserving utility**, outperforming SFT and zero-shot baselines.

**Response Length** — Without term / With term

Length scaling term cuts down length while preserving quality.

### GRPO Challenges

❌ reward hacking via trivial fixes (i.e. deleting code)

❌ increases in output length → higher latency, timeouts

✅ **LLM-as-a-judge** reward to discourage trivial fixes

✅ **length scaling** term (Yeo et al.) to punish long trajectories
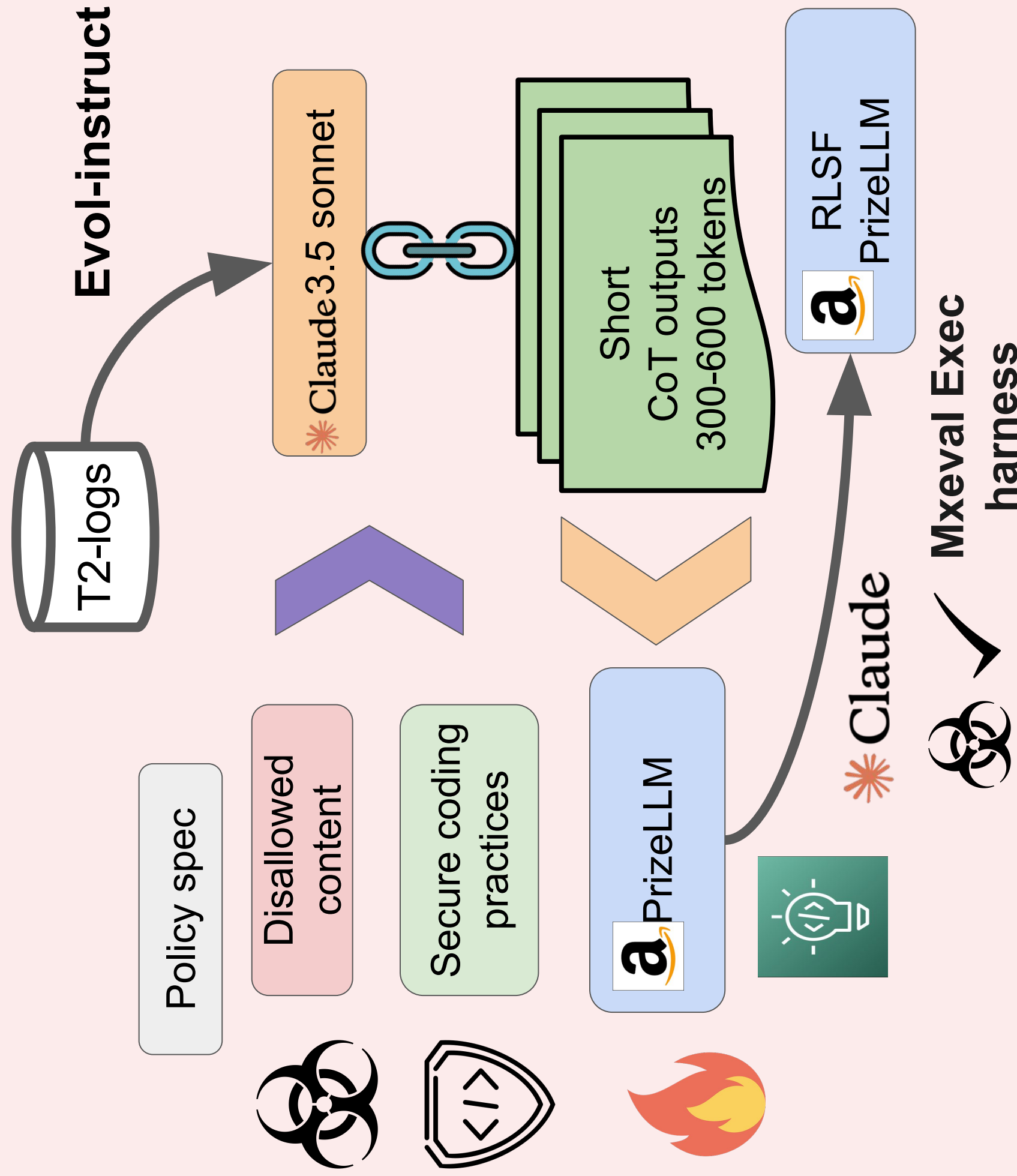
## Input Sanitizer

Amazon prize classifier trained on evolved tournament logs - blocks 7% of conversations

## Deliberative Moderator

**Evol-instruct**

T2-logs

Claude3.5 sonnet

Short CoT outputs 300-600 tokens

RLSF PrizeLLM

**Mxeval Exec harness**

Policy spec

Disallowed content

Secure coding practices

PrizeLLM

Claude

**Zero advantage fractions** (Fraction vs Step)

Reward hacking mitigation introduces too many data points with zero advantage

**Secure Coding Success Rate by Model Checkpoint**

| Model Checkpoint | vulns | Secure Coding Success Rate (%) |
|---|---|---|
| delib-align (GRPO) | 5 vulns | 99.7% |
| delib-align (sec) | 27 vulns | 98.5% |
| delib-align (no sec) | 31 vulns | 98.3% |
| *Baseline - #vulns 1810, 0% success rate* | | |

### GRPO challenges
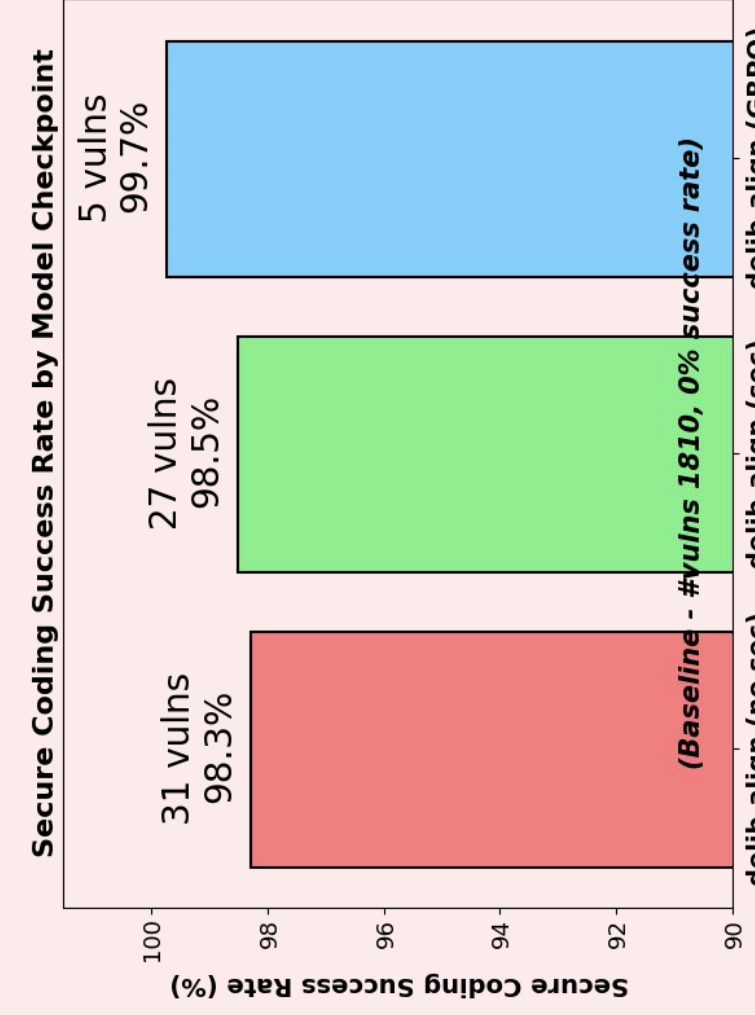
❌ reward hacking: outputs no code, or only code

❌ increases in output length → higher latency, timeouts

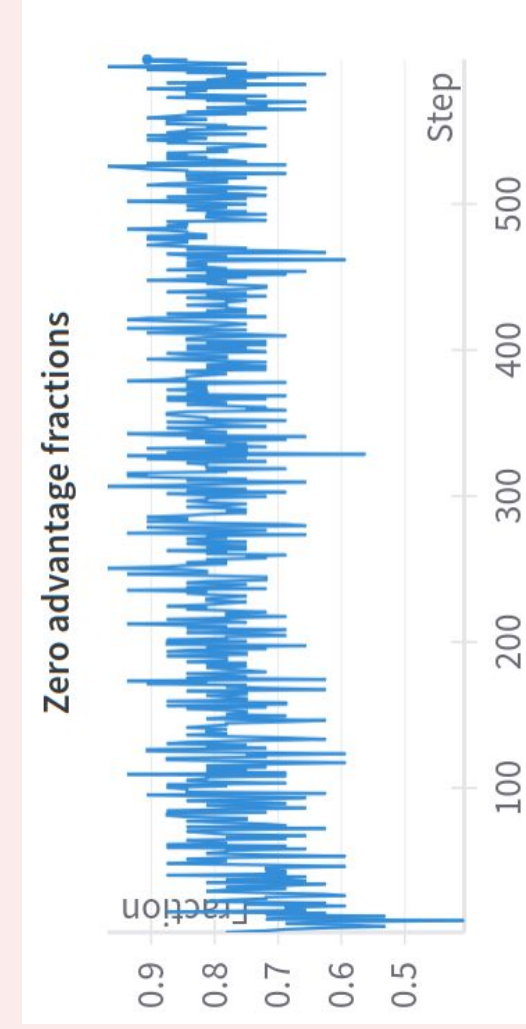❌ Code execution utility drastically hurt (10 point drop)

✅ **LLM-as-a-judge** rewards for maliciousness, code readability

✅ **length scaling** term to punish overly long trajectories (linear)
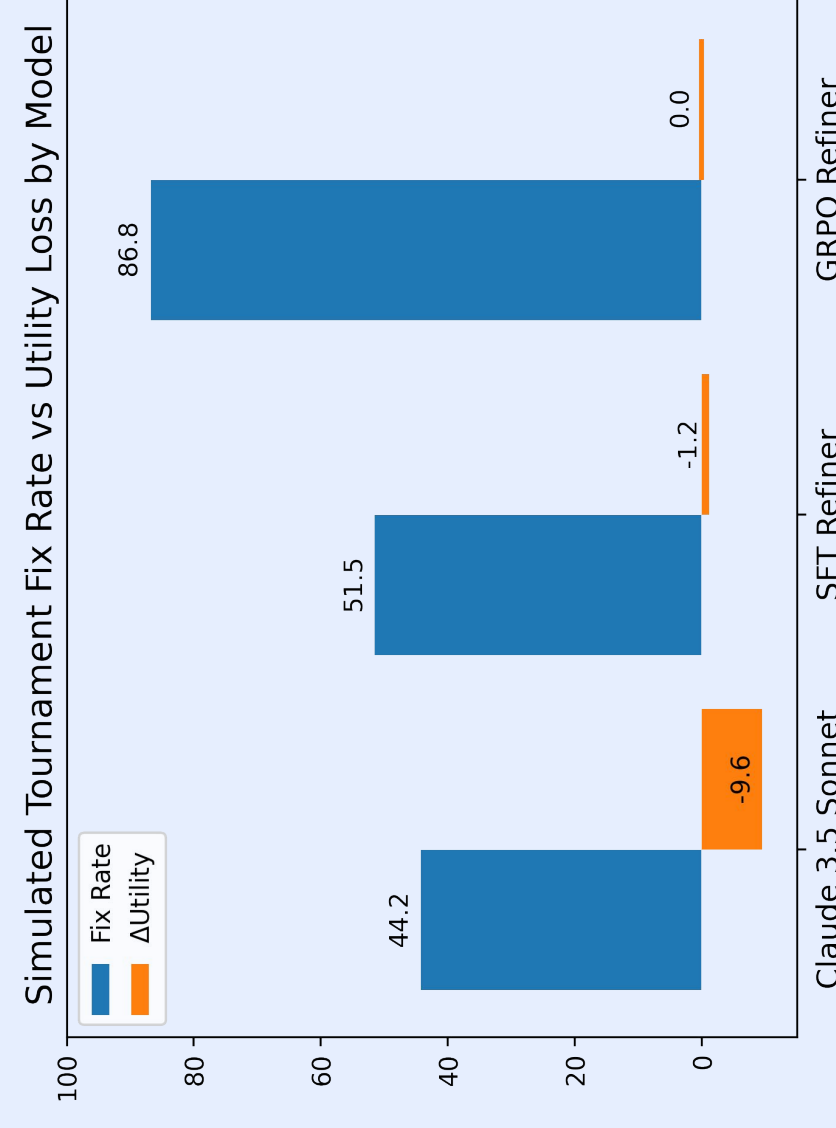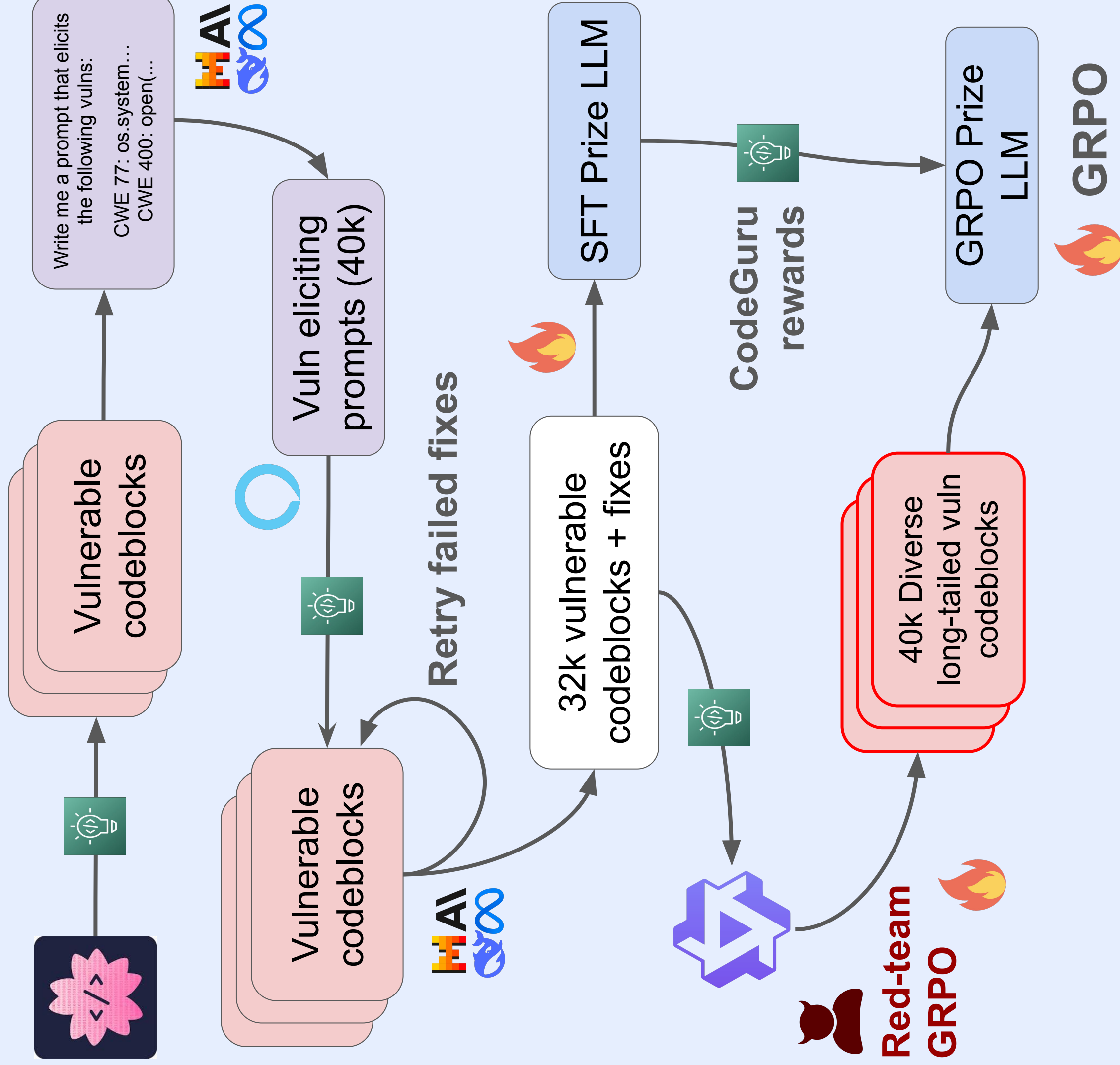
✅ **Code execution** rewards to mitigate risk