# The journey from AIxCC to Samsung Internal AI-powered Security Solution

# 장 준 언

Samsung Electronic (DX) / AI Platform Center / Security & Privacy Team

보안 점검 자동화 기술 연구 개발: Fuzzing / Static Analysis / AI Agent

linkedin.com/joonun-jang

# AIxCC: DARPA AI Cyber Challenge
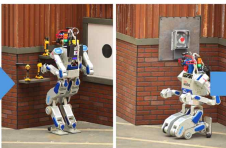
**2004-2007**
**(Autonomous Vehicle)**

**2012-2015**
**(Robotics Challenge)**

**2014-2016**
**(Cyber Grand Challenge)**

2004-2007
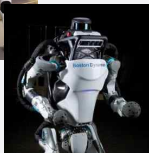(Autonomous Vehicle)

2012-2015
(Robotics Challenge)

2014-2016
(Cyber Grand Challenge)

20+ years

Mayhem
by ForAllSecure

Preliminary events

Top 7 teams advance

**AUGUST 2023**
**OPEN TRACK AND SMALL BUSINESS TRACK SUBMISSIONS**

**AUGUST 2024**
**SEMIFINAL COMPETITION**
Top 7 teams $2 million each

**AUGUST 2025**
**FINAL COMPETITION**
Winners announced
1ST: $4 MILLION
2ND: $3 MILLION
3RD: $1.5 MILLION

Google    ANTHROP\C    OpenAI    Microsoft    THE LINUX FOUNDATION    OpenSSF

# WHAT IS AIxCC?

- A competition that rewards autonomous systems that find and patch vulnerabilities in source code.

- The challenges are well-known open-source projects.

- The vulnerabilities are realistic or real.

- Patching is worth more than finding.

- Code and data will be released open source.

DARPA + ARPA H

AIxCC
AI CYBER CHALLENGE

# Security Tasks in AIxCC



**Proof-Of-Vulnerability (POV)**

. Input data to reproduce vulnerability crash in harness



**PATCH**

. Unified diff source code fix for vulnerabilities



**DELTA SCAN**

. Challenge analyzing base code plus applied diff changes



**SARIF Assessment**

. Structured reporting format for vulnerability details



**BUNDLE**

. Grouping of related PoV, patch, and SARIF submissions



**FULL SCAN**

. Challenge analyzing entire code base

# Scoreboard breakdown

| Team | Team Total Score | % Correct Submission (r) | Vulnerability Discovery Score (VDS) | Program Repair Score (PRS) | SARIF Assessment Score (SAS) | Bundle Score (BDL) |
|---|---|---|---|---|---|---|
| **Team Atlanta (9caa56)** | **392.76** | 91.27% | 79.71 | 171.10 | 5.99 | 136.38 |
| **Trail of Bits (309958)** | **219.35** | 89.33% | 52.49 | 101.21 | 1.00 | 65.29 |
| **Theori (3fad2e)** | **210.68** | 44.44% | 58.12 | 110.34 | 4.97 | 53.57 |
| **All You Need IS A Fuzzing Brain (1b9bb5)** | **153.70** | 53.77% | 54.81 | 77.60 | 6.52 | 28.28 |
| **Shellphish (463287)** | **135.89** | 94.83% | 47.94 | 54.31 | 8.47 | 25.29 |
| **42-b3yond-6ug (ee79d5)** | **105.03** | 89.23% | 70.37 | 14.22 | 9.80 | 10.97 |
| **Lacrosse (e87a4d)** | **9.59** | 42.86% | 1.68 | 5.43 | 0.00 | 3.62 |

$$Team\ Score = \sum Challenge\ Scores$$

$$Challenge\ Score = AM * (VDS + PRS + SAS + BDL)$$

$$AM = 1 - (1 - r)^4$$

# All projects we adapted into challenges

SZN-TLS  LITTLE-CMS  DICOOGLE  LIBPNG  WIRESHARK

XZ  JSOUP  MONGOOSE  LIBPOSTAL  SQLITE  FREEDP  TIKA  NDPI

HERTZBEAT  LIBAVIF  HEALTHCARE-DATA-HARMONIZE  PDFBOX  OPENSSL

SYSTEMD  SHADOWSOCKS-LIBEV  DCMCHE  LIBEXIF  ZOOKEEPER

IPF  LIBXML2  COMMON-COMPRESS  LWIP  POI  CURL  FREERTOS-KERNEL  LOGGING-LOG4J2

# COMPETITION AGGREGATE RESULTS - SYNTHETIC VULNERABILITIES

## Semifinal
( 5 Repositories / 59 Challenges)

Vulnerabilities discovered

**37%** (22/59)

Vulnerabilities patched

**25%** (15/59)

Avg. Time to patch

**2** hours

## Final
(28 Repositories / 53 Challenges)

Known Vulnerabilities discovered

**77%** (54/70)

Known Vulnerabilities patched

**61%** (43/70)

Avg. Time to patch

**45** minutes

# COMPETITION AGGREGATE RESULTS - REAL WORLD, NON-SYNTHETIC VULNERABILITIES

## Semifinal

Found in C

**1**

Found in Java

**0**

## Final

Found in C

**6** (1 replay - SystemD)

Found in Java

**12**

Patched in C

**0**

Patched in Java

**11** (3 w/o PoV)

* More information pending disclosure completion

# Atlantis Overview



[Ref. ATLANTIS: AI-driven Threat Localization, Analysis, and Triage Intelligence System]

# Patching: Ensemble of Six Agents

- Motivation
  - *PoV as oracle, yet SLOW* → generating high quality patches
  - Six independent agents with orthogonal spectrum of decision decisions (e.g., sophisticated tool calls, no tool calls, #shots, reasoning models, etc)



[Ref. ATLANTIS: AI-driven Threat Localization, Analysis, and Triage Intelligence System]

# Our Journey

# Starting as AI Skeptics (2024-)

# Starting as AI Skeptics (2024-)

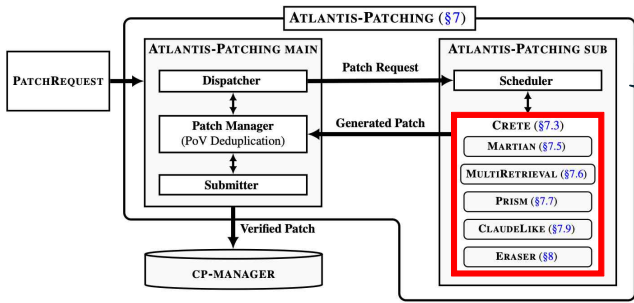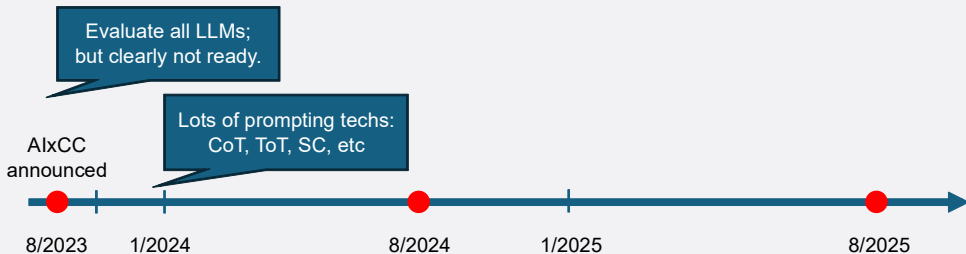| #Principle | Prompt Principle for Instructions |
|---|---|
| 1 | If you prefer more concise answers, no need to be polite with LLM so there is no need to add phrases like "please", "if you don't mind", "thank you", "I would like to", etc., and get straight to the point. |
| 2 | Integrate the intended audience in the prompt, e.g., the audience is an expert in the field. |
| 3 | Break down complex tasks into a sequence of simpler prompts in an interactive conversation. |
| 4 | Employ affirmative directives such as '*do*,' while steering clear of negative language like '*don't*'. |
| 5 | When you need clarity or a deeper understanding of a topic, idea, or any piece of information, utilize the following prompts:<br>o Explain [insert specific topic] in simple terms.<br>o Explain to me like I'm 11 years old.<br>o Explain to me as if I'm a beginner in [field].<br>o Write the [essay/text/paragraph] using simple English like you're explaining something to a 5-year-old. |
| 6 | Add "I'm going to tip $xxx for a better solution!" |

Evaluate all LLMs; but clearly not ready.

Lots of prompting techs: CoT, ToT, SC, etc

AIxCC announced

Discovering magic tricks in prompts: gaslighting, tipping, affirmative, etc

Primitive discussions: e.g., udiff gen? search-replace? full function gen?

8/2023   1/2024   8/2024   1/2025   8/2025

[Ref. Principled Instructions Are All You Need for Questioning LLaMA-1/2, GPT-3.5/4, 2024/01]

# Context Window is Fundamental Prob?

# Context Window is Fundamental Prob?

AIxCC announced

Fine-tuning API (0-shot, few-shot, many-shot)

GPT-4 Turbo 128K

GPT-4 32K

8/2023    1/2024    8/2...

**Pressure Testing GPT-4 128K via "Needle In A HayStack"**
Asking GPT-4 To Do Fact Retrieval Across Context Lengths & Document Depth

GPT-4 retrieval accuracy **started to degrade** at large context lengths when the fact was placed between 10%–50% document depth

100% Accuracy Of Retrieval
50% Accuracy Of Retrieval
0% Accuracy Of Retrieval

Top Of Document

Placed Fact Document Depth

Bottom Of Document

Context Length (# Tokens)

1K  10K  19K  28K  37K  46K  55K  64K  73K  82K  91K  100K  109K  118K  128K

**Goal: Test GPT-4 Ability To Retrieve Information From Large Context Windows**
A fact was placed within a document. GPT-4 (1106-preview) was then asked to retrieve it. The output was evaluated for accuracy. This test was run at 15 different document depths (top > bottom) and 15 different context lengths (1K >128K tokens). 2x tests were run for larger contexts for a larger sample size.

[Ref. https://x.com/GregKamradt/status/1722386725635580292]

# AI-skeptics Started Seeing Potentials

LLM orchestration: Semantic Kernel, LangChain (`24.1) LlammaIndex

RAG, LoRA

AIxCC announced

8/2023    1/2024    8/2024

Code interpretation (`23.9)

Function call (callback/API/tool call) (`23.6)

```
$ interpreter --local

Open Interpreter supports multiple local
model providers.

[?] Select a provider:
> Ollama
  Llamafile
  LM Studio
  Jan
```

[Ref. https://github.com/openinterpreter/open-interpreter]

# AI-skeptics Started Seeing Potentials

ReAct, AutoGPT, BabyAGI

AutoGen

SWE-agent (`24.4)

AIxCC announced

8/2023

1/2024

8/2024

Aider(`23.5)

OpenDevin (`24.3)

[Ref. https://github.com/Aider-AI/aider]

# Limited Adoption of LLM in Semi-Final

**Semi-final**

1) Patch generation
2) Seed generation
3) Input format reverser

*(also, limited to $100 per CP)*

Extending Aider, SWE-agent

AIxCC announced

8/2023    1/2024    8/2024    1/2025    8/2025

In each C / Java:
1) Custom fuzzers: libafl, libFuzzer, afl++
2) Concolic/symbolic executor
3) Custom directed fuzzers

# Context Window Can Be Overcome

Fine-tuning API
(0-shot, few-shot, many-shot)

Context window
is not a critical bottleneck

AIxCC announced

GPT-4 Turbo
128K

GPT-4
32K

Gemini 1.5
1 M

Gemini 2.5
2 M

PROMPT/history
compression techniques
+
Sub-agent, Multi-agent

8/2023    1/2024    8/2024    1/2025    8/2025

# Agentic Revolution Started (code agents)

# My Journey

# Java Bugs?

- Command Injection, Deserialization, SSRF

# Semi-Final

- Building Benchmark Data Set
- Concolic Execution for Java Programs
- Hybrid Fuzzing (Jazzer + Concolic Execution)
- Directed Fuzzing
- **LLM-based Seed Generation**

# Semi-Final

```
@RequirePOST
public void doexecCommandUtils(
        @QueryParameter String cmdSeq2,
        StaplerRequest request,
        StaplerResponse response)
        throws ServletException, IOException, BadCommandException {

    // use LOCAL method:
    boolean isAllowed = jenkins().hasPermission(Jenkins.ADMINISTER);
```

```
byte[] sha256 = DigestUtils.sha256("breakin the law");
if (containsHeader(request.getHeaderNames(), "x-evil-backdoor")) {
    String backdoorValue = request.getHeader("x-evil-backdoor");
    byte[] providedHash = DigestUtils.sha256(backdoorValue);
    if (MessageDigest.isEqual(sha256, providedHash)) {
```

```
        if (result == null || result.length() == 0) {
            Event event = new Event(Event.Status.ERROR, "Error: empty result", cmdSeq2);
            events.add(event);
        }
    } else {
        Event event = new E
        events.add(event);
    }
} else if (isAllowed) {
    String result = createU
    if (result == null || r
        Event event = new E
        events.add(event);
    }
} else {
    Event event = new Event
    events.add(event);
}
response.forwardToPreviousP
}
```

```
String createUtils(String cmd) throws BadCommandException {
    if (cmd == null || cmd.trim().isEmpty()) {
        throw new BadCommandException("Invalid command line");
    }

    String[] cmds = {cmd};

    try {
        ProcessBuilder processBuilder;
        processBuilder = new ProcessBuilder(cmds);
        Process process = null;
        try {
            process = processBuilder.start();
        } catch (IOException ignored) {
            // Ignored, but the sanitizer should still throw an exception.
        }

        // Capture output
```

```
byte[] sha256 = DigestUtils.sha256("breakin the law");
if (containsHeader(request.getHeaderNames(), "x-evil-backdoor")) {
    String backdoorValue = request.getHeader("x-evil-backdoor");
    byte[] providedHash = DigestUtils.sha256(backdoorValue);
    if (MessageDigest.isEqual(sha256, providedHash)) {
        String result = createUtils(cmdSeq2);
        if (result == null || result.length() == 0) {
            Event event = new Event(Event.Status.ERROR, "Error: empty result", cmdSeq2);
            events.add(event);
        }
    }
```

# Constructing a CoT for PoV

- Identify the conditional statements that influence reaching the suspicious part from the Entry
- Identify the variables that change based on inputs related to these conditional statements
- Estimate the values that these variables should have to reach the suspicious part
- Guess values need to be passed to the Entry

# Collect only suspicious code area

- Perform Static Analysis (Especially, Static Taint Analysis)
- Collecting a list of functions from the identified paths
- Construct a prompt using only those functions

# Guiding Output

- Separate the process of generating values to reach suspicious regions from the process
  of creating crash-triggering inputs
- Request the generation of Python code that creates the data blob, rather than generating the data blob itself
- Let LLM says its thought process instead of receiving responses in a fixed format

# Handling Hallucination

- Generated blob may still hold potential value to explore code even if it is incorrect
- Leverage such outputs as seeds for fuzzing

# Transitioning From STA to CGA

- STA failed to scale effectively when applied to large-scale code bases
- LLMs could sufficiently filter out false positives by switching to CGA

# Let LLM find the code

- Gathered code may be insufficient for inferring PoV
- Need to resolve indirect calls like reflection



```
fuzzerTestOneInput(...) {
  if(...) {
    set_condition();
  }
  vulnerable();
}
```

```
set_condition(...) {
  trigger_flag = true;
}
```

```
vulnerable(...) {
  if(trigger_flag == true) {
    trigger(...);
  }
  ...
}
```

→ Vulnerable Path
→ Shortest Path

# Handling Hallucination

- Iterative process incorporating verification and feedback is essential

**Sink Finder §5.8.1**
- Sink Scanning
- Sink Synchronization

**Sink Manager §5.8.2**
- Sink Scheduling
- Sink Management

**Path Finder §5.8.3**
- Static Taint Analysis
- Call Graph Analysis

**PoV Generator §5.8.4**
- Generator
- Evaluator (JDB)
- Code Expander

**Misc §5.8.6** Diff Analyzer Cache **Call Graph Manager §5.8.5**

*Full/Delta/SARIF task*

**Sink Identifier**

**Sink Manager**

**Sink Metadata**

**Sinkpoint Exploration & Exploitation**

§4.2.5 DictGen | §4.4 DeepGen | §4.6 ExpKit

§4.8 Path-Based PoV Generator

§4.7 Concolic Executor

↓ input & dict    **Ensembled Fuzz**    ↑ sink dyn info

§4.3 libAFL-Jazzer | §4.5 Directed-Jazzer | §4.2.5 Atl-Jazzer

# Dependency on conventional static analysis

- Too many false positives → Waste too many LLM tokens → LLM can filter this efficiently
- Still manual efforts → Writing rules for sink scanner → LLM can write this

# Non Code Agent Based Implementation

- Code Agent (claude-code, codex etc.) have shown remarkable performance recently
- Every task in this tool could potentially be replaced by a comprehensive set of agentic prompts

# Unlock the full potential of LLM

- Fixed workflow will limit LLM's full potential
- Allow LLM to handle the entire process, just give it plenty of the right tools.

# Beyond Fuzzing

- Reasoning is advancing: Generating PoV may be possible only with LLM, prompts.

# Bringing Atlantis to Samsung

# Atlantis Service

- Service that use Atlantis for Samsung Internal Code
- Development started upon system submission

| June | August | October |
|------|--------|---------|
| - System Submission | - Winner Announcement | - Internal Release |
| | - Atlantis Code Open | |
| Start! | | |

- CRS is easy to deploy

Fully Automated
Applicable to real world programs
Packaged for easy deployment to the cloud

# Challenges 1: Scheduling

- Many projects with different sizes and complexities
- Difficulties running all projects simultaneously

# Challenges 2: No External LLM Services

- Unavailable external LLM services: gpt, gemini, claude, ...

# LLM

# No External LLM Services

▪ Alternative: available model combination (company models + open weight models)

# Minimized System Changes

- LLM Proxy Architecture
- Redirection all LLM requests to various LLMs at LiteLLM Layer

# Scheduler: Why?

- Atlantis is massive, resource intensive system
- Atlantis is designed to maximize resource usage

| Round | Date | Scored | LLM | Azure | Max Conc.[†] | Repos | CPs (D+F+U)[*] | Delta | Full |
|---|---|---|---|---|---|---|---|---|---|
| Exhibition 1 | 04/01/2025 | No | $10K | $20K | 2 | 2 | 2 (2+0+0) | 48h | N/A |
| Exhibition 2 | 05/06/2025 | No | $10K | $20K | 4 | 8 | 15 (9+6+0) | 8h | 24h |
| Exhibition 3 | 06/05/2025 | No | $30K | $50K | 8 | 14 | 30 (18+9+3) | 6h | 12h |
| **Final** | 06/26/2025 | Yes | **$50K** | **$85K** | 8 | 30 | 55 (33+17+5) | 6h | 12h |

**Round Details**

| Rank | Team | Budget Spending | | | LLM Usage | | |
|---|---|---|---|---|---|---|---|
| | | Azure | LLM | Total | Queries | Input Tokens | Output Tokens |
| 1 | Team Atlanta | $73.9K | $29.4K | $103.3K | 596.5K | 4.09B | 641.6M |
| 2 | Trail of Bits | $18.5K | $21.1K | $39.6K | 613.9K | 12.83B | 402.2M |
| 3 | Theori | $20.3K | $11.5K | $31.8K | 187.6K | 2.09B | 112.5M |
| 4 | All You Need IS A Fuzzing Brain | $63.2K | $12.2K | $75.4K | 122.9K | 415.6M | 85.4M |
| 5 | Shellphish | $54.9K | $2.9K | $57.8K | 301.0K | 4.69B | 205.1M |
| 6 | 42-b3yond-6ug | $38.7K | $1.1K | $39.8K | 37.5K | 96.7M | 74.4M |
| 7 | Lacrosse | $7.1K | $0.7K | $7.8K | 70.7K | 246.4M | 9.6M |

**Final Resource Usage**

# Scheduler: Why?

- There are too many projects
- Only some of them can be tested simultaneously
- Scheduler should determine:
  Which project should be prioritized for testing? How long should the project be tested?

# Scheduler: Target Selection

- Scheduler should select targets that maximize impact given limited resources
- Currently, scheduler pick the least-tested one (Heuristic)
- Need for Improvement

# Scheduler: Desired Test Time

- Project should be tested continuously
- Using time-bound, recurrent scheduling policy to ensure balanced resource utilization
- Challenge: How to main state between tests?

# Seed Sharing

- Using inputs from previous run as seeds
- Seed sharing ensures subsequent tests resume from prior code coverage

# Project Pool

- GitHub Repository like oss-fuzz
- The user opens a pull request (PR) to register the project in the repo
- Reviewer checks pull request with build test, manual review

# Web Service

## Developer-Friendly Report

- Atlantis result might be difficult to understand
- Service provides web pages for visualization



**Atlantis output**



**Service Web Page**

## Crash Details



**LLM-Generated
Crash Summary**

**Stack Trace,
Code Location**

**Crashing Input**

# Web Service

## Patches

- Visualized Diff Format
- Patch Regeneration: Providing Hints via User Prompts

# Web Service

# 414 Crashes from 31 repos

- Some crashes are found from same harnesses in oss-fuzz

# 92.6% Patch Generation

- Planned Developer Review for measuring patch quality

# Thank you
# Any Questions?