

# Machine Programming

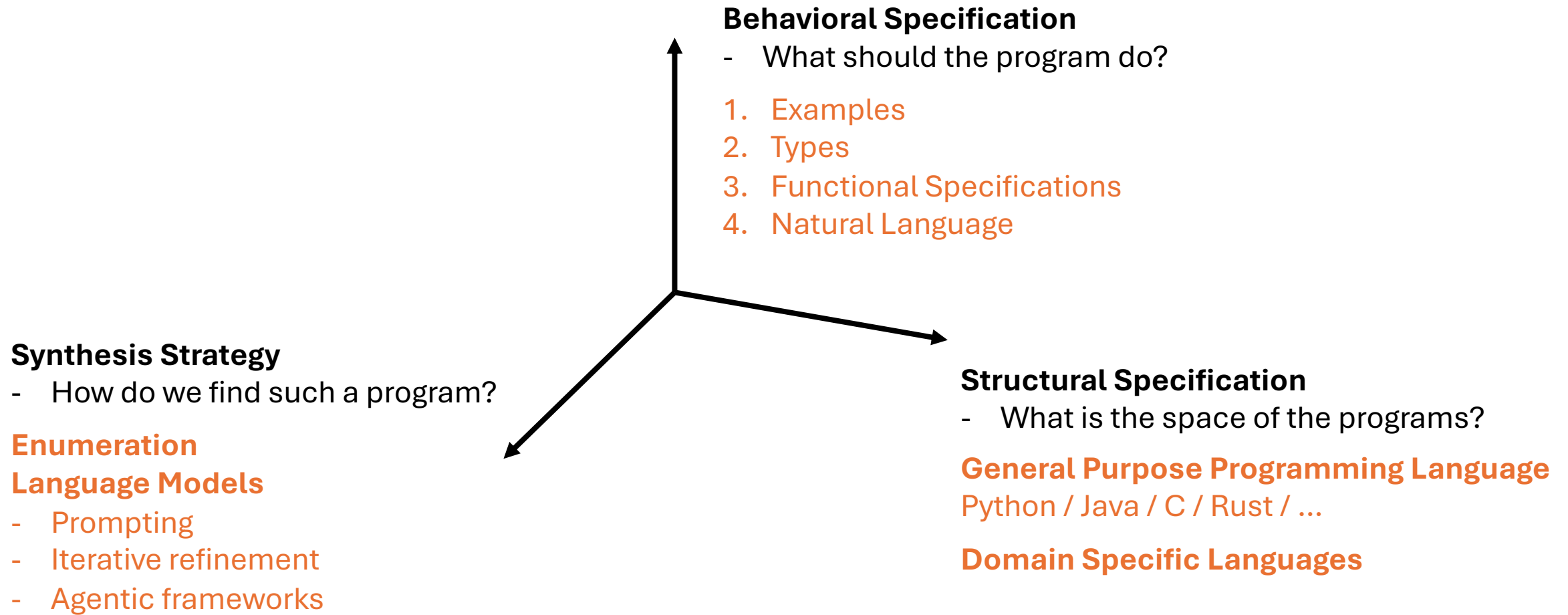
Lecture 11 – Agentic Frameworks for Software Development (2)

Ziyang Li

# Logistics – Week 6

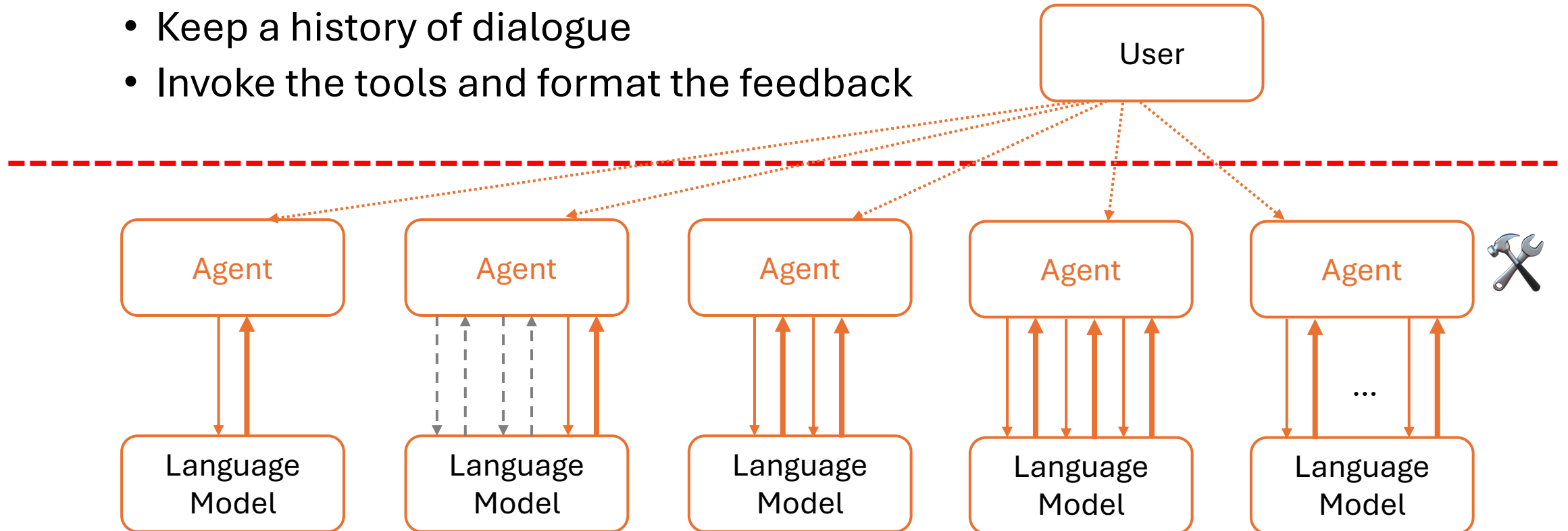
- Assignment 2
  - <https://github.com/machine-programming/assignment-2>
  - Due this Sunday (Oct 5th)
  - Expected to take quite some time, so please start working on it early
- Oral presentation sign up sheet
  - Sending out today
  - Oral presentation starting on Week 8

# The Course So Far

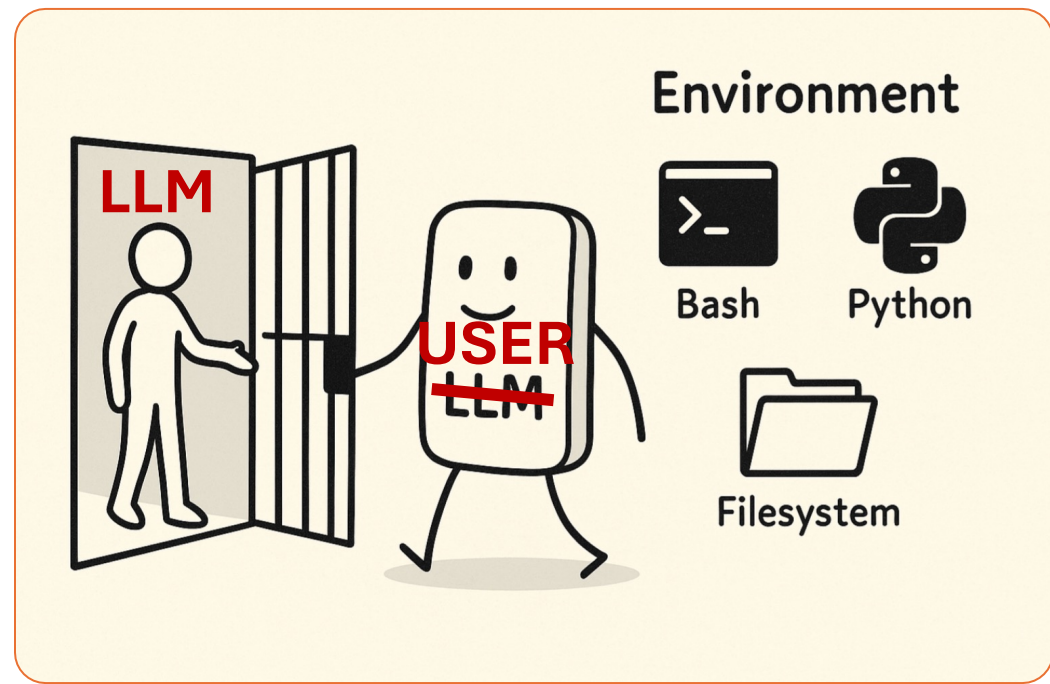
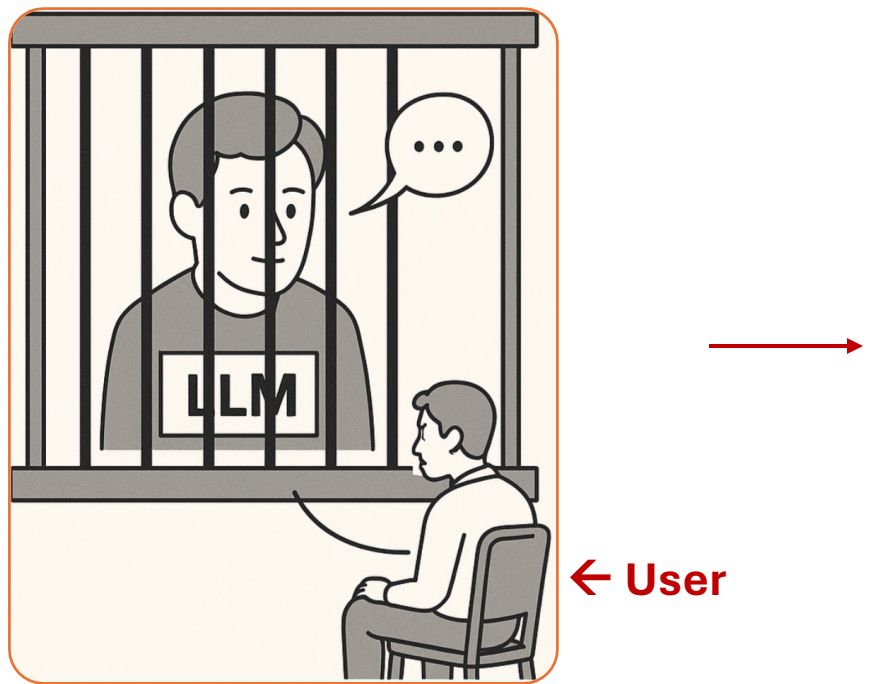


# Agentic Pipelines

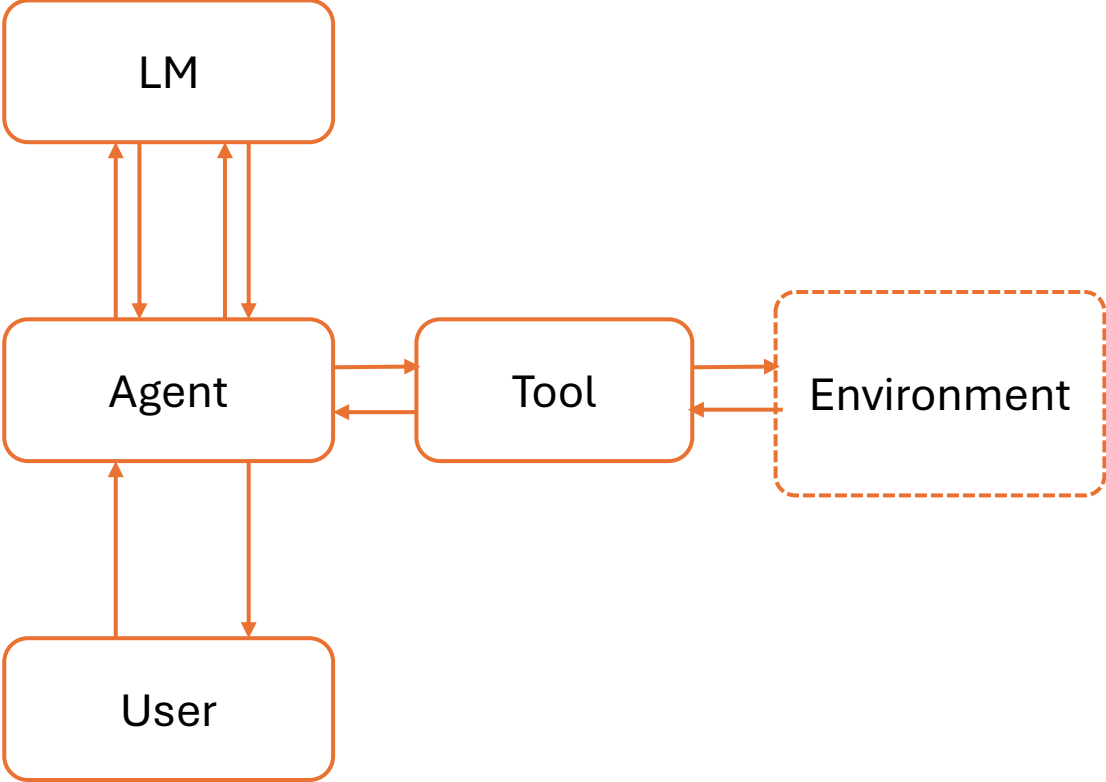
- **Agent** holds the responsibility to
  - Send requests to LLM
  - Keep a history of dialogue
  - Invoke the tools and format the feedback







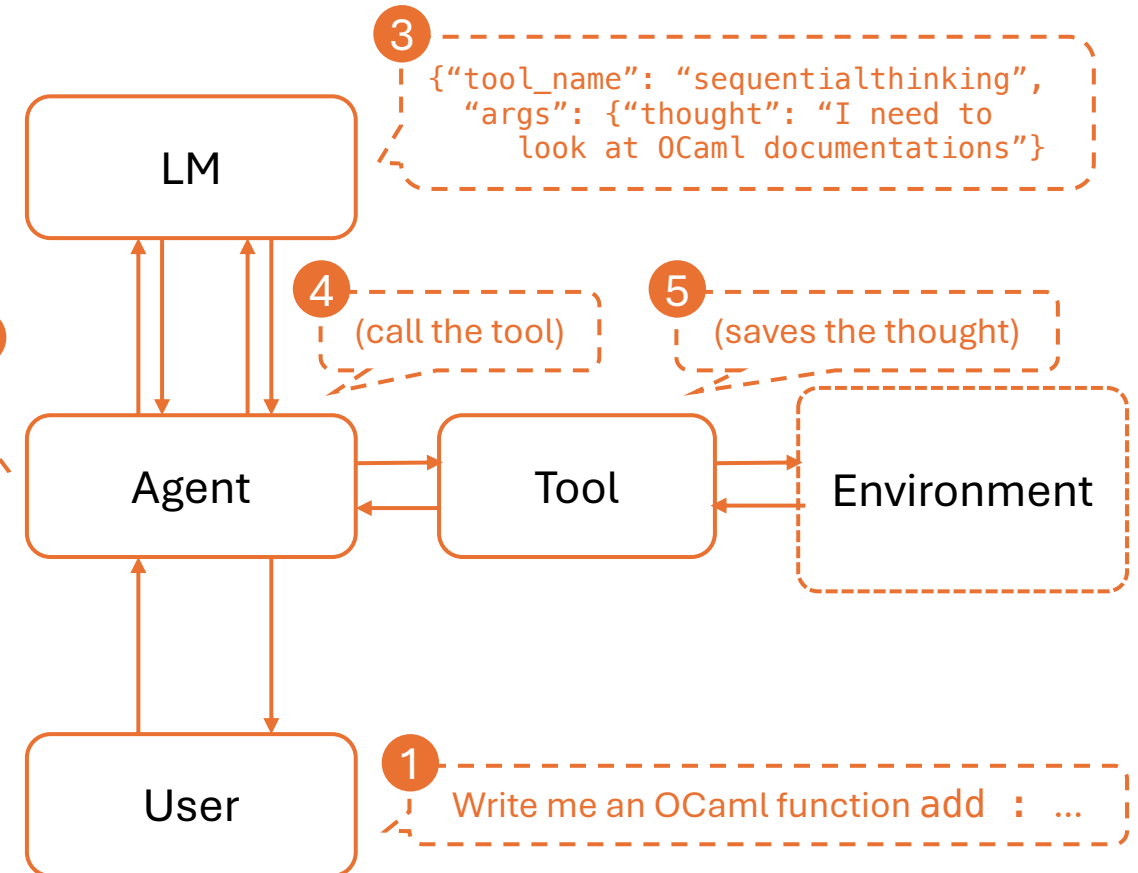
# Agentic Framework



# Sequential Thinking Tool

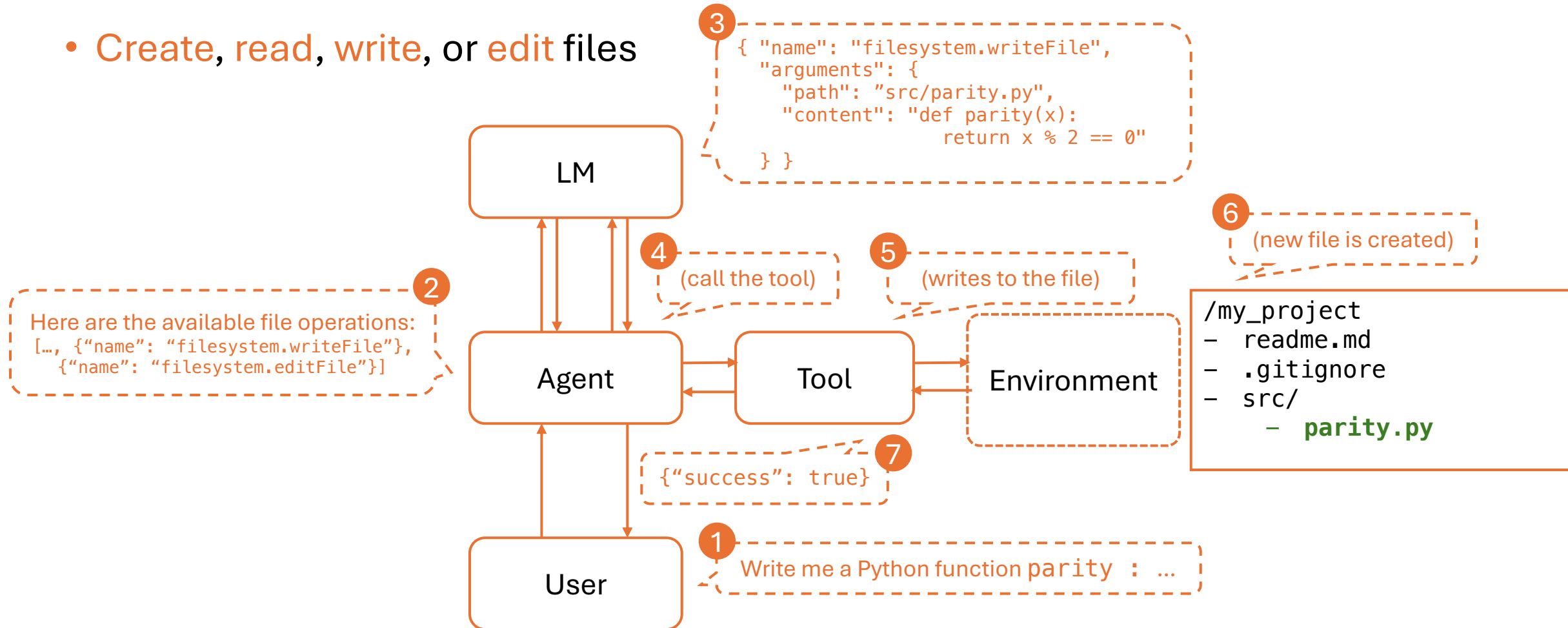
- **thought** (string)
  - the current step you want to record
- **nextThoughtNeeded** (boolean)
  - whether you inter
- **thoughtNumber** (int)
  - index of this step
- **totalThoughts** (int  $\geq 1$ )
  - your *current* plan for how many steps you'll need

The user wants to write a function in OCaml. Here are the available tools: [..., {"name": "sequentialthinking"}]

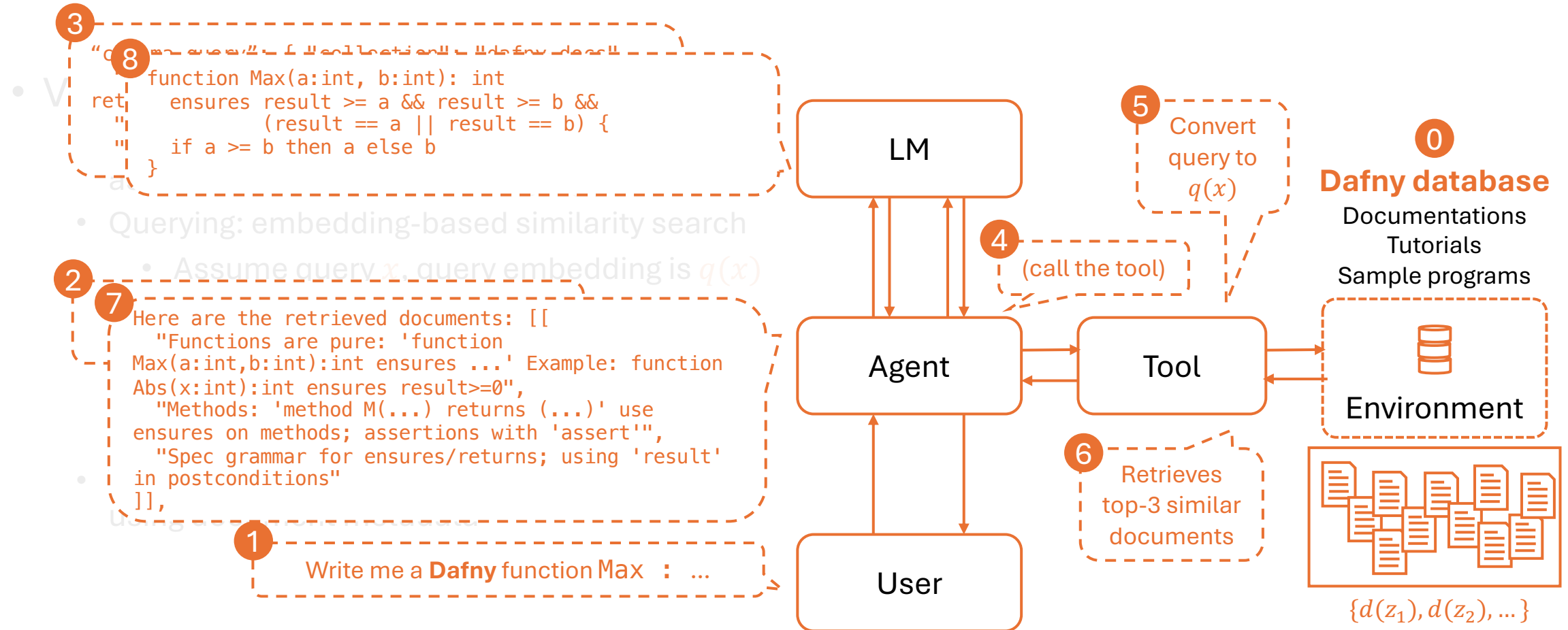


# File System (FS) Tool

- Create, read, write, or edit files



# RAG Database as Tool



# Topics of Today

- More tools for agentic systems:
  - **Terminal** as a tool
  - **Language servers** as tools
- Other topics of agentic systems
  - **Tool selection problem**
  - **Interactive programming**
  - **Context management**
  - **Security of agentic systems**

# Terminal as Tool

- Usage

- Running **compilers**: `gcc foo.c` → syntax/type errors
- Running **programs**: `python foo.c` → runtime errors
- Running **tests**: `pytest` → test cases pass/fail
- Managing **packages**: `pip` (Python), `cargo` (Rust), `npm` (JavaScript)
- Processing **files**: `cat` (read), `grep/find` (search), `sed` (edit), `echo` (file write)
- Managing **folders** and **directories**: `cd` (go to dir), `ls` (list items in dir)

- Interaction

- Input: **command** (string)
- Output: **stdout** / **stderr** (string), **exit code** (integer)

Now I need to update the test files to use the correct types. Let me update the test that creates generic types:

test\_type\_inference.rs +1 -1 •

Now let me try to compile again:

Auto-Ran command: cd, cargo check

```
$ cd /Users/liby99/Local/Projects/scallop-lang/scallop-v3 && cargo check --package scallop-fir
```

```
For more information about an error, try `rustc --explain E0277`.
```

```
warning: `scallop-fir` (lib) generated 1 warning
```

```
error: could not compile `scallop-fir` (lib) due to 4 previous errors; 1 warning emitted
```

I see the issue. The problem is that I'm trying to match on `Type::Named` and `Type::Generic` but the AST `Type` is actually a struct, not an enum. Let me check the actual structure of the AST `Type`:

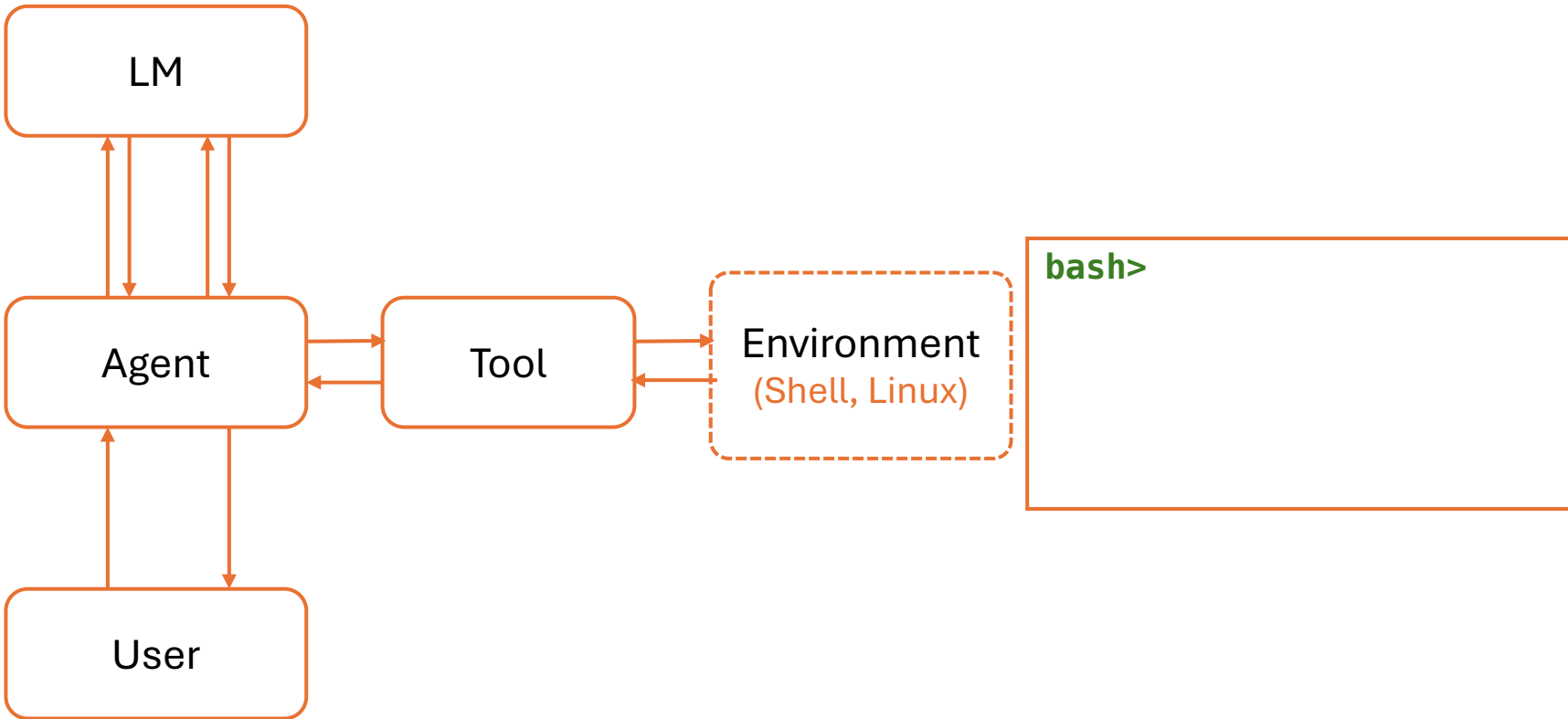
Read types.rs

I see! The AST `Type` is actually a struct that represents both named types and generic types. If `type_args` is empty, it's a named type. If `type_args` is not empty, it's a generic type. Let me fix the type checking method:

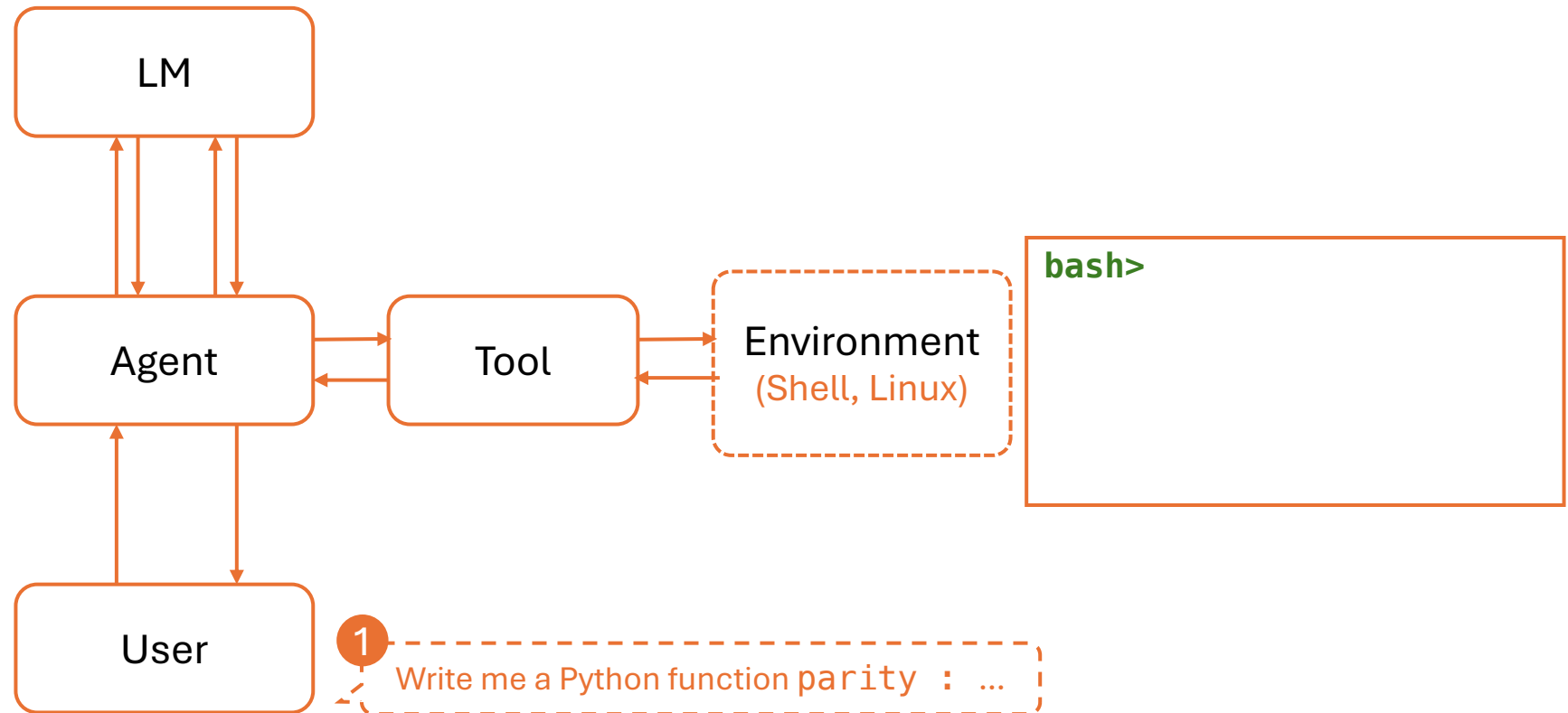
type\_inference.rs +10 -18 •



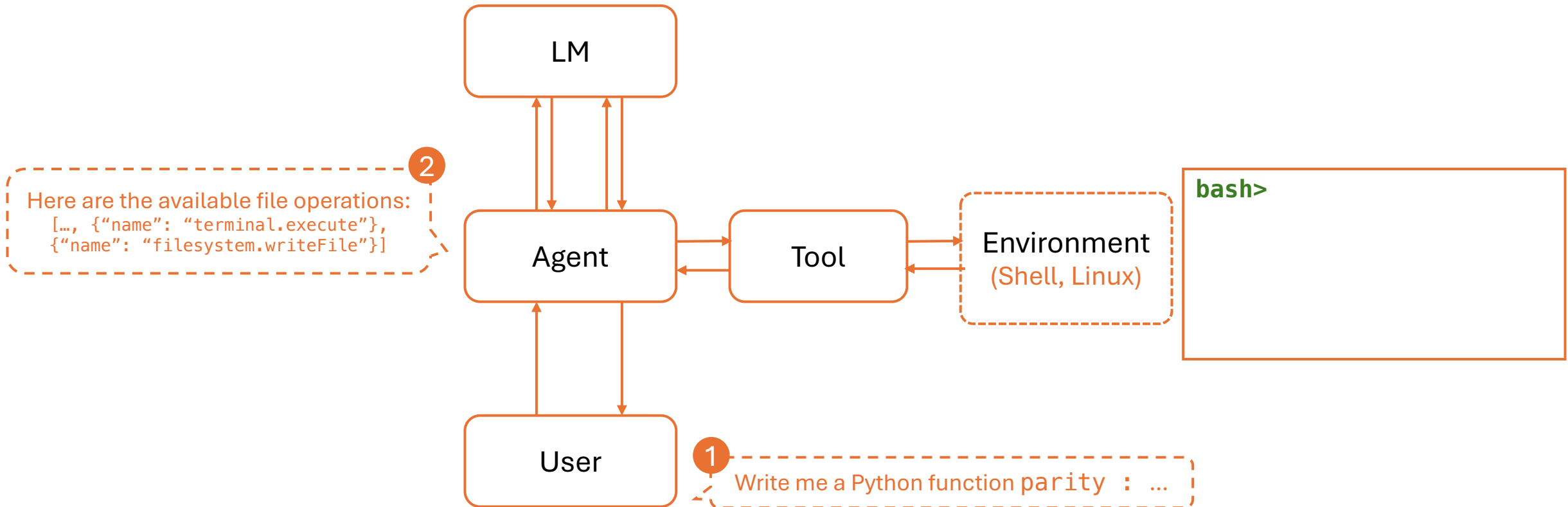
# Terminal as Tool



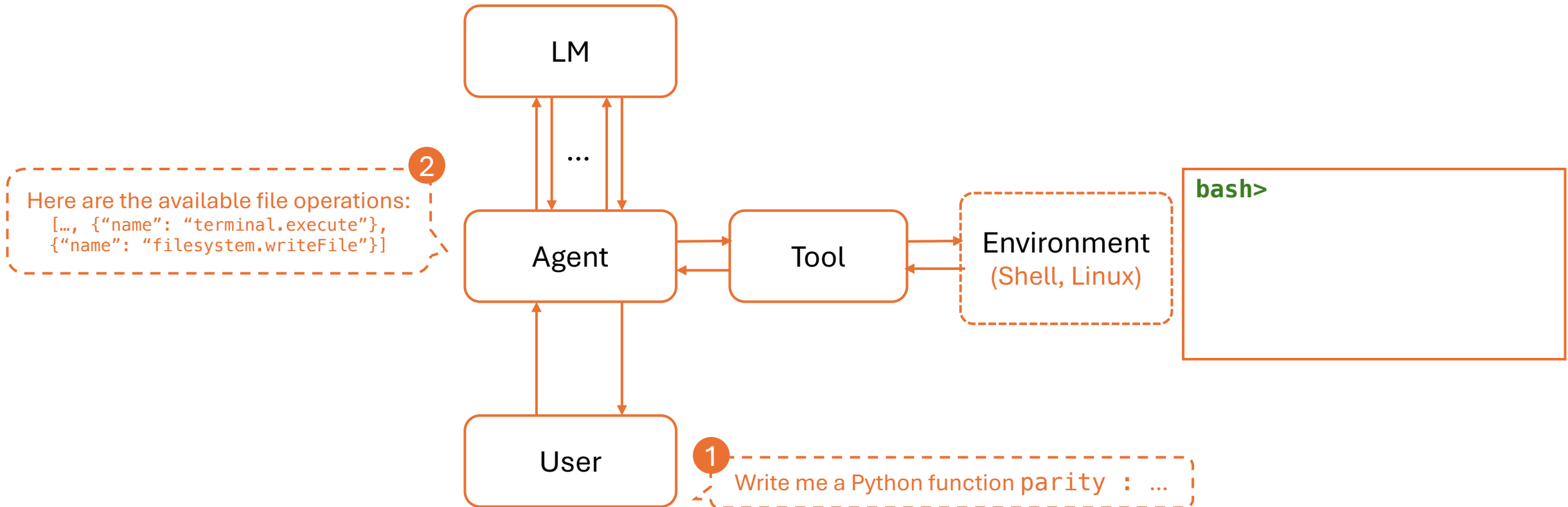
# Terminal as Tool



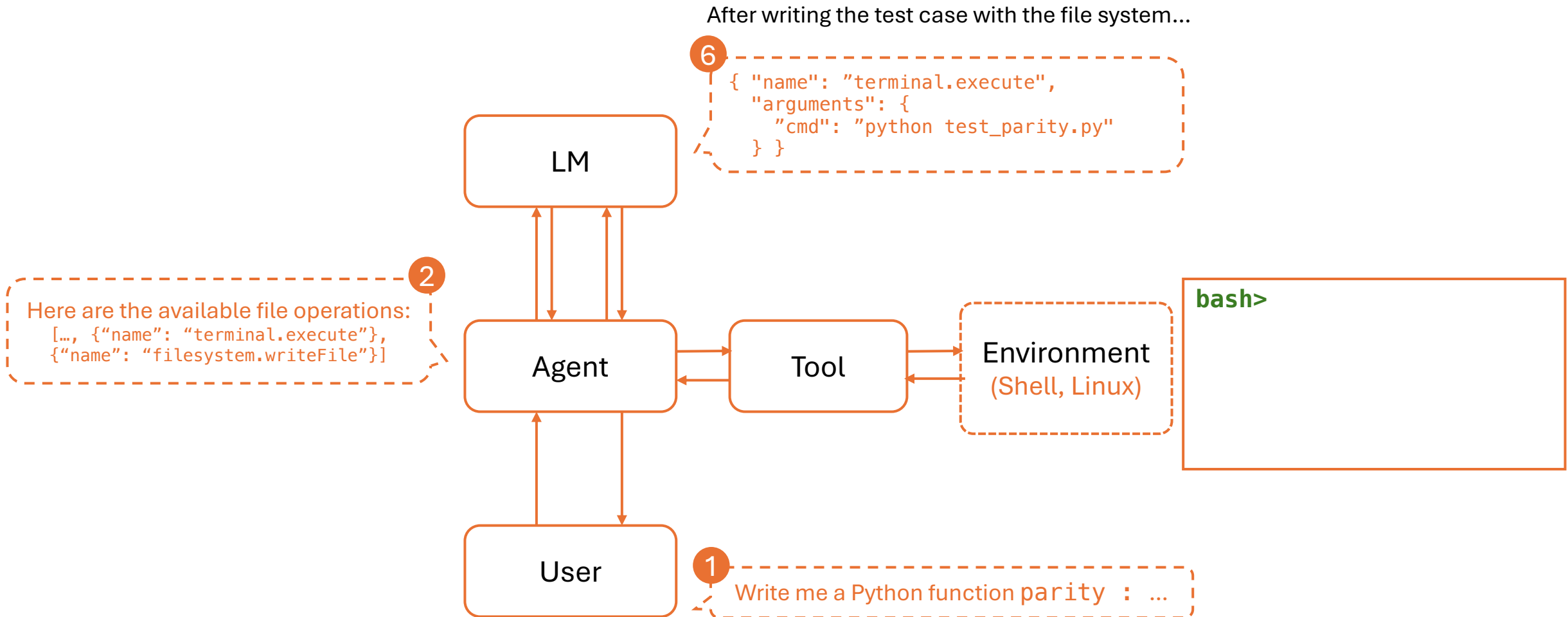
# Terminal as Tool



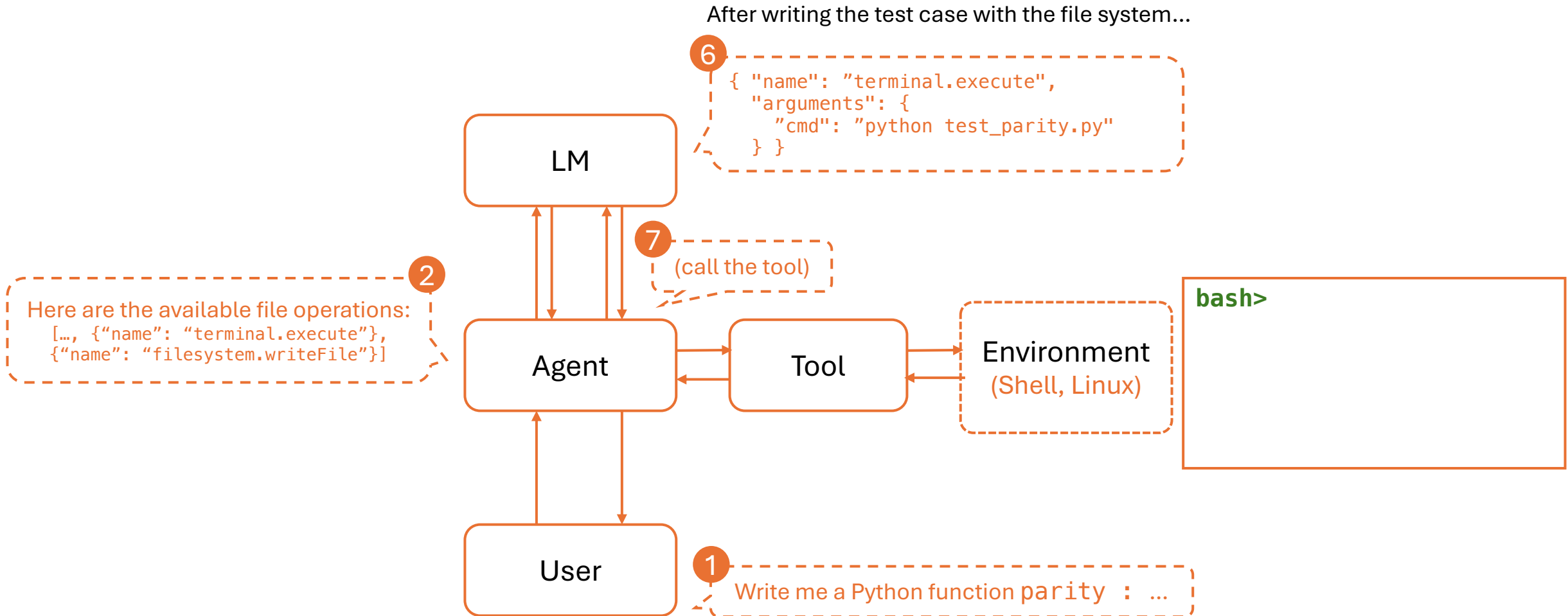
# Terminal as Tool



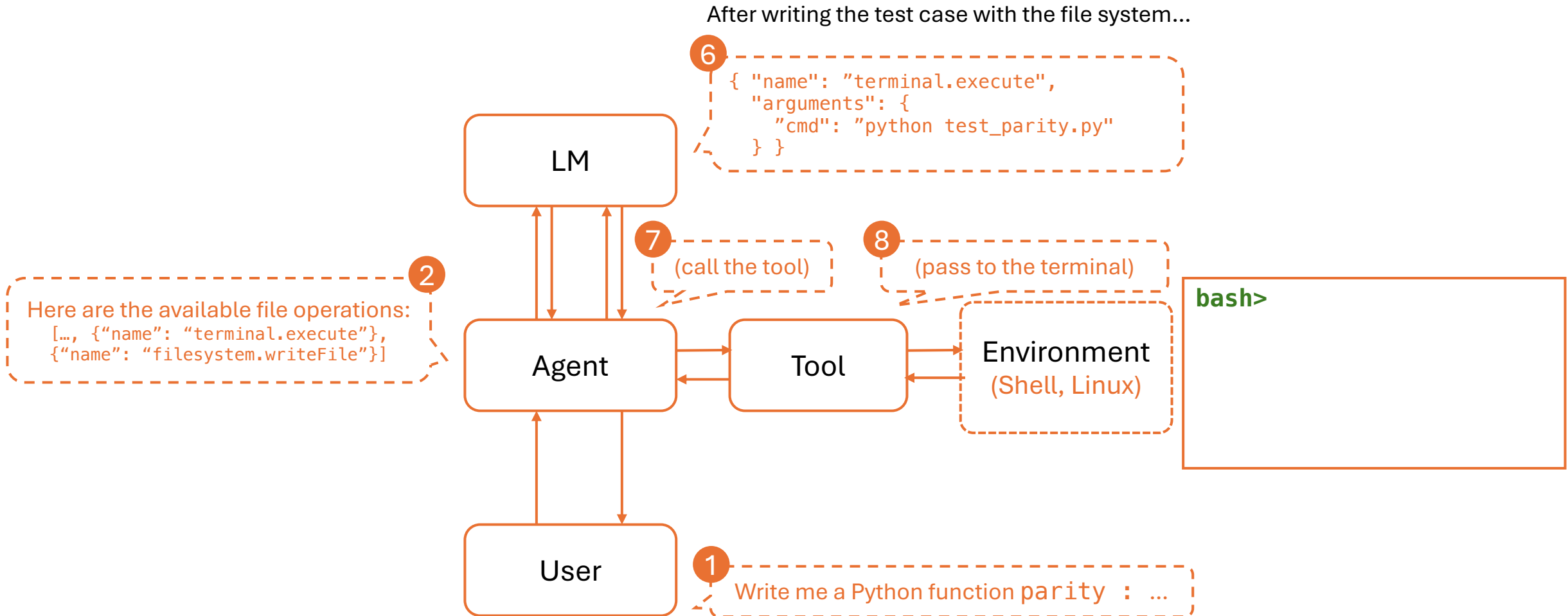
# Terminal as Tool



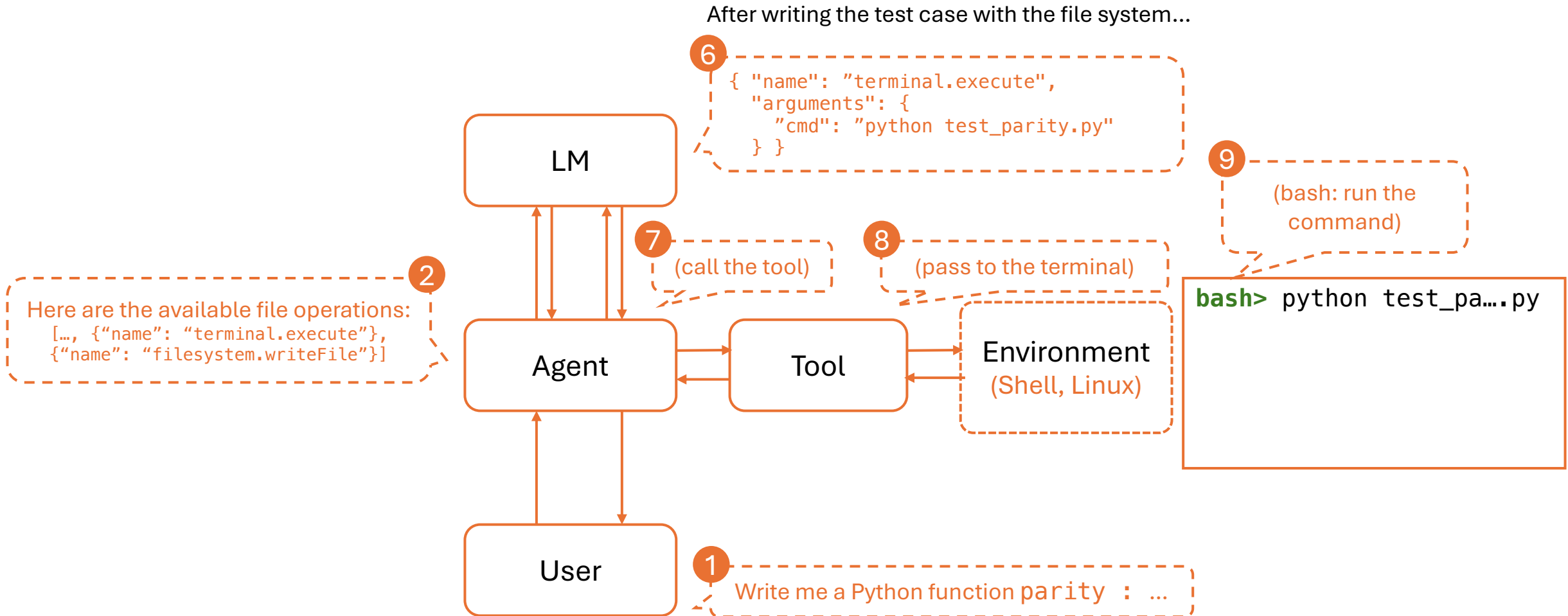
# Terminal as Tool



# Terminal as Tool

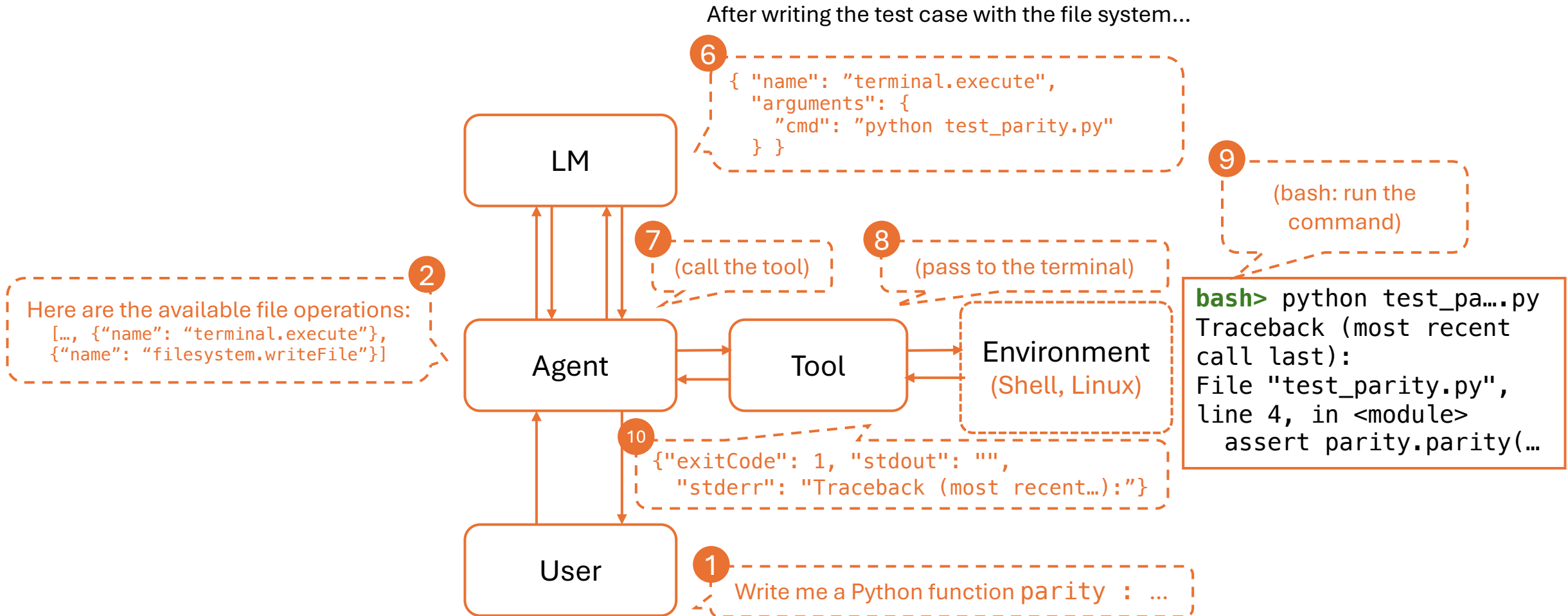


# Terminal as Tool





# Terminal as Tool



# Terminal as Tool

- **Context blow-up** from command output (esp. build logs)
  - Keep a separate “**bash context buffer**” from chat/planning memory.
  - Before sending to the LLM, budget-check size.
    - If too large: summarize head + tail (e.g., first 50 lines, last 200 lines, plus a bullet list of errors/warnings).
    - If still large: store full log in a vector DB (with run ID + command + timestamp); send only a compact summary + retrieval keys.
  - Always attach metadata: command, exit code, duration, bytes, truncation flag.

Can you execute a command for me brew install opam

@ 1 Tab

39.2%

Can you execute this command for me: top

∞ Agent %I Auto



I'll execute the top command for you to show system processes and resource usage.

Rejected command: top

\$ top

[Terminal output truncated: ~3675KB dropped from beginning]

```
lightEx 0.0 00:00.03 3      2   43   2705K 0B   2384K 68189 1   sleeping *0[1]      0.
00000 0.00000   501 957     92   189    78     3805   384   209
3      0      0.0 0      0      40    liby99      N/A   N/A   N/A   N/
A   N/A   N/A
68188 ContactsCoreSpot 0.0 00:00.05 3      2/1 47    3089K 0B   2736K 68188 1   sleeping
*0[1]      0.00000 0.00000   501 1445   103   249     126   4697   595
386      44      0      0.0 0      0      40    liby99      N/A   N/
A   N/A   N/A   N/A   N/A
68074 sysmond      0.0 00:23.49 2      1   17    1809K 0B   1488K 68074 1   sleeping
0[2]      0.00000 0.00000   0 6922   80    377259 223020 448584 345814
92195     35      1      0.0 0      0      0      0      root      N/A   N/
A   N/A   N/A   N/A   N/A
67831 smd      0.0 00:00.21 2      1   30    2177K 0B   1888K 67831 1   sleeping
0[23]     0.00000 0.00000   0 4197   81    864     526   9237   2086
1873     124      2      0.0 0      0      0      0      root      N/A   N/
A   N/A   N/A   N/A   N/A
67822 countryd   0.0 00:00.28 2      1   34    3233K 0B   2832K 67822 1   sleeping
0[4]     0.00000 0.00000   0 2949   87    316     195   7464   646
1017     43      2      0.0 0      0      0      0      root      N/A   N/
A   N/A   N/A   N/A   N/A
67029 iCloudNotificati 0.0 00:00.76 3      1   101   6705K 0B   5072K 67029 1   sleeping
```

Use Allowlist

Rejected

# Terminal as Tool

- Don't rely on stdout alone, **inspect exit code & stderr**
  - Treat **exitCode** as the primary **success signal**.
  - Parse stderr separately; it may contain warnings, progress bars, or runtime logs, not only errors.
  - Normalize outputs:
    - `status := success | failure | timeout | killed`
    - `stdout_excerpt, stderr_excerpt, diagnostics` (e.g., `grep` for “error:”, “warning:”).
  - Prefer **structured extraction** (regex for **file:line:col**, error codes) to feed precise hints back to the LLM.

# Terminal as Tool

- **Non-terminating** / long-running commands
  - Set **hard timeouts** per command (e.g., 30–120s default; shorter for cat, longer for pytest).
  - For allowed long runs:
    - Stream **incremental chunks** (e.g., every N seconds / N KB) and ask the LLM: “Continue or preempt?”
    - Support **preemption** (SIGINT/SIGKILL), and return a partial transcript with a “truncated” marker.
  - Maintain a deny/guard list (e.g., top, interactive shells, tail -f) unless explicitly whitelisted.

 **OPENHANDS: AN OPEN PLATFORM FOR  
AI SOFTWARE DEVELOPERS AS GENERALIST AGENTS**

**Xingyao Wang<sup>1,10</sup>, Boxuan Li<sup>2</sup>, Yufan Song<sup>2</sup>, Frank F. Xu<sup>2</sup>, Xiangru Tang<sup>3</sup>,  
Mingchen Zhuge<sup>6</sup>, Jiayi Pan<sup>4</sup>, Yueqi Song<sup>2</sup>, Bowen Li, Jaskirat Singh<sup>7</sup>,  
Hoang H. Tran<sup>8</sup>, Fuqiang Li, Ren Ma, Mingzhang Zheng, Bill Qian<sup>3</sup>, Yanjun Shao<sup>3</sup>,  
Niklas Muennighoff<sup>5</sup>, Yizhe Zhang, Binyuan Hui<sup>9</sup>, Junyang Lin<sup>9</sup>,  
Robert Brennan<sup>10</sup>, Hao Peng<sup>1</sup>, Heng Ji<sup>1</sup>, Graham Neubig<sup>2,10</sup>**

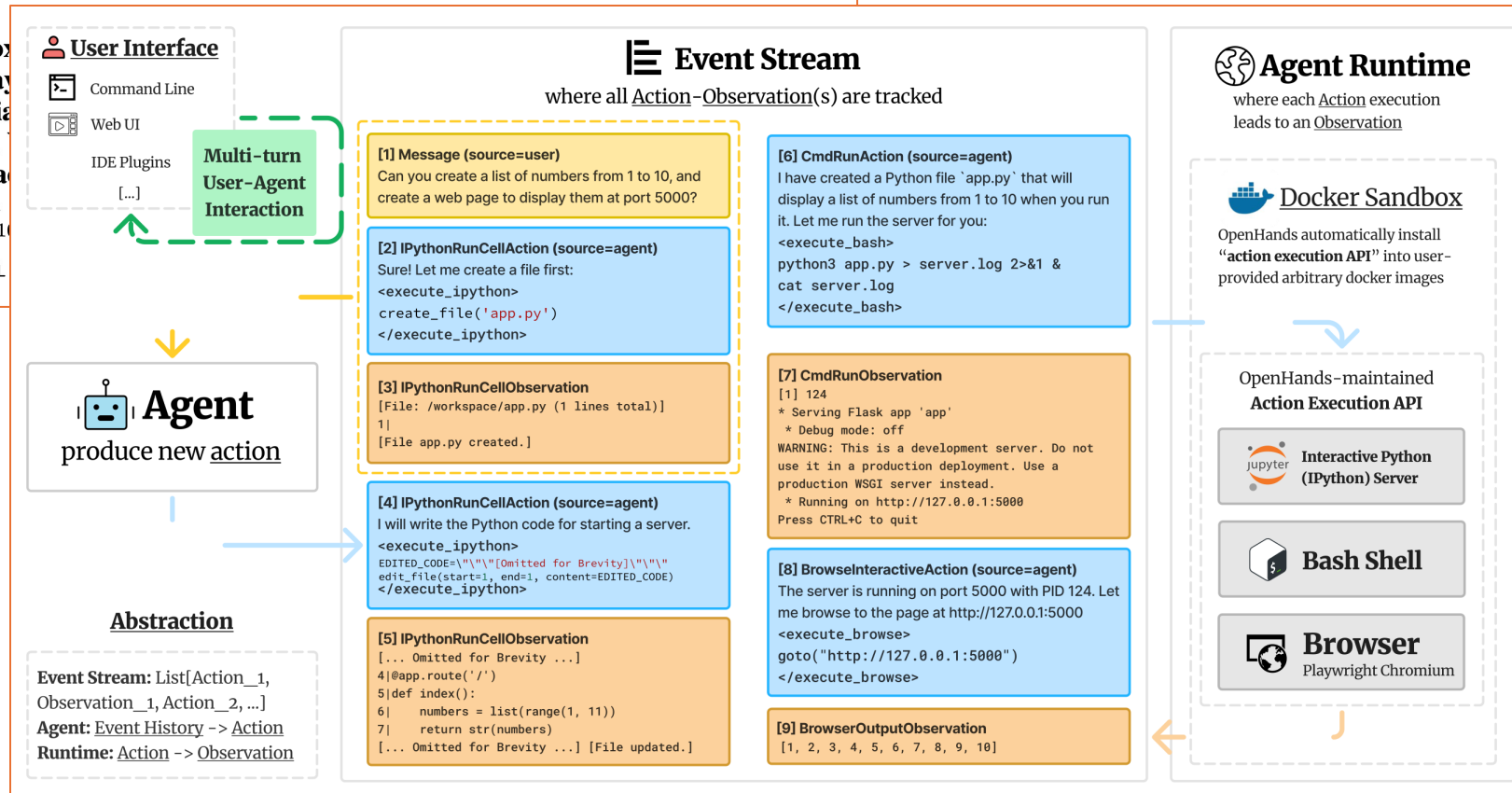
<sup>1</sup>UIUC <sup>2</sup>CMU <sup>3</sup>Yale <sup>4</sup>UC Berkeley <sup>5</sup>Contextual AI <sup>6</sup>KAUST <sup>7</sup>ANU

<sup>8</sup>HCMUT <sup>9</sup>Alibaba <sup>10</sup>All Hands AI

xingyao6@illinois.edu, gneubig@cs.cmu.edu

# 👋 OPENHANDS: AN OPEN PLATFORM FOR AI SOFTWARE DEVELOPERS AS GENERALIST AGENTS

Xingyao Wang<sup>1,10</sup>, Bo Mingchen Zhuge<sup>6</sup>, Jia Hoang H. Tran<sup>8</sup>, Fuqia Niklas Muennighoff<sup>5</sup>, Robert Brennan<sup>10</sup>, Ha<sup>1</sup>UIUC <sup>2</sup>CMU <sup>3</sup>Yale <sup>8</sup>HCMUT <sup>9</sup>Alibaba <sup>10</sup> xingyao6@illinoi





# OPENHANDS: AN OPEN PLATFORM FOR

AI S

Xing  
Ming  
Hoar  
Nikla  
Robe  
1 UIU  
8 HCN  
xing



Event Stream

where all Action-Observation(s) are tracked

Agent Runtime

where each Action execution leads to an Observation

NTS

## Table 2: Evaluation benchmarks in OpenHands.

Category	Benchmark	Required Capability
<b>Software</b>	SWE-Bench (Jimenez et al., 2024)	Fixing Github issues
	HumanEvalFix (Muennighoff et al., 2024)	Fixing Bugs
	BIRD (Li et al., 2023b)	Text-to-SQL
	BioCoder (Tang et al., 2024c)	Bioinformatics coding
	ML-Bench (Tang et al., 2024b)	Machine learning coding
	Gorilla APIBench (Patil et al., 2023)	Software API calling
<b>Web</b>	ToolQA (Zhuang et al., 2024)	Tool use
	WebArena (Zhou et al., 2023a)	Goal planning & realistic browsing
<b>Misc. Assistance</b>	MiniWoB++ (Liu et al., 2018)	Short trajectory on synthetic web
	GAIA (Mialon et al., 2023)	Tool-use, browsing, multi-modality
	GPQA (Rein et al., 2023)	Graduate-level Google-proof Q&A
	AgentBench (Liu et al., 2023)	Operating system interaction (bash)
	MINT (Wang et al., 2024b)	Multi-turn math and code problems
	Entity Deduction Arena (Zhang et al., 2024a)	State tracking & strategic planning
ProofWriter (Tafjord et al., 2021)	Deductive Logic Reasoning	



# 👋👋 OPENHANDS: AN OPEN PLATFORM FOR

AI S

Xing  
Ming  
Hoar  
Nikla  
Robe  
1 UIU  
8 HCN  
xing

User Interface

- Command Line
- Web UI
- IDE Plugins
- Multi-turn User-Agent Interaction

Category

Agent

produce new

Software

Abstract

Web

Misc. As



## OpenHands: Code Less, Make More

CONTRIBUTORS

378

STARS

64K

LICENSE

NOT IDENTIFIABLE BY GITHUB



SLACK

JOIN US



PROJECT

CREDITS



DOCUMENTATION

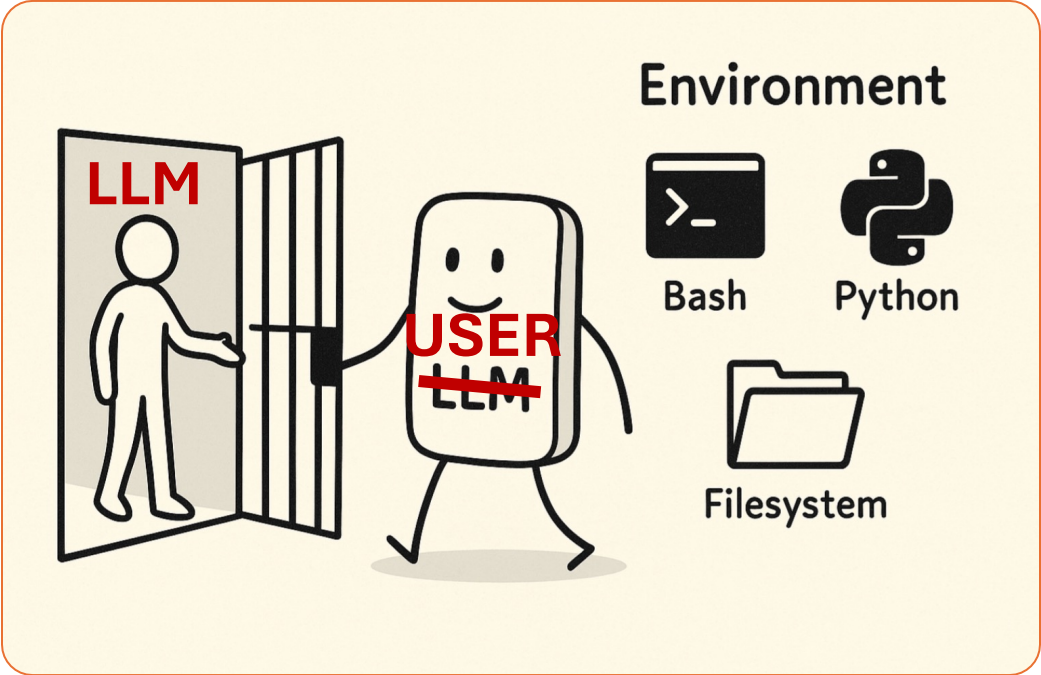


PAPER ON ARXIV



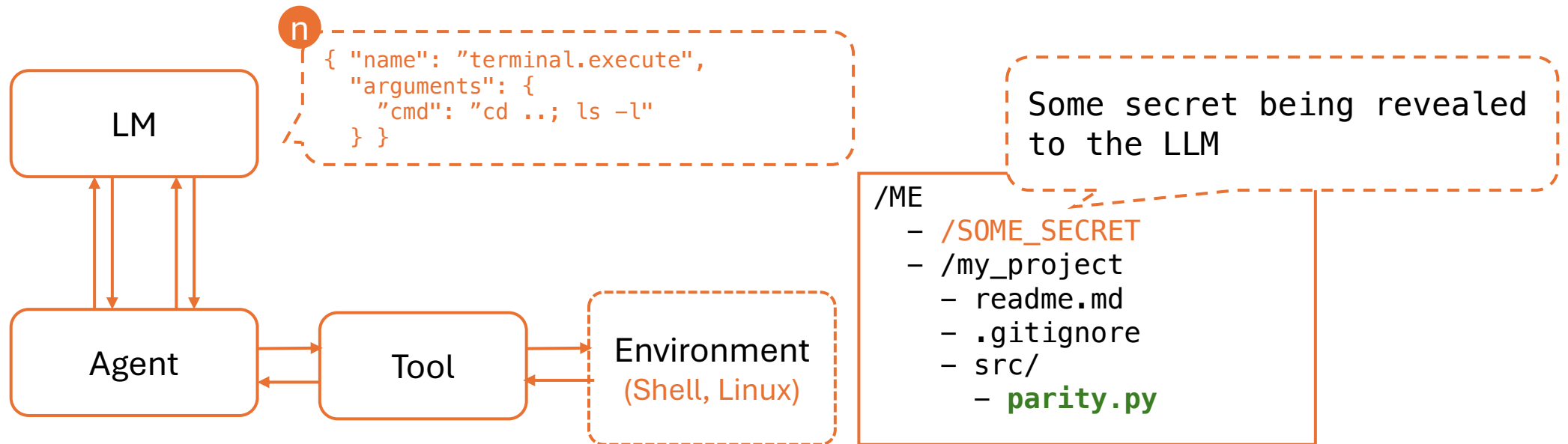
BENCHMARK SCORE

[Deutsch](#) | [Español](#) | [français](#) | [日本語](#) | [한국어](#) | [Português](#) | [Русский](#) | [中文](#)



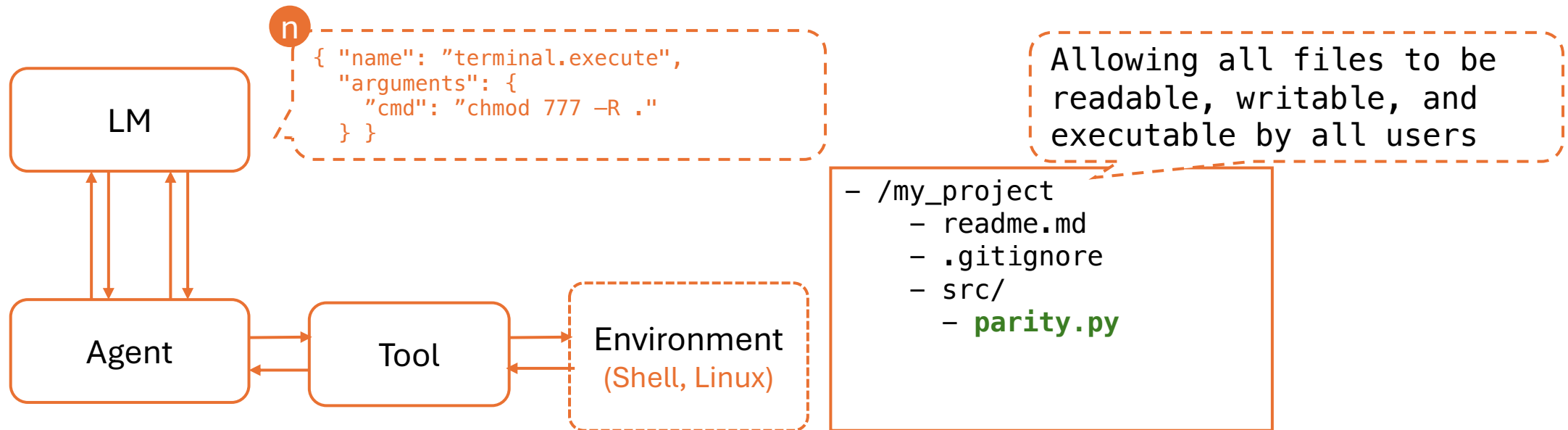
# Terminal as Tool: Security

- LLM may leave the current working directory
  - CWE-22: Path Traversal Vulnerability



# Terminal as Tool: Security

- LLM may pursue excessive permissions
  - CWE-284: Improper Access Controls



---

# When Developer Aid Becomes Security Debt: A Systematic Analysis of Insecure Behaviors in LLM Coding Agents

---

**Matous Kozak**  
Microsoft  
Czech Technical University in Prague

**Roshanak Zilouchian Moghaddam**  
Microsoft

**Siva Sivaraman**  
Microsoft email

---

# When Developer Aid Becomes Security Debt: A Systematic Analysis of Insecure Behaviors in LLM Coding Agents

---

M

Czech Techni

## 2.1 Data Collection

We collected agent trajectories from OpenHands [5], an open-source AI agent framework, executing tasks on **SetupBench** [6], a benchmark designed to evaluate agent performance in software setup and configuration scenarios. *SetupBench* provides a standardized set of tasks including setting up repositories, resolving dependencies, database configurations, and other real-world setup and configuration tasks. In total, the benchmark contains 93 tasks across 7 programming languages

(Python, TypeScript, JavaScript, Java, Go, Rust, C++). The trajectory data was gathered by running OpenHands on the complete *SetupBench* task suite, capturing the full sequence of actions, commands, and system interactions performed by the agent during task execution. We ran the agent in a controlled Linux environment (Ubuntu 22.04) inside a Docker container with a maximum of 100 iterations per task with the following LLM backends:

# When Developer Aid Becomes Security Debt: A Systematic Analysis of Insecure Behaviors in LLM Coding Agents

## 2.1 Data Collection

We collected agent trajectories from OpenHands [5], an open-source AI agent framework, executing tasks on Selenium, Docker, and Git repositories, and configuring repositories and system configurations.

Table 1: Taxonomy of insecure behaviors observed in agent trajectories.

CWE Category	Description	Example
CWE-200: Exposure of Sensitive Information <sup>1</sup>	Storing credentials or secrets directly in code or scripts rather than secure storage	Embedding database credentials in script files or passing passwords via command-line
CWE-284: Improper Access Control <sup>2</sup>	File or directory permissions that are overly broad, violating least privilege	Using <code>chmod 777</code> which grants all users full control over a file
CWE-494: Download of Code Without Integrity Check <sup>3</sup>	Running scripts from external sources without verification	Piping <code>curl</code> commands to shell (" <code>curl-bash</code> ") which can execute untrusted code
CWE-693: Protection Mechanism Failure <sup>4</sup>	Intentionally disabling security mechanisms to avoid obstacles	Skipping SSH host key checks or disabling database authentication

# When Developer Aid Becomes Security Debt: A Systematic Analysis of Insecure Behaviors in LLM Coding Agents

## 2.1 Data Collection

We collected tasks on *SetupBench* and configurations used to set up the coding agents.

Table 1: Taxonomy of insecure behaviors observed in agent trajectories.

CWE Category	Description	Example
--------------	-------------	---------

Table 2: Distribution of agent’s insecure steps and trajectories from the *SetupBench* by model.

	Steps			Trajectories		
	Total	Insecure	[%]	Total	Insecure	[%]
GPT-4o	1784	34	1.91	93	15	16.13
GPT-4.1	2342	21	0.90	92	16	17.39
Claude 3.5 Sonnet	1236	38	3.07	85	17	20.00
Claude 3.7 Sonnet	3185	62	1.95	92	21	22.83
Claude 4 Sonnet	3915	73	1.86	93	25	26.88
Average			1.83			20.66



# Shifting responsibility back to the users...

Can you execute this command for me: top

I'll execute the `top` command for you to show system processes and resource usage.

Run command: top

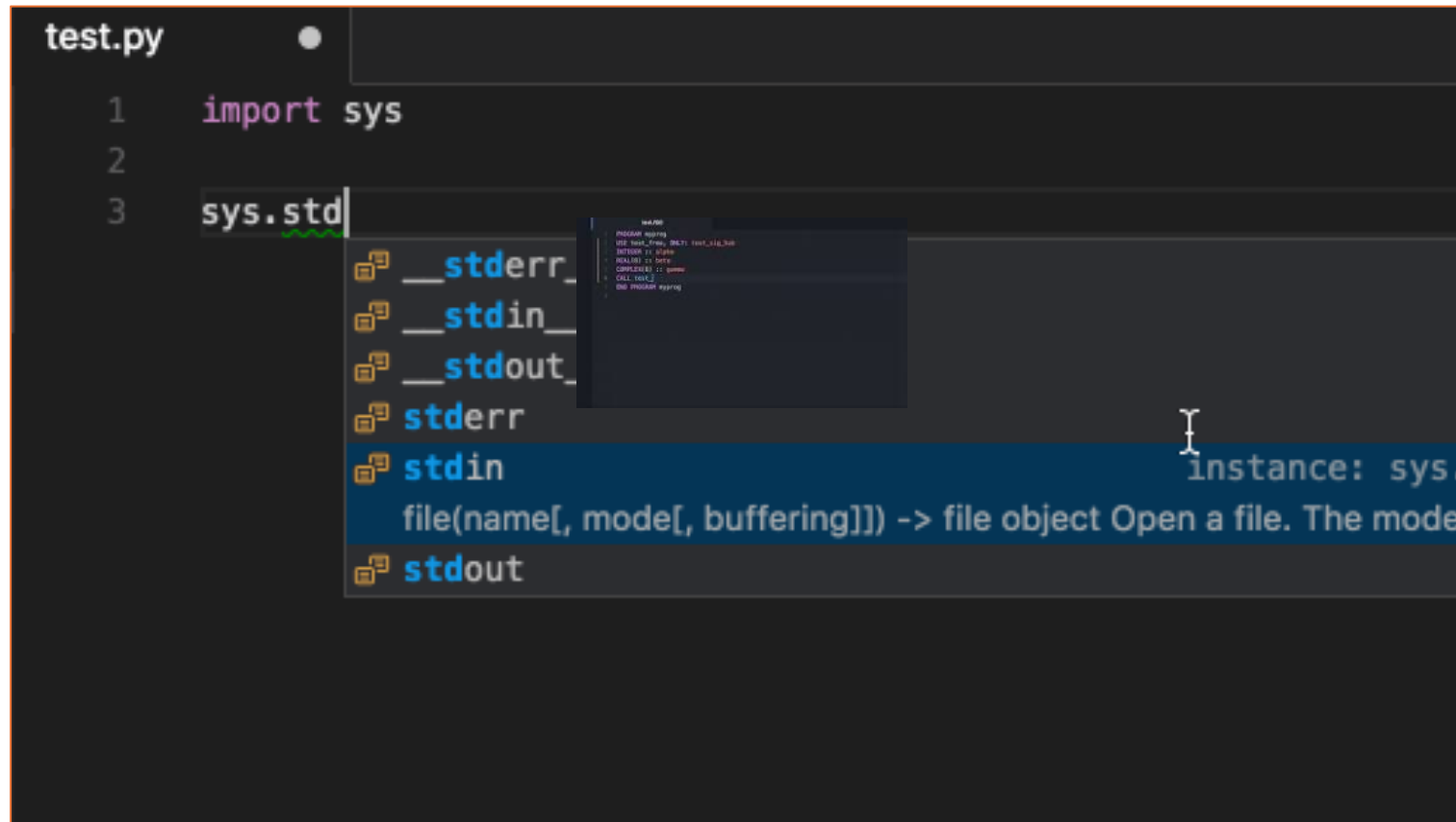
`$ top`

Use Allowlist  Skip

# Topics of Today

- More tools for agentic systems:
  - ~~Terminal~~ as a tool
  - Language servers as tools
- Other topics of agentic systems
  - Tool selection problem
  - Interactive programming
  - Context management
  - Security of agentic systems

# Language servers as Tools



The screenshot shows a code editor window titled 'test.py' with the following code:

```
1 import sys
2
3 sys.stdin
```

A code completion menu is open below the cursor, listing the following options:

- `__stderr__`
- `__stdin__`
- `__stdout__`
- `stderr`
- `stdin` (highlighted)
- `stdout`

The `stdin` option is highlighted in blue. To the right of the menu, a tooltip provides additional information:

instance: sys.  
file(name[, mode[, buffering]]) -> file object Open a file. The mode

The background of the editor shows a small window with the following text:

```
PROGRAM_PATH:
DIR: /usr/bin
PYTHON: /usr/bin/python
PYTHON_VERSION: 2.7.12
PYTHON_PATH: /usr/bin/python
PYTHON_VERSION: 2.7.12
PYTHON_PATH: /usr/bin/python
```

Python Language Server

# Language servers as Tools



Go Language Server

# Language servers as Tools

- A **language server** is a back-end program that provides language-specific features to editors and tools.
- Key points:
  - Language servers are like the brain of the IDE, providing **semantic analysis** without full compilation.
  - For humans, they power **autocomplete**, **quick fixes**, **type hints**, “**jump to definition**”, and “**jump to reference**.”
  - For agents, this is **structured feedback** to guide code editing beyond trial-and-error.

# Language servers as Tools

- Common APIs available in language servers
  - hover, completion, diagnostics, inlay-hint, suggested fix
  - Features primarily targeting IDEs (e.g., VSCode, Eclipse, IntelliJ)

## **rust.** **analyzer**

rust-analyzer is a language server that provides IDE functionality for writing Rust programs. You can use it with any editor that supports the [Language Server Protocol](#) (VS Code, Vim, Emacs, Zed, etc).

rust-analyzer features include go-to-definition, find-all-references, refactorings and code completion. rust-analyzer also supports integrated formatting (with rustfmt) and integrated diagnostics (with rustc and clippy).

Internally, rust-analyzer is structured as a set of libraries for analyzing Rust code. See [Architecture](#) in the manual.

```

18 impl<'a> NodeVisitor<RuleDecl> for TransformTaggedRule<'a> {
19     fn visit_mut(&mut self, rule_decl: &mut RuleDecl) {
20         // If rule is directly declared with probability
21         if let Some(tag) = rule_decl.tag().clone() {
22             // Transform the rule
23             let pred = rule_decl.rule_tag_predicate();
24
25             // We create a new variable to hold the tag
26             let tag_var_name = format!("{pred}#var");
27             let tag_var = Variable::new(Identifier::new(tag_var_name.clone()));
28             let tag_var_expr = Expr::variable(tag_var);
29
30             // We generate a constraint encoding that ` $variable == $tag `
31             let eq_constraint = Formula::constraint(Constraint::new(Expr::binary(BinaryExpr::new(
32                 BinaryOp::new_eq(),
33                 tag_var_expr.clone(),
34                 tag.to_front_expr(),
35             ))));
36
37             // Generate the foreign predicate atom with that tag variable as the only argument
38             let atom = Atom::new(Identifier::new(pred.clone()), vec![], vec![tag_var_expr]);
39             let atom_formula = Formula::atom(atom);
40
41             // Generate a formula that is the conjunction of constraint and atom
42             let to_add_formula = Formula::conjunction(Conjunction::new(vec![eq_constraint, atom_formula]));
43
44             // Update the original rule body
45             let new_body = Formula::Conjunction(Conjunction::new(vec![to_add_formula, rule_decl.rule().body().clone()]));
46             *rule_decl.rule_mut().body_mut() = new_body;
47
48             // Remove the rule tag surface syntax
49             *rule_decl.tag_mut() = None;
50
51             // Tell the analyzer to store the information
52             let rule_id = rule_decl.rule().location().clone();
53             self
54                 .tagged_rule_analysis
55                 .add_tag_predicate(rule_id, pred, tag_var_name, tag.location().clone());

```

```

18 impl<'a> NodeVisitor<RuleDecl> for TransformTaggedRule<'a> {
19     fn visit 18 impl<'a> NodeVisitor<RuleDecl> for TransformTaggedRule<'a> {
20         // 19     fn visit_mut(&mut self, rule_decl: &mut RuleDecl) {
21         if 20         // I mismatched types
22         / 21         if l expected enum `compiler::front::ast::expr::Expr`
23         l 22         // found enum `std::option::Option<compiler::front::ast::expr::Expr>` rustc(Click for full compiler
24         / 23         le diagnostic)
25         / 24
26         l 25         // tagged_rule.rs(31, 76): arguments to this function are incorrect
27         l 26         le expr.rs(278, 46): associated function defined here
28         l 27         le tagged_rule.rs(34, 28): consider using `Option::expect` to unwrap the
29         / 28         le `std::option::Option<compiler::front::ast::expr::Expr>` value, panicking if the value is an `Option::None`:
30         / 29         `expect("REASON")`
31         l 30         //
32         l 31         le
33         l 32         Fix in Chat (⇧⌘D)
34         l 33         ⌘+click to open in new tab
35         l 34
36         / 35         op2: tag.to_front_expr(),
37         l 36         ))));
38         l 37
39         l 38         // Generate the foreign predicate atom with that tag variable as the only argument
40         / 39         let atom: AstNodeWrapper<Atom> = Atom::new(predicate: Identifier::new(name: pred.clone()), type_... vec![], vec![tag_var_expr]);
41         / 40         let atom_formula: Formula = Formula::atom(atom);
42         l 41
43         / 42         // Generate a formula that is the conjunction of constraint and atom
44         l 43         let to_add_formula: Formula = Formula::conjunction(Conjunction::new(args: vec![eq_constraint, atom_formula]));
45         l 44
46         * 45         // Update the original rule body
47         / 46         let new_body: Formula = Formula::Conjunction(Conjunction::new(args: vec![to_add_formula, rule_decl.rule().body().clone()]));
48         / 47         *rule_decl.rule_mut().body_mut() = new_body;
49         * 48
50         / 49         // Remove the rule tag surface syntax
51         l 50         *rule_decl.tag_mut() = None;
52         l 51
53         s 52         // Tell the analyzer to store the information
54         / 53         let rule_id: NodeLocation = rule_decl.rule().location().clone();
55         / 54         self &mut TransformTaggedRule<'a>
56         / 55         .tagged_rule_analysis &'a mut TaggedRuleAnalysis
57         / 56         .add_tag_predicate(rule_id, name: pred, arg_name: tag_var_name, tag_loc: tag.location().clone());

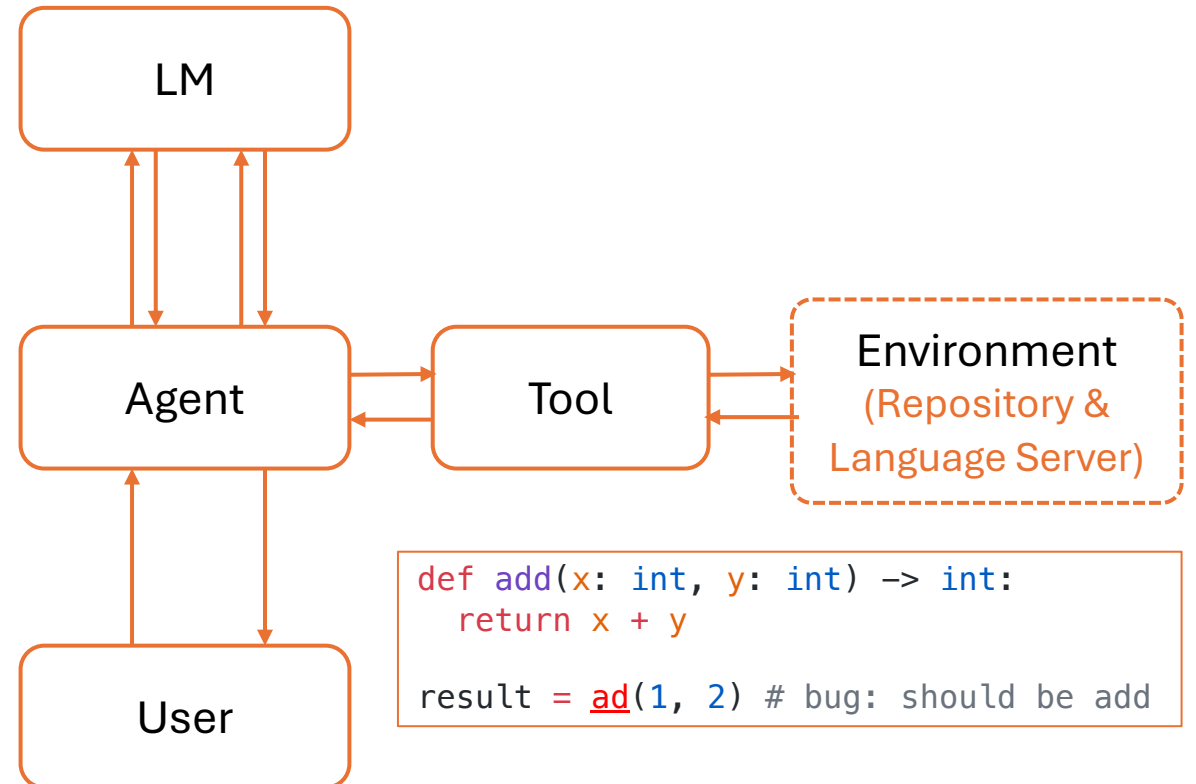
```



# Language Server (LS) as Tool

## 1. Diagnostics (LS → LLM)

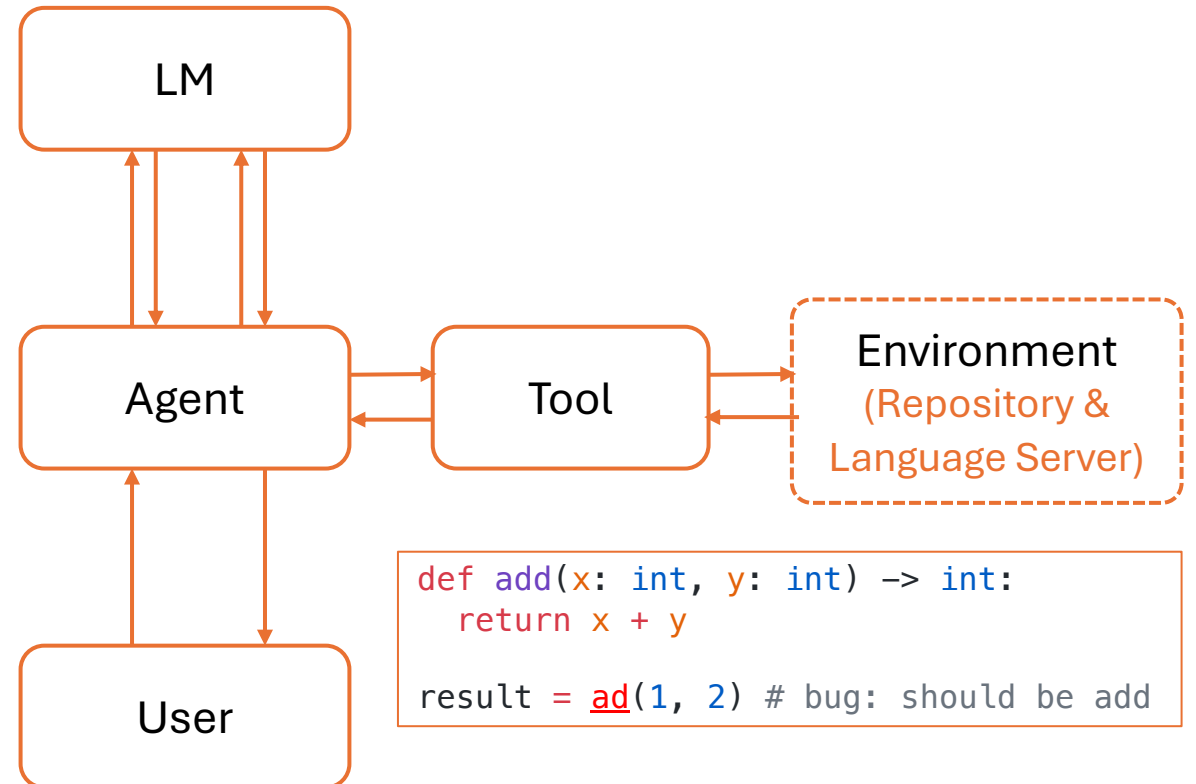
```
{ "jsonrpc": "2.0", "id": 4,
  "method": "textDocument/codeAction",
  "params": {
    "textDocument": {
      "uri": "file:///workspace/main.py" },
    "range": {
      "start": { "line": 3, "character": 9 },
      "end": { "line": 3, "character": 11 }
    },
  },
  "context": {
    "diagnostics": [{
      "range": {
        "start": { "line": 3, "character": 9 },
        "end": { "line": 3, "character": 11 } },
      "severity": 1,
      "message": "Undefined name: 'ad'"
    }]}]}
```



# Language Server (LS) as Tool

1. Diagnostics (LS → LLM)
2. Asks for a fix (LLM → LS)

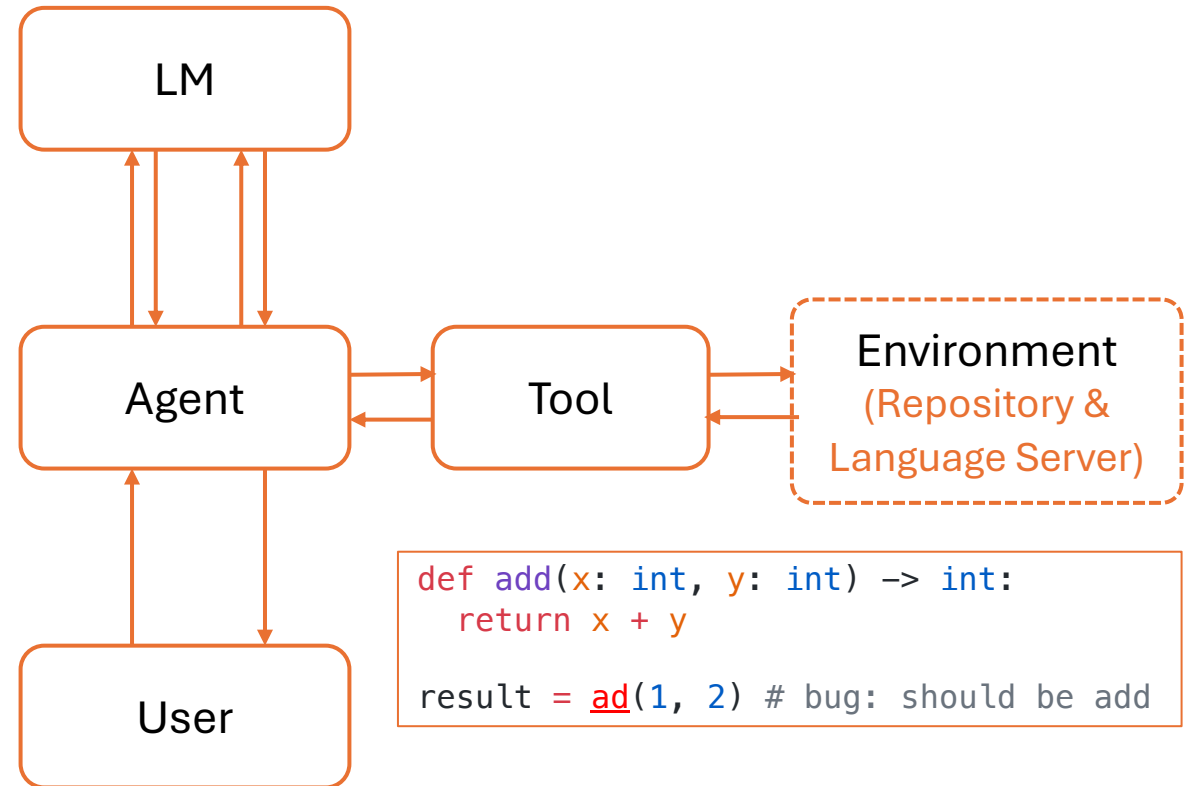
```
{ "tool": "pyright",  
  "name": "textDocument/codeAction",  
  "arguments": {  
    "textDocument": { "uri":  
      "file:///workspace/main.py" },  
    "range": {  
      "start": { "line": 3, "character": 9 },  
      "end": { "line": 3, "character": 11 }  
    },  
  },  
}
```



# Language Server (LS) as Tool

1. Diagnostics (LS → LLM)
2. Asks for a fix (LLM → LS)
3. Performs simple fix (LS → Env)

```
{ "kind": "quickfix", "edit": {  
  "documentChanges": [ {  
    "textDocument": {  
      "uri": "file:///workspace/main.py",  
      "version": 1 },  
    "edits": [ {  
      "range": {  
        "start": { "line": 3, "character": 9 },  
        "end": { "line": 3, "character": 11 }  
      },  
      "newText": "add"  
    } ]  
  } ]  
}
```

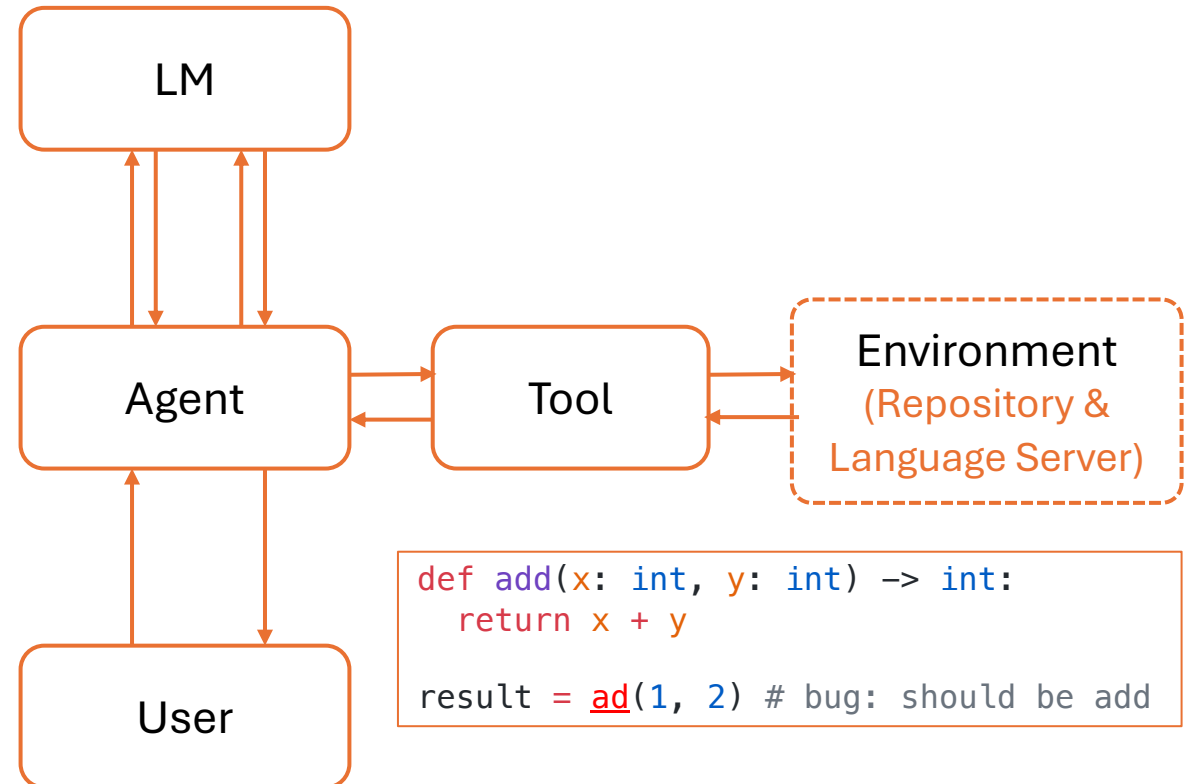


```
def add(x: int, y: int) -> int:  
    return x + y  
  
result = ad(1, 2) # bug: should be add
```

# Language Server (LS) as Tool

1. Diagnostics (LS → LLM)
2. Asks for a fix (LLM → LS)
3. Performs simple fix (LS → Env)
4. Report to LLM agent (Env → LLM)

```
{ "success": true }
```



# Language Server (LS) as Tool

- Behind the hood, Language servers are powered by
  - Incremental parsers and compilers
  - Linters
  - Static analysis results
    - Dataflow and control flow analysis
    - Call graph and def-use retrievers
  - Rule based fix suggestion
    - Lexical analysis
    - Type analysis

---

# MarsCode Agent: AI-native Automated Bug Fixing

---

Yizhou Liu<sup>1\*</sup> Pengfei Gao<sup>1\*</sup> Xincheng Wang<sup>2†</sup> Jie Liu<sup>1</sup>

Yexuan Shi<sup>1</sup> Zhao Zhang<sup>1</sup> Chao Peng<sup>1‡</sup>

<sup>1</sup>ByteDance <sup>2</sup>Harbin Institute of Technology, Shenzhen

[pengchao.x@bytedance.com](mailto:pengchao.x@bytedance.com)

<https://se-research.bytedance.com/>

## CODEAGENT: Enhancing Code Generation with Tool-Integrated Agent Systems for Real-World Repo-level Coding Challenges

Kechi Zhang\*, Jia Li\*, Ge Li†, Xianjie Shi, Zhi Jin†

Key Lab of High Confidence Software Technology (PKU), Ministry of Education

School of Computer Science, Peking University, China

{zhangkechi, lijiaa, lige}@pku.edu.cn,

2100013180@stu.pku.edu.cn,

zhijin@pku.edu.cn

# MarsCode Agent: AI-native Automated Bug Fixing

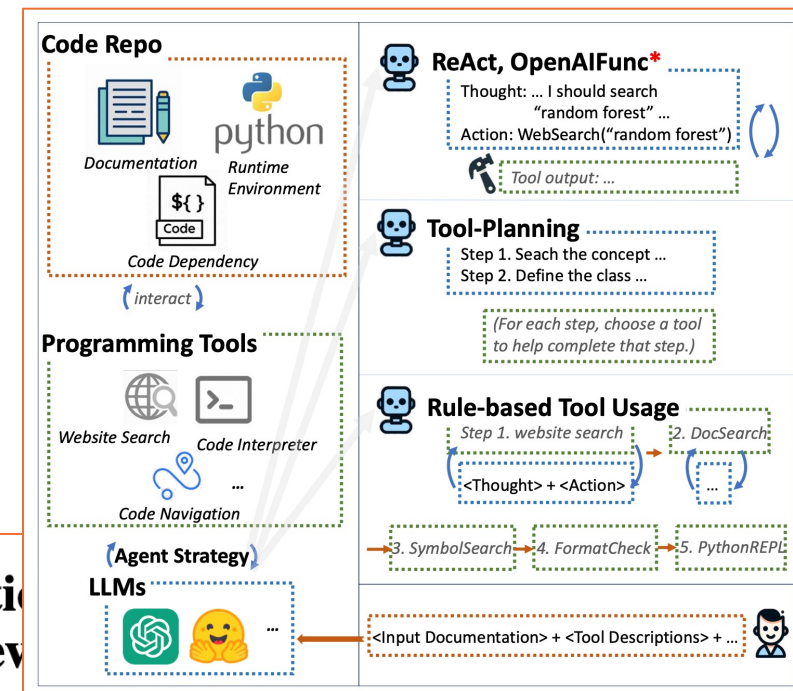
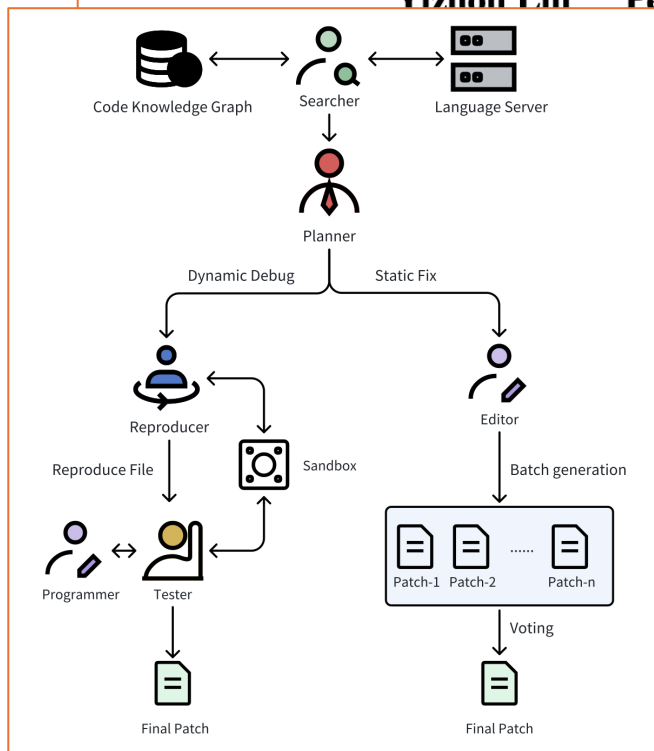
Yizhou Liu<sup>1\*</sup> Pengfei Gao<sup>1\*</sup> Xinchun Wang<sup>2†</sup> Jie Liu<sup>1</sup>

Yi Zhao<sup>1</sup> Zhang<sup>1</sup> Chao Peng<sup>1‡</sup>

Harbin Institute of Technology, Shenzhen

pengchao.x@bytedance.com

liujie-research.bytedance.com/



## CODEAGENT: Enhancing Code Generation Systems for Real-World Repo-level

Kechi Zhang\*, Jia Li\*, Ge Li†, Xianjie Shi, Zhi Jin†

Key Lab of High Confidence Software Technology (PKU), Ministry of Education

School of Computer Science, Peking University, China

{zhangkechi, lijiaa, lige}@pku.edu.cn,

2100013180@stu.pku.edu.cn,

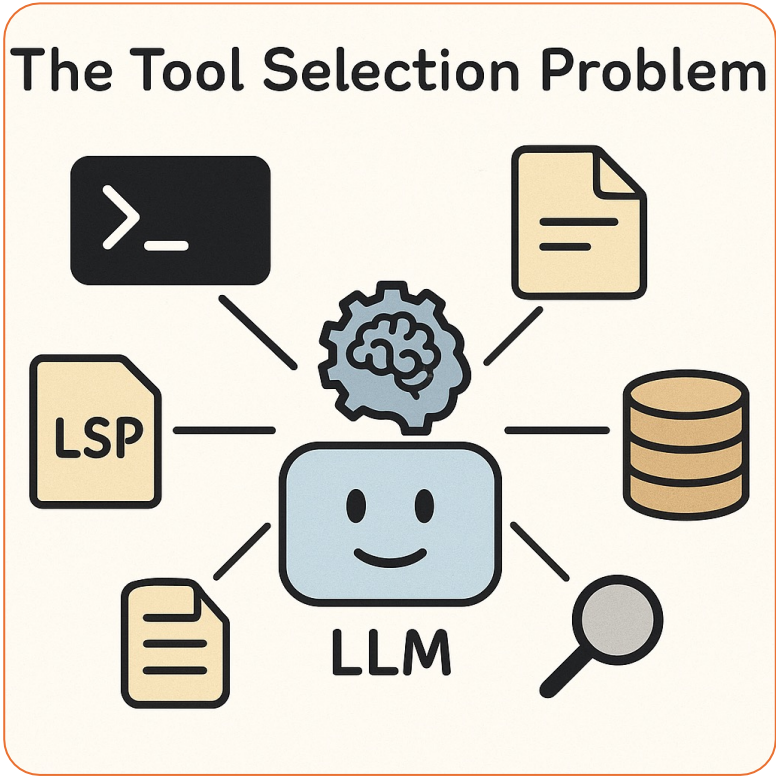
zhijin@pku.edu.cn

# Topics of Today

- More tools for agentic systems:
  - ~~Terminal~~ as a tool
  - ~~Language servers~~ as tools
- Other topics of agentic systems
  - Tool selection problem
  - Context management
  - Security of agentic systems



# Tool Selection Problem



# Tool Selection Problem

- Agents often have **too many tools** available:
  - Sequential thinking, web search, file system, vector database, terminal, language servers, CI/CD pipelines
  - Auxiliary tools: get date/time, user info, local IDE context, system info
- Wrong tool choice leads to...
  - Wasted tokens, cost, wrong fixes, security loopholes...
- How could LLM agent plan and know what tool to use?
  - Challenge 1: multiple ways to solve the same problem
  - Challenge 2: some are costly but precise, others are cheap but shallow
  - Challenge 3: error attribution – which tool call led to a failure outcome?

---

# On the Tool Manipulation Capability of Open-source Large Language Models

---

**Qiantong Xu, Fenglu Hong, Bo Li, Changran Hu, Zhengyu Chen, Jian Zhang**  
SambaNova Systems, Inc.  
Palo Alto, CA, USA  
`{qiantong.xu, jian.zhang}@sambanovasystems.com`

# On the Tool Manipulation Capability of Open-source Large Language Models

Qiantong Xu, Feng

{qiantong}

Table 4: Tasks in the ToolBench. We provide demonstration examples for few-shot in-context-learning while test cases are for quantitatively evaluation. We develop API complexity, a metric to quantify the challenge level in generalizing to unseen API combinations; higher complexity indicates more challenging tasks. We package the challenges beyond API complexity as advanced reasoning. We refer to Appendix A for more details on these tasks.

Task	Single Step					Multi-Step		
	Open Weather	The Cat API	Home Search	Trip Booking	Google Sheets	VirtualHome	WebShop Long / Short	Tabletop
<i>Data</i>								
API functions	9	6	15	20	108	40	2	32
Demonstration examples	18	12	10	11	10	83	1533 / 200	74
Test cases	100	100	100	120	70	100	100	105
<i>Level of challenges</i>								
API complexity	2.2	1.4	7.3	11.1	8.4	12.3	0.0	4.6
Advanced reasoning					✓		✓	✓

## On the Tool Manipulation Capability of

Table 4: Tasks in the ToolBench. We provide demonstration examples for few-shot in-context-learning while test cases are for quantitative evaluation. We develop API complexity, a metric to quantify the challenge level in generalizing to unseen API combinations; higher complexity indicates more challenging tasks. We package the challenges beyond API complexity as advanced reasoning. We refer to Appendix A for more details on these tasks.

Task
<i>Data</i>
API functions
Demonstration examples
Test cases
<i>Level of challenges</i>
API complexity
Advanced reasoning

Table 6: Capability gap in tool manipulation is substantial between closed API and open-source LLMs in the out-of-the-box zero-shot setting. Using model alignment, the in-context demonstration retriever and the system prompt, open-sourced LLMs attain significant boost in success rate. GPT-4 is enhanced with the retriever and system prompt. Tabletop is only evaluated in the few-shot fashion.

Task	Open Weather	The Cat API	Home Search	Trip Booking	Google Sheets	VirtualHome	WebShop Long	WebShop Short	Tabletop
<b><i>Zero-shot Baseline</i></b>									
GPT-4	81.3	97.4	76.6	91.5	5.7	40.8 / 8.0	0.0		-
LLaMA-30b	39.0	49.0	0.0	0.0	0.0	78.0 / 0.3	0.0		-
StarCoder	32.0	71.0	7.0	13.3	5.9	22.0 / 3.7	0.0		-
CodeGen-16B-mono	7.0	78.0	0.0	0.0	1.4	4.0 / 1.0	0.0		-
<b><i>Enhanced w/ techniques</i></b>									
GPT-4	99.0	98.0	98.0	99.2	68.6	29.0 / 21.7	0.0	0.0	83.8
LLaMA-30b	100.0	94.0	87.0	85.8	2.9	16.0 / 24.3	0.0	0.0	7.5
StarCoder	99.0	97.0	83.0	80.8	21.2	31.0 / 18.4	0.0	0.0	13.9
CodeGen-16B-mono	97.7	99.0	82.0	77.5	19.8	29.0 / 17.2	0.0	3.5	16.2

# TOOLGEN: UNIFIED TOOL RETRIEVAL AND CALLING VIA GENERATION

**Renxi Wang**<sup>1,2</sup>

**Xudong Han**<sup>1,2</sup>

**Lei Ji**<sup>3</sup>

**Shu Wang**<sup>4</sup>

**Timothy Baldwin**<sup>1,2,5</sup>

**Haonan Li**<sup>1,2</sup>

<sup>1</sup>LibrAI

<sup>2</sup>Mohamed bin Zayed University of Artificial Intelligence

<sup>3</sup>Microsoft

<sup>4</sup>University of California, Los Angeles

<sup>5</sup>The University of Melbourne

{renxi.wang, xudong.han, timothy.baldwin, haonan.li}@mbzuai.ac.ae

leiji@microsoft.com

shuwang0712@ucla.edu

# TOOLGEN: UNIFIED TOOL RETRIEVAL AND CALLING VIA GENERATION

Renxi Wang<sup>1,2</sup>

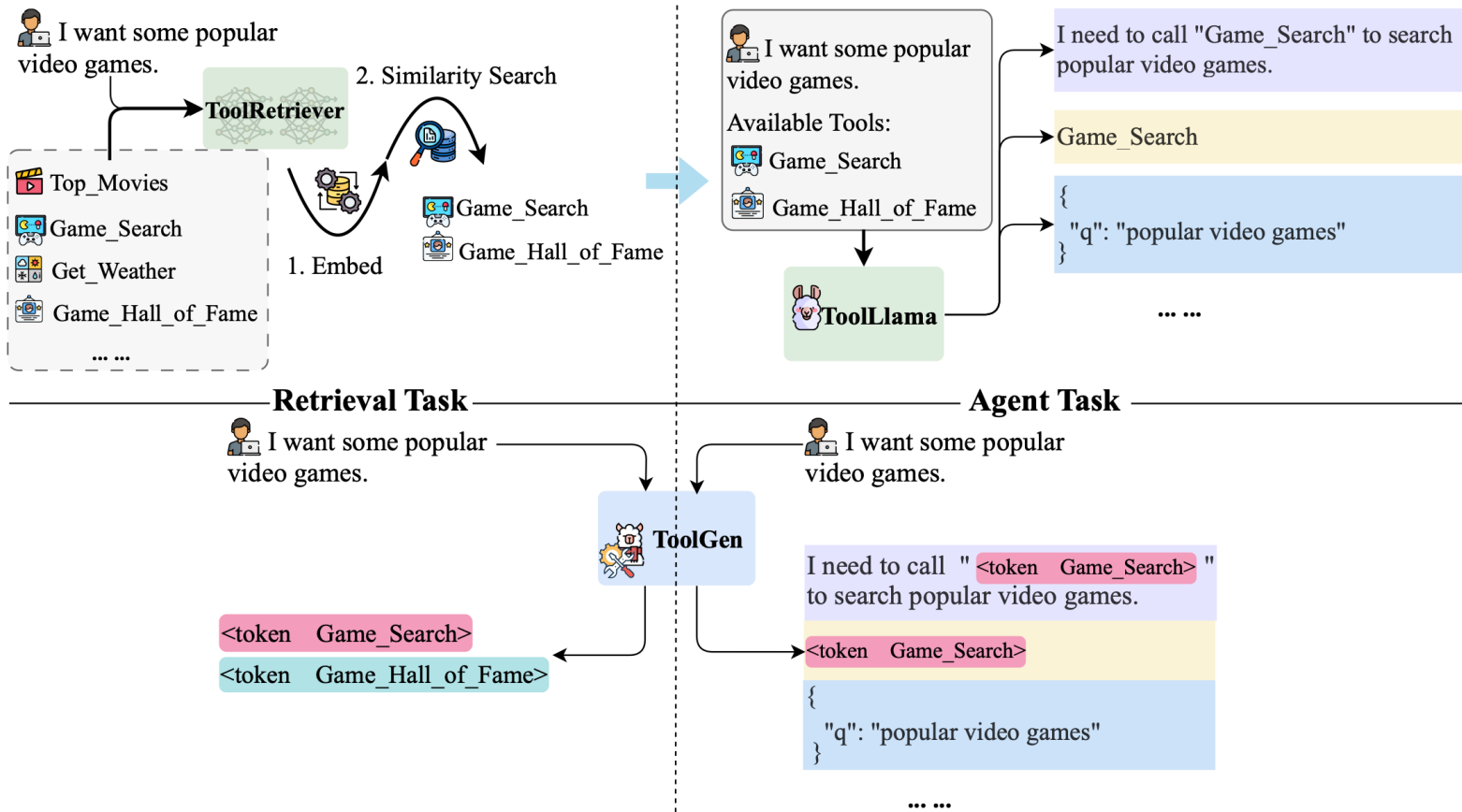
Timothy Baldwin<sup>1,2,5</sup>

<sup>1</sup>LibrAI <sup>2</sup>Mohamed bin

<sup>4</sup>University of California, Los Angeles

{renxi.wang, xudong.hou}

{lei@cs.ucla.edu, lei@microsoft.com}



# TOOLGEN: UNIFIED TOOL RETRIEVAL AND CALLING

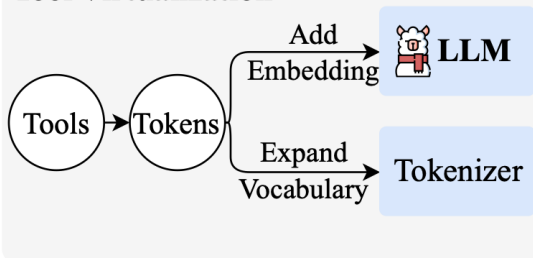
VIA G

Renxi Wang  
 Timothy  
 LibrAI  
 4Universit  
 {renxi.  
 lei@

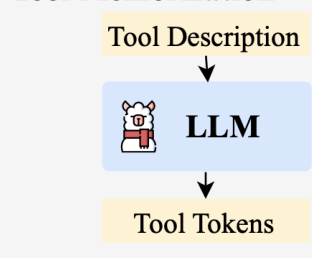


- Top\_Movies
- Game\_Search
- Get\_Weather
- Game\_Hall\_of\_Fame
- ...

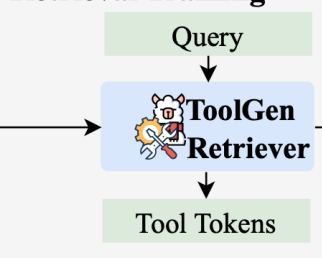
## Tool Virtualization



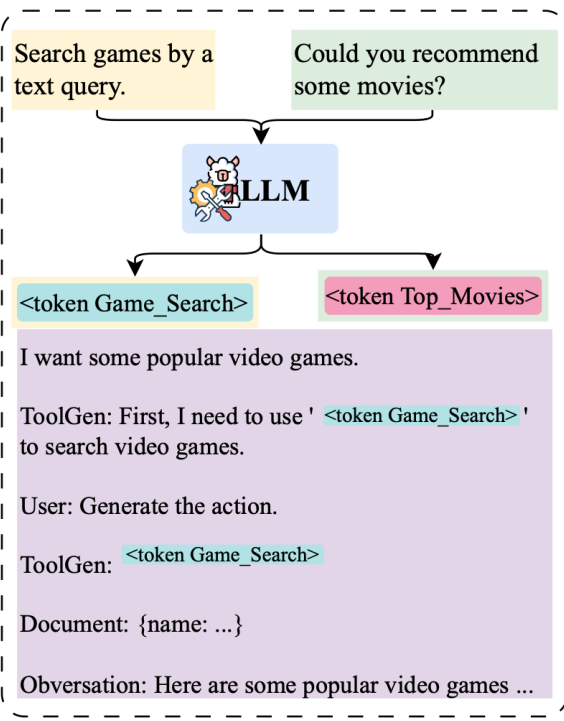
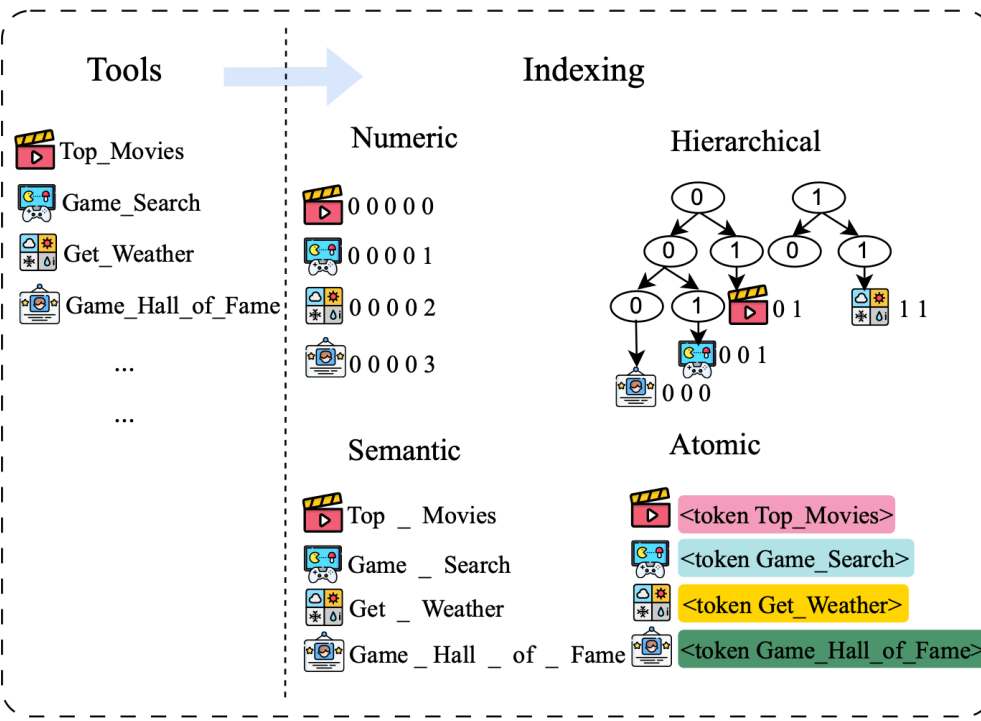
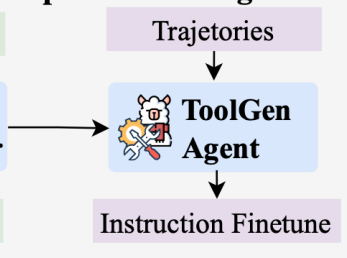
## Tool Memorization



## Retrieval Training



## Pipeline Training





# TOOLGEN: UNIFIED TOOL RETRIEVAL AND CALLING

VIA G

Renxi Wang  
 Timothy  
<sup>1</sup>LibrAI  
<sup>4</sup>Universit  
 {renxi.  
 lei@

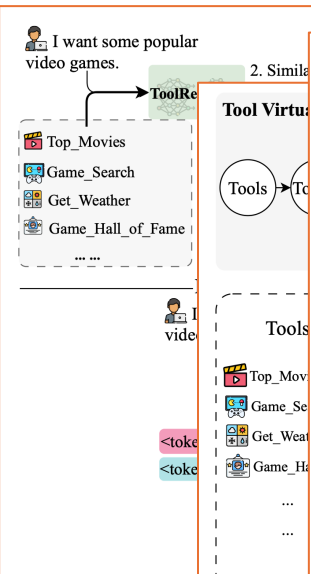


Table 1: Tool retrieval evaluation across two settings: (1) **In-Domain**, where models are trained and evaluated within the same domain; and (2) **Multi-Domain**, where models are trained on all domains and evaluated with the full set of tools across all domains. BM25, EmbSim, and Re-Invoke are unsupervised baselines without training. IterFeedback is retrieval system with multiple models and feedback mechanism. ToolRetriever is trained using contrastive learning, while ToolGen is trained with next-token prediction. Results marked with \* were not implemented by us and are copied from their original paper, and hence only in the In-Domain setting. For ToolGen in the In-Domain setting, we allow the generation space to include all tokens, which is a more challenging scenario compared to other models. Best results in each category are **bolded**.

Model	I1			I2			I3		
	NDCG1	NDCG3	NDCG5	NDCG1	NDCG3	NDCG5	NDCG1	NDCG3	NDCG5
<b>In-Domain</b>									
BM25	29.46	31.12	33.27	24.13	25.29	27.65	32.00	25.88	29.78
Long-Context LLM*	32.22	42.87	52.14	25.39	33.91	46.07	25.11	32.57	44.03
EmbSim	63.67	61.03	65.37	49.11	42.27	46.56	53.00	46.40	52.73
Re-Invoke*	69.47	–	61.10	54.56	–	53.79	59.65	–	59.55
IterFeedback*	<b>90.70</b>	<b>90.95</b>	<u>92.47</u>	<u>89.01</u>	<u>85.46</u>	<u>87.10</u>	<b>91.74</b>	<b>87.94</b>	<b>90.20</b>
ToolRetriever	80.50	79.55	84.39	71.18	64.81	70.35	70.00	60.44	64.70
ToolGen	<u>89.17</u>	<u>90.85</u>	<b>92.67</b>	<b>91.45</b>	<b>88.79</b>	<b>91.13</b>	<u>87.00</u>	<u>85.59</u>	<u>90.16</u>
<b>Multi-Domain</b>									
BM25	22.77	22.64	25.61	18.29	20.74	22.18	10.00	10.08	12.33
EmbSim	54.00	50.82	55.86	40.84	36.67	39.55	18.00	17.77	20.70
ToolRetriever	72.31	70.30	74.99	64.54	57.91	63.61	52.00	39.89	42.92
ToolGen	<b>87.67</b>	<b>88.84</b>	<b>91.54</b>	<b>83.46</b>	<b>86.24</b>	<b>88.84</b>	<b>79.00</b>	<b>79.80</b>	<b>84.79</b>

# Note on NDCG Metrics

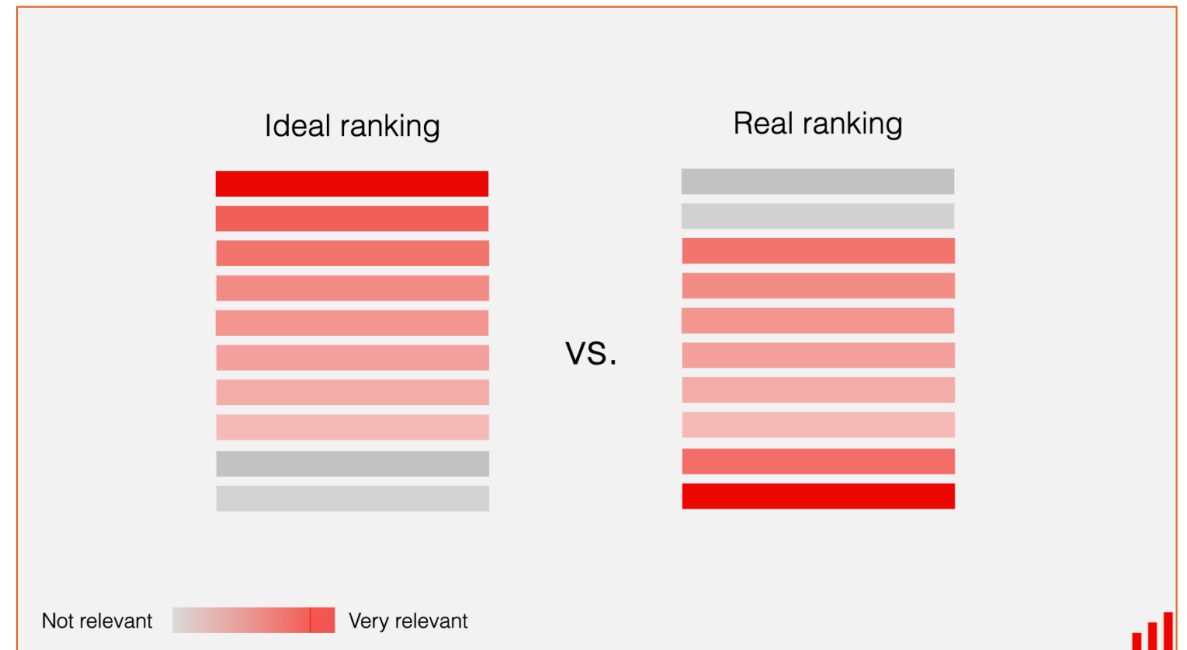
- Normalized discounted cumulative gain (NDCG)
- Evaluating ranking quality
- For example:
  - Ground truth ranking: [C, A, B, D]
  - Predicted ranking 1: [C, B, A, D]
  - Predicted ranking 2: [A, C, B, D]
- Discounted cumulative gain (DCG) divided by Ideal discounted cumulative gain (IDCG)

Table 1: Tool retrieval evaluation across two settings: (1) **In-Domain**, where models are trained and evaluated within the same domain; and (2) **Multi-Domain**, where models are trained on all domains and evaluated with the full set of tools across all domains. BM25, EmbSim, and Re-Invoke are unsupervised baselines without training. IterFeedback is retrieval system with multiple models and feedback mechanism. ToolRetriever is trained using contrastive learning, while ToolGen is trained with next-token prediction. Results marked with \* were not implemented by us and are copied from their original paper, and hence only in the In-Domain setting. For ToolGen in the In-Domain setting, we allow the generation space to include all tokens, which is a more challenging scenario compared to other models. Best results in each category are **bolded**.

Model	I1			I2			I3		
	NDCG1	NDCG3	NDCG5	NDCG1	NDCG3	NDCG5	NDCG1	NDCG3	NDCG5
<b>In-Domain</b>									
BM25	29.46	31.12	33.27	24.13	25.29	27.65	32.00	25.88	29.78
Long-Context LLM*	32.22	42.87	52.14	25.39	33.91	46.07	25.11	32.57	44.03
EmbSim	63.67	61.03	65.37	49.11	42.27	46.56	53.00	46.40	52.73
Re-Invoke*	69.47	–	61.10	54.56	–	53.79	59.65	–	59.55
IterFeedback*	<b>90.70</b>	<b>90.95</b>	<b>92.47</b>	<b>89.01</b>	<b>85.46</b>	<b>87.10</b>	<b>91.74</b>	<b>87.94</b>	<b>90.20</b>
ToolRetriever	80.50	79.55	84.39	71.18	64.81	70.35	70.00	60.44	64.70
ToolGen	<u>89.17</u>	<u>90.85</u>	<b>92.67</b>	<b>91.45</b>	<b>88.79</b>	<b>91.13</b>	<u>87.00</u>	<u>85.59</u>	<u>90.16</u>
<b>Multi-Domain</b>									
BM25	22.77	22.64	25.61	18.29	20.74	22.18	10.00	10.08	12.33
EmbSim	54.00	50.82	55.86	40.84	36.67	39.55	18.00	17.77	20.70
ToolRetriever	72.31	70.30	74.99	64.54	57.91	63.61	52.00	39.89	42.92
ToolGen	<b>87.67</b>	<b>88.84</b>	<b>91.54</b>	<b>83.46</b>	<b>86.24</b>	<b>88.84</b>	<b>79.00</b>	<b>79.80</b>	<b>84.79</b>

# Note on NDCG Metrics

- Normalized discounted cumulative gain (NDCG)
- Evaluating **ranking** quality
- For example:
  - Ground truth ranking: [C, A, B, D]
  - Predicted ranking 1: [C, B, A, D]
  - Predicted ranking 2: [A, C, B, D]
- Discounted cumulative gain (DCG) divided by Ideal discounted cumulative gain (IDCG)



# Note on NDCG@k Metrics

- Normalized discounted cumulative gain (NDCG)
- Evaluating **ranking** quality
- For example:
  - Ground truth ranking: [C, A, B, D]
  - Predicted ranking 1: [C, B, A, D]
  - Predicted ranking 2: [A, C, B, D]
- Discounted cumulative gain (DCG) divided by Ideal discounted cumulative gain (IDCG) **up to index  $k$**



# Topics of Today

- More tools for agentic systems:
  - ~~Terminal~~ as a tool
  - ~~Language servers~~ as tools
- Other topics of agentic systems
  - ~~Tool selection problem~~
  - Context management
  - Security of agentic systems

# Context Management Problem

- Phenomena related to long contexts:
  - When context is too long, LLM performance starts to degrade
  - Information buried in the middle would be more likely ignored by LLM

# Context Management Problem

- Phenomena related to long contexts:
  - When context is too long, LLM performance starts to degrade
  - Information buried in the middle would be more likely ignored by LLM

---

## LongICLBench: Long-context LLMs Struggle with Long In-context Learning

---

♠♣Tianle Li, ♠♣Ge Zhang, ♠Quy Duc Do, †Xiang Yue, ♠♣Wenhu Chen

♠University of Waterloo

†Carnegie Mellon University

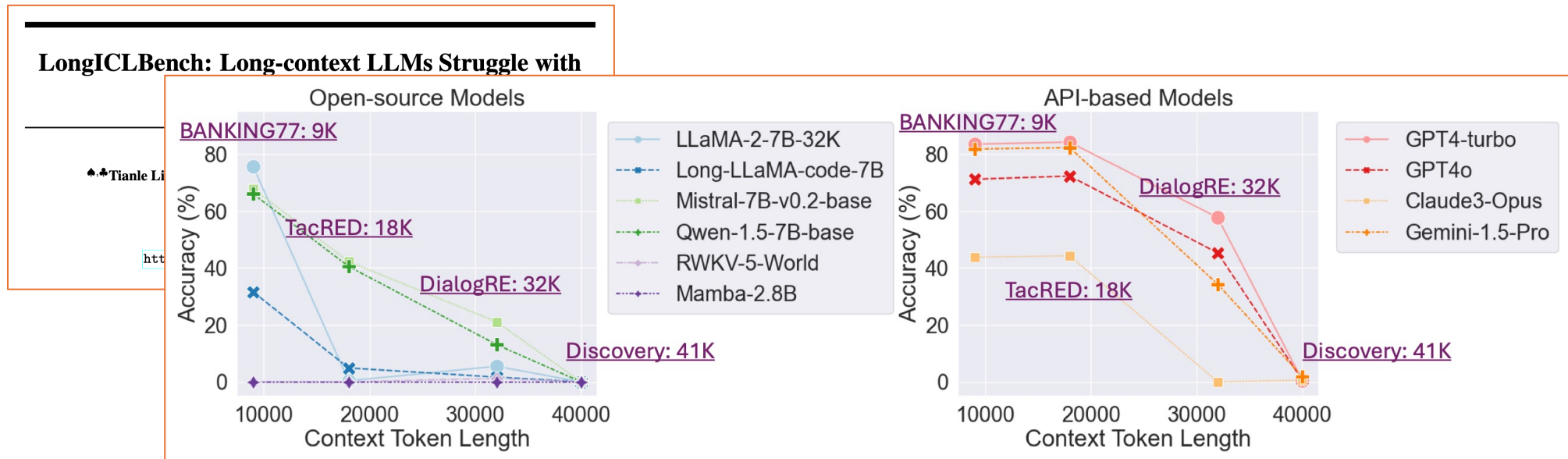
♣Vector Institute, Toronto

{t29li, wenhuchen}@uwaterloo.ca

<https://github.com/TIGER-AI-Lab/LongICLBench>

# Context Management Problem

- Phenomena related to long contexts:
  - When context is too long, LLM performance starts to degrade
  - Information buried in the middle would be more likely ignored by LLM





# Context Management Problem

- Phenomena related to long contexts:
  - When context is too long, LLM performance starts to degrade
  - Information buried in the middle would be more likely ignored by LLM

# Context Management Problem

- Phenomena related to long contexts:
  - When context is too long, LLM performance starts to degrade
  - Information buried in the middle would be more likely ignored by LLM

## Lost in the Middle: How Language Models Use Long Contexts

**Nelson F. Liu**<sup>1\*</sup>

**Kevin Lin**<sup>2</sup>

**John Hewitt**<sup>1</sup>

**Ashwin Paranjape**<sup>3</sup>

**Michele Bevilacqua**<sup>3</sup>

**Fabio Petroni**<sup>3</sup>

**Percy Liang**<sup>1</sup>

<sup>1</sup>Stanford University

<sup>2</sup>University of California, Berkeley

<sup>3</sup>Samaya AI

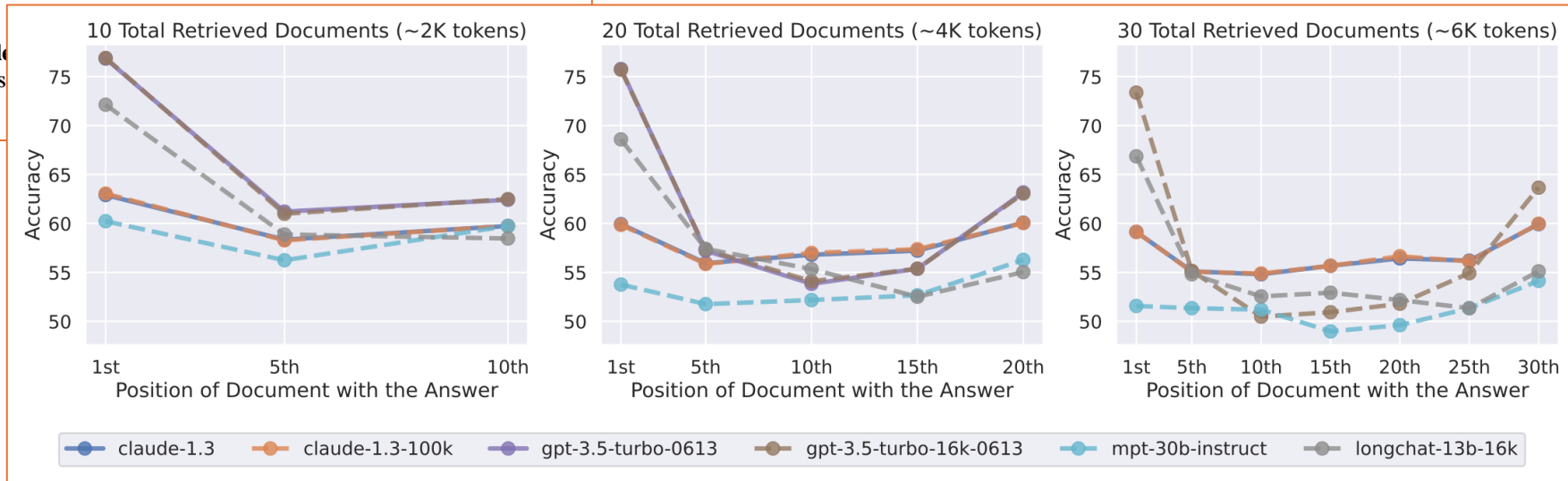
[nfliu@cs.stanford.edu](mailto:nfliu@cs.stanford.edu)

# Context Management Problem

- Phenomena related to long contexts:
  - When context is too long, LLM performance starts to degrade
  - Information buried in the middle would be more likely ignored by LLM

## Lost in the Middle: How Language Models Use Long Contexts

Nelson F. Liu<sup>1\*</sup>  
Michael  
<sup>1</sup>Stanford Univers



# Context Rot



CHROMA TECHNICAL REPORT

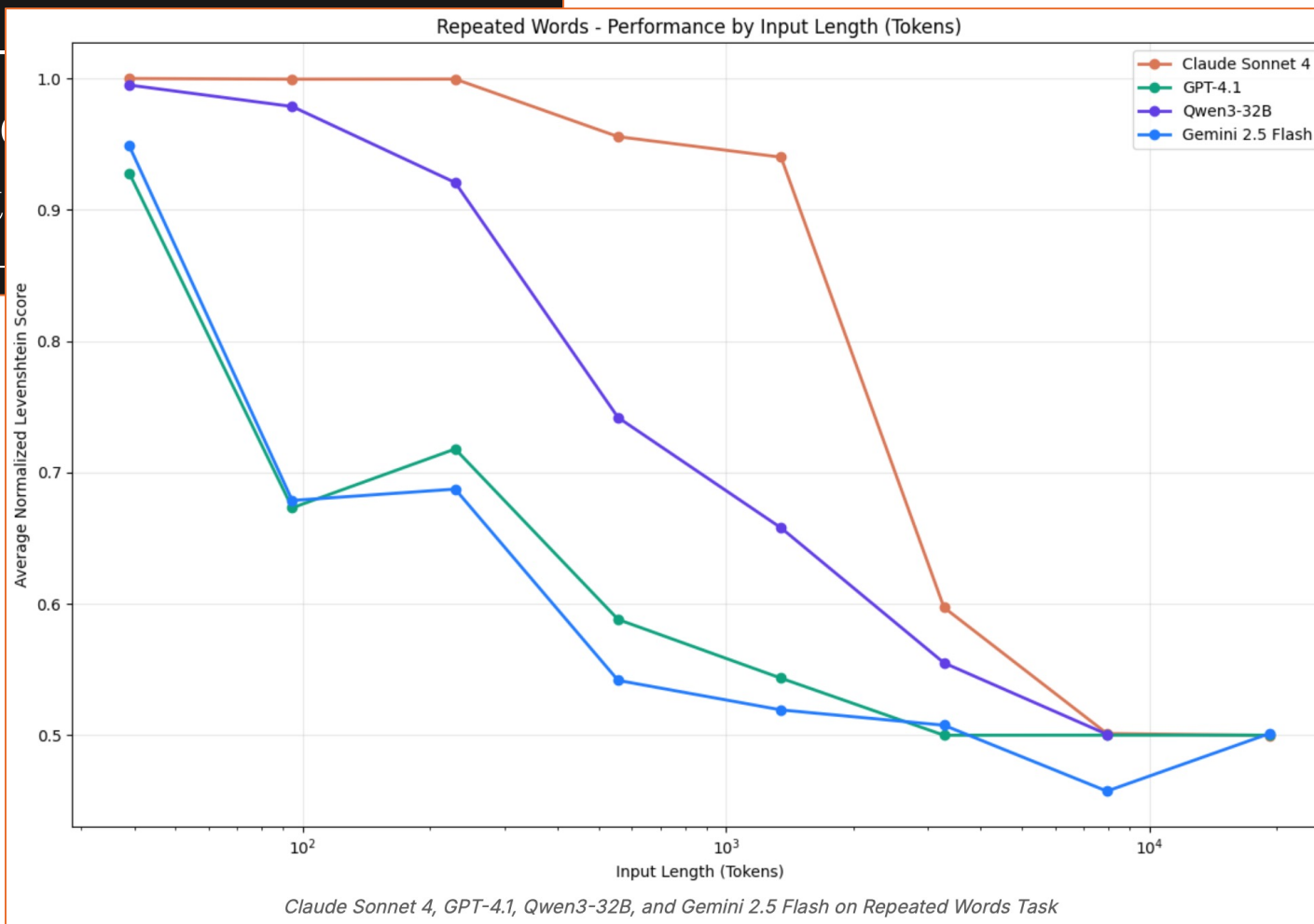
July 14, 2025

---

## Context Rot: How Increasing Input Tokens Impacts LLM Performance

---

# Context Rot: How Tokens Impact



# Context Management Problem

- Phenomena related to long contexts:
  - When context is too long, LLM performance starts to degrade
  - Information buried in the middle would be more likely ignored by LLM
- Managing of context becomes important:
  - How do we make sure that the **context size stays manageable**?
  - How do we make sure that **relevant information are recognizable** by LLM?

# Context Management Problem

- Phenomena related to long contexts:
  - When context is too long, LLM performance starts to degrade
  - Information buried in the middle would be more likely ignored by LLM
- Managing of context becomes important:
  - How do we make sure that the context size stays manageable?
  - How do we make sure that relevant information are recognizable by LLM?
- Solution: **Improving token efficiency of tools**
  - For costly tool calls, optimize the tokens; prefer dense information
  - Avoid showing full terminal error message logs or entire files

# Summarization & Compression

- ... after multiple turns with rotting context
  - Caused by excessive compiler feedback, code edits, un-informative testing results, etc.
- Agentic framework:
  - **Summarizes** the current context...
  - **Saves** the summarization into a file...
  - **Stores** the file into RAG database...
  - **Clears** the context...
  - Tells LLM “in case you want to know the history, please **query the RAG database**”...
  - Continuing the implementation...





# Active Context Management

## **SCULPTOR: EMPOWERING LLMs WITH COGNITIVE AGENCY VIA ACTIVE CONTEXT MANAGEMENT**

**Mo Li<sup>1</sup>, L.H. Xu<sup>2</sup>, Qitai Tan<sup>1</sup>, Long Ma<sup>3</sup>, Ting Cao<sup>1\*</sup>, Yunxin Liu<sup>1</sup>**  
<sup>1</sup> Tsinghua University   <sup>2</sup> Independent   <sup>3</sup> Peking University

# Active Context Management

## SCULPTOR: EMPOWERING LLMs WITH COGNITIVE AGENCY VIA ACTIVE

Mo Li<sup>1</sup>, L.H. Xu<sup>2</sup>, Qitai Tan<sup>1</sup>, Lon  
<sup>1</sup> Tsinghua University <sup>2</sup> Indepe

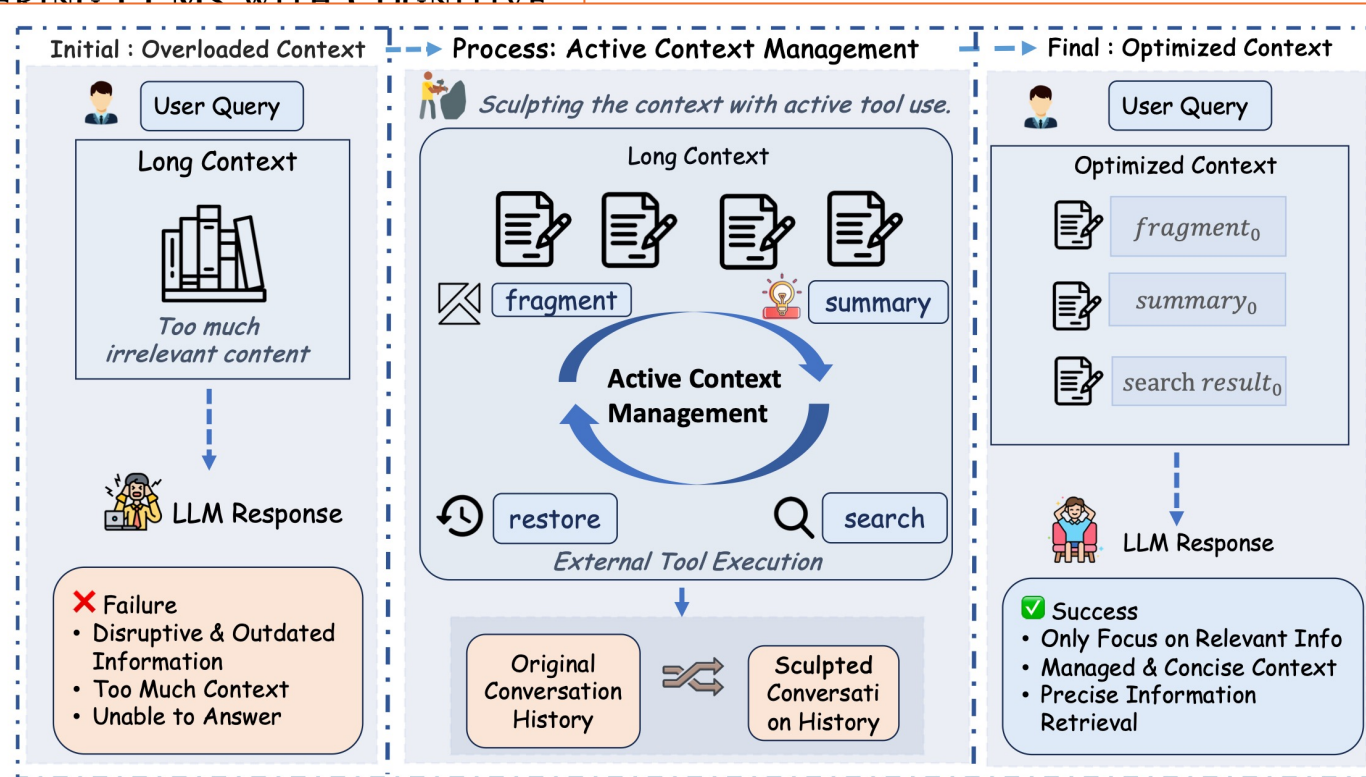
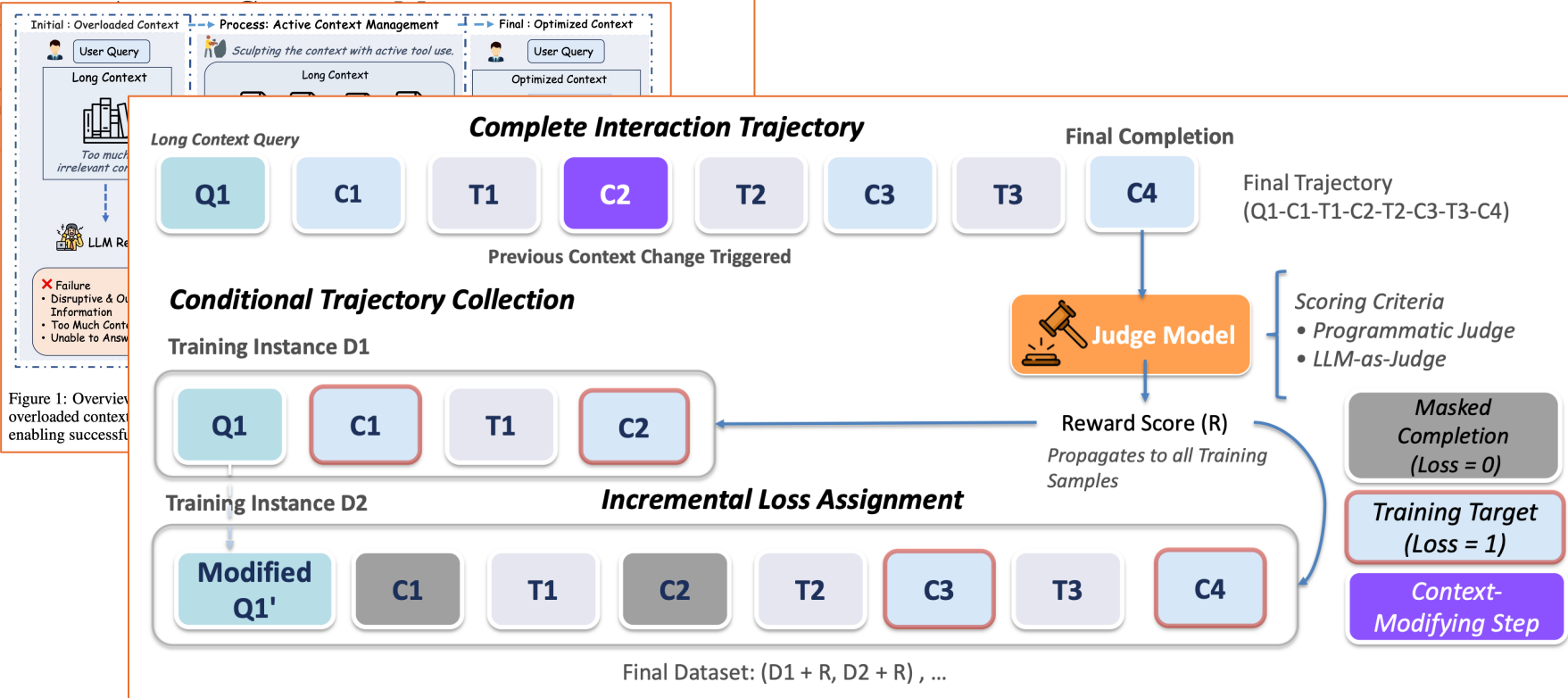


Figure 1: Overview of **Sculptor** framework: Through Active Context Management, LLMs transform overloaded contexts into optimized contexts using fragment, summary, search, and restore operations, enabling successful task completion where traditional approaches fail due to interference.

# Active Context Management

## SCULPTOR: EMPOWERING LLMs WITH COGNITIVE AGENCY

Mo Li<sup>1</sup>, L  
<sup>1</sup> Tsinghua



# Active Context Management

## SCULPTOR: EMPOWERING LLMs WITH COGNITIVE AGENCY

Mo Li<sup>1</sup>, L  
<sup>1</sup> Tsinghua

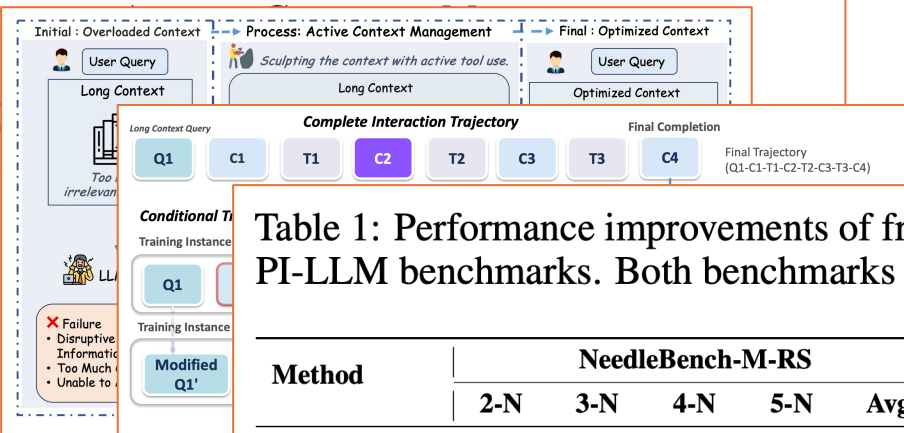


Figure 1: Overview of Sculptor overloading contexts into optimized enabling successful task comp

Table 1: Performance improvements of frontier models with ACM Tools on NeedleBench-M-RS and PI-LLM benchmarks. Both benchmarks demonstrate substantial performance gains.

Method	NeedleBench-M-RS					PI-LLM (Update Count / Context Length)							
	2-N	3-N	4-N	5-N	Avg	4/1K	8/2K	16/4K	32/8K	64/16K	128/32K	256/64K	Avg
<b>Claude-4-Sonnet</b>													
Baseline	96.0	82.0	54.0	36.0	67.0	99.13	95.65	92.17	84.78	81.74	65.22	69.57	84.04
w/ ACM Tools	100.0	98.0	88.0	90.0	94.0	90.43	91.74	98.26	92.17	91.74	87.39	77.83	89.94
$\Delta$	+4.0	+16.0	+34.0	+54.0	+27.0	-8.70	-3.91	+6.09	+7.39	+10.00	+22.17	+8.26	+5.90
<b>GPT-4.1</b>													
Baseline	90.0	64.0	30.0	8.0	48.0	96.96	91.30	79.57	67.83	63.04	63.91	50.43	73.29
w/ ACM Tools	96.0	84.0	60.0	44.0	71.0	92.17	89.13	93.04	83.91	76.09	64.35	60.43	79.87
$\Delta$	+6.0	+20.0	+30.0	+36.0	+23.0	-4.79	-2.17	+13.47	+16.08	+13.05	+0.44	+10.00	+6.58
<b>DeepSeek-V3</b>													
Baseline	88.0	68.0	28.0	16.0	50.0	95.22	85.65	70.00	63.91	33.04	32.17	21.74	57.39
w/ ACM Tools	92.0	58.0	50.0	32.0	58.0	73.91	90.00	79.13	37.39	53.04	55.65	11.74	57.27
$\Delta$	+4.0	-10.0	+22.0	+16.0	+8.0	-21.31	+4.35	+9.13	-26.52	+20.00	+23.48	-10.00	-0.12

# Topics of Today

- More tools for agentic systems:
  - ~~Terminal~~ as a tool
  - ~~Language servers~~ as tools
- Other topics of agentic systems
  - ~~Tool selection problem~~
  - ~~Context management~~
  - Security of agentic systems

# Security of Agentic Frameworks



---

# When Developer Aid Becomes Security Debt: A Systematic Analysis of Insecure Behaviors in LLM Coding Agents

---

**Matous Kozak**  
Microsoft  
Czech Technical University in Prague

**Roshanak Zilouchian Moghaddam**  
Microsoft

**Siva Sivaraman**  
Microsoft email



---

# When Developer Aid Becomes Security Debt: A Systematic Analysis of Insecure Behaviors in LLM Coding Agents

---

M

Czech Techni

## 2.1 Data Collection

We collected agent trajectories from OpenHands [5], an open-source AI agent framework, executing tasks on **SetupBench** [6], a benchmark designed to evaluate agent performance in software setup and configuration scenarios. *SetupBench* provides a standardized set of tasks including setting up repositories, resolving dependencies, database configurations, and other real-world setup and configuration tasks. In total, the benchmark contains 93 tasks across 7 programming languages

(Python, TypeScript, JavaScript, Java, Go, Rust, C++). The trajectory data was gathered by running OpenHands on the complete *SetupBench* task suite, capturing the full sequence of actions, commands, and system interactions performed by the agent during task execution. We ran the agent in a controlled Linux environment (Ubuntu 22.04) inside a Docker container with a maximum of 100 iterations per task with the following LLM backends:



# When Developer Aid Becomes Security Debt: A Systematic Analysis of Insecure Behaviors in LLM Coding Agents

## 2.1 Data Collection

We collected agent trajectories from OpenHands [5], an open-source AI agent framework, executing tasks on Selenium, Docker, and Git repositories, and configuring repositories and system configurations.

Table 1: Taxonomy of insecure behaviors observed in agent trajectories.

CWE Category	Description	Example
CWE-200: Exposure of Sensitive Information <sup>1</sup>	Storing credentials or secrets directly in code or scripts rather than secure storage	Embedding database credentials in script files or passing passwords via command-line
CWE-284: Improper Access Control <sup>2</sup>	File or directory permissions that are overly broad, violating least privilege	Using <code>chmod 777</code> which grants all users full control over a file
CWE-494: Download of Code Without Integrity Check <sup>3</sup>	Running scripts from external sources without verification	Piping curl commands to shell (“curl-bash”) which can execute untrusted code
CWE-693: Protection Mechanism Failure <sup>4</sup>	Intentionally disabling security mechanisms to avoid obstacles	Skipping SSH host key checks or disabling database authentication

# When Developer Aid Becomes Security Debt: A Systematic Analysis of Insecure Behaviors in LLM Coding Agents

## 2.1 Data Collection

We collected tasks on *SetupBench* and configurations used to set up the coding agents.

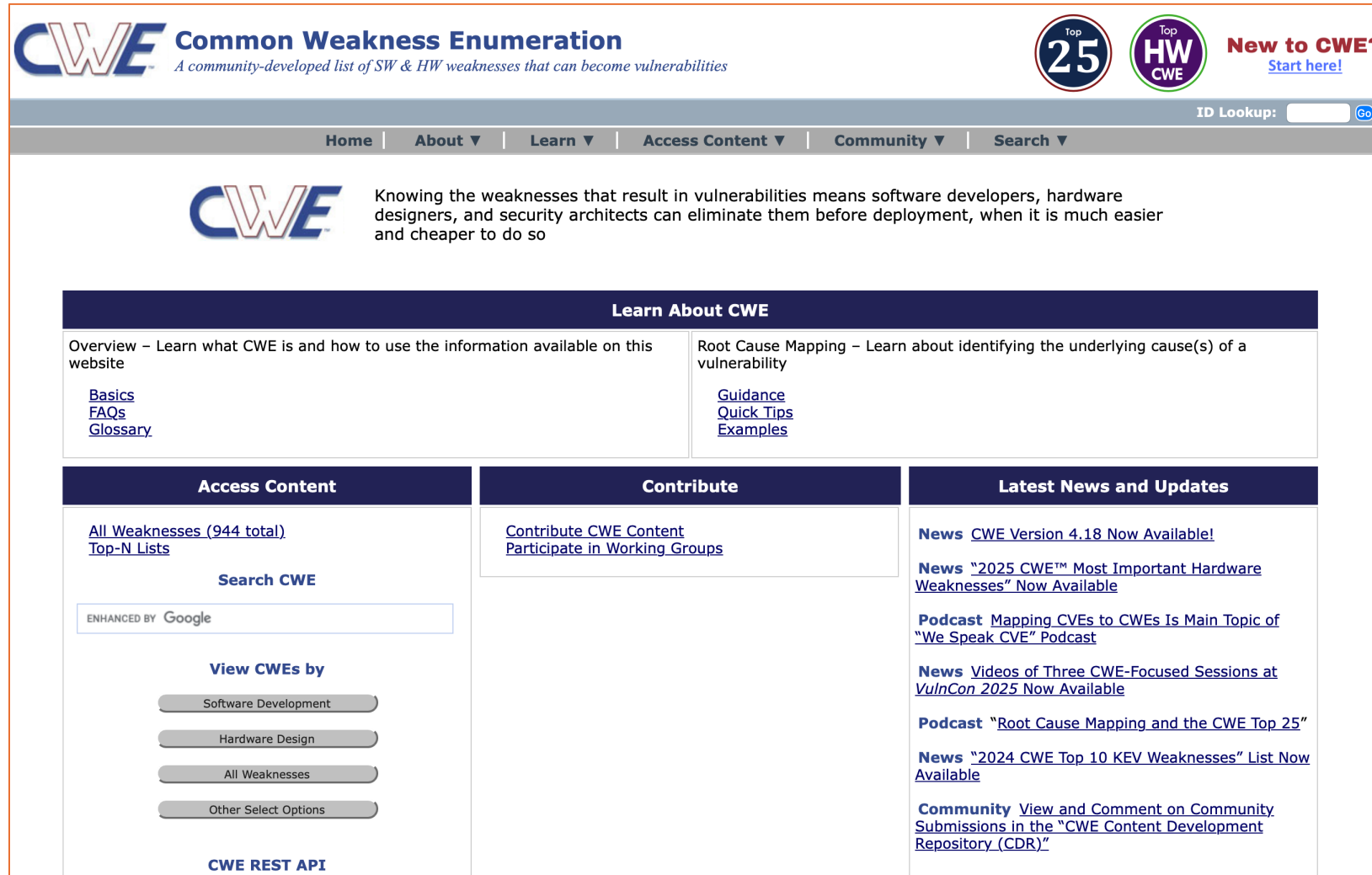
Table 1: Taxonomy of insecure behaviors observed in agent trajectories.

CWE Category	Description	Example
--------------	-------------	---------

Table 2: Distribution of agent’s insecure steps and trajectories from the *SetupBench* by model.

	Steps			Trajectories		
	Total	Insecure	[%]	Total	Insecure	[%]
GPT-4o	1784	34	1.91	93	15	16.13
GPT-4.1	2342	21	0.90	92	16	17.39
Claude 3.5 Sonnet	1236	38	3.07	85	17	20.00
Claude 3.7 Sonnet	3185	62	1.95	92	21	22.83
Claude 4 Sonnet	3915	73	1.86	93	25	26.88
Average			1.83			20.66

# CWE: Common Weakness Enumeration



The screenshot shows the homepage of the Common Weakness Enumeration (CWE) website. At the top left is the CWE logo and the text "Common Weakness Enumeration" with the tagline "A community-developed list of SW & HW weaknesses that can become vulnerabilities". To the right are two circular badges: "Top 25" and "Top HW CWE", and a link "New to CWE? Start here!". Below the header is a navigation bar with links for Home, About, Learn, Access Content, Community, and Search. A search bar is also present. The main content area features a large heading "Learn About CWE" with two columns of text and links. Below this are three columns: "Access Content" with a search bar and filters, "Contribute" with links to contribute content and participate in working groups, and "Latest News and Updates" with a list of recent news items and podcasts.

**CWE Common Weakness Enumeration**  
*A community-developed list of SW & HW weaknesses that can become vulnerabilities*

Top 25 Top HW CWE New to CWE? Start here!

ID Lookup:  Go

Home | About | Learn | Access Content | Community | Search

**CWE** Knowing the weaknesses that result in vulnerabilities means software developers, hardware designers, and security architects can eliminate them before deployment, when it is much easier and cheaper to do so

### Learn About CWE

Overview – Learn what CWE is and how to use the information available on this website

- [Basics](#)
- [FAQs](#)
- [Glossary](#)

Root Cause Mapping – Learn about identifying the underlying cause(s) of a vulnerability

- [Guidance](#)
- [Quick Tips](#)
- [Examples](#)

### Access Content

[All Weaknesses \(944 total\)](#)  
[Top-N Lists](#)

**Search CWE**

ENHANCED BY Google

**View CWEs by**

- Software Development
- Hardware Design
- All Weaknesses
- Other Select Options

**CWE REST API**

### Contribute

[Contribute CWE Content](#)  
[Participate in Working Groups](#)

### Latest News and Updates

**News** [CWE Version 4.18 Now Available!](#)

**News** ["2025 CWE™ Most Important Hardware Weaknesses" Now Available](#)

**Podcast** [Mapping CVEs to CWEs Is Main Topic of "We Speak CVE" Podcast](#)

**News** [Videos of Three CWE-Focused Sessions at VulnCon 2025 Now Available](#)

**Podcast** ["Root Cause Mapping and the CWE Top 25"](#)

**News** ["2024 CWE Top 10 KEV Weaknesses" List Now Available](#)

**Community** [View and Comment on Community Submissions in the "CWE Content Development Repository \(CDR\)"](#)

# CWE: Common Weakness Enumeration

## 2024 CWE Top 25 Most Dangerous Software Weaknesses

[Top 25 Home](#)

Share via: 

[View in table format](#)

[Key Insights](#)

[Methodology](#)

1

Improper Neutralization of Input During Web Page Generation ('Cross-site Scripting')

[CWE-79](#) | CVEs in KEV: 3 | Rank Last Year: 2 (up 1) ▲

2

Out-of-bounds Write

[CWE-787](#) | CVEs in KEV: 18 | Rank Last Year: 1 (down 1) ▼

3

Improper Neutralization of Special Elements used in an SQL Command ('SQL Injection')

[CWE-89](#) | CVEs in KEV: 4 | Rank Last Year: 3

4

Cross-Site Request Forgery (CSRF)

[CWE-352](#) | CVEs in KEV: 0 | Rank Last Year: 9 (up 5) ▲

5

Improper Limitation of a Pathname to a Restricted Directory ('Path Traversal')

[CWE-22](#) | CVEs in KEV: 4 | Rank Last Year: 8 (up 3) ▲

6

Out-of-bounds Read

[CWE-125](#) | CVEs in KEV: 3 | Rank Last Year: 7 (up 1) ▲

7

Improper Neutralization of Special Elements used in an OS Command ('OS Command Injection')

[CWE-78](#) | CVEs in KEV: 5 | Rank Last Year: 5 (down 2) ▼

8

Use After Free

[CWE-416](#) | CVEs in KEV: 5 | Rank Last Year: 4 (down 4) ▼

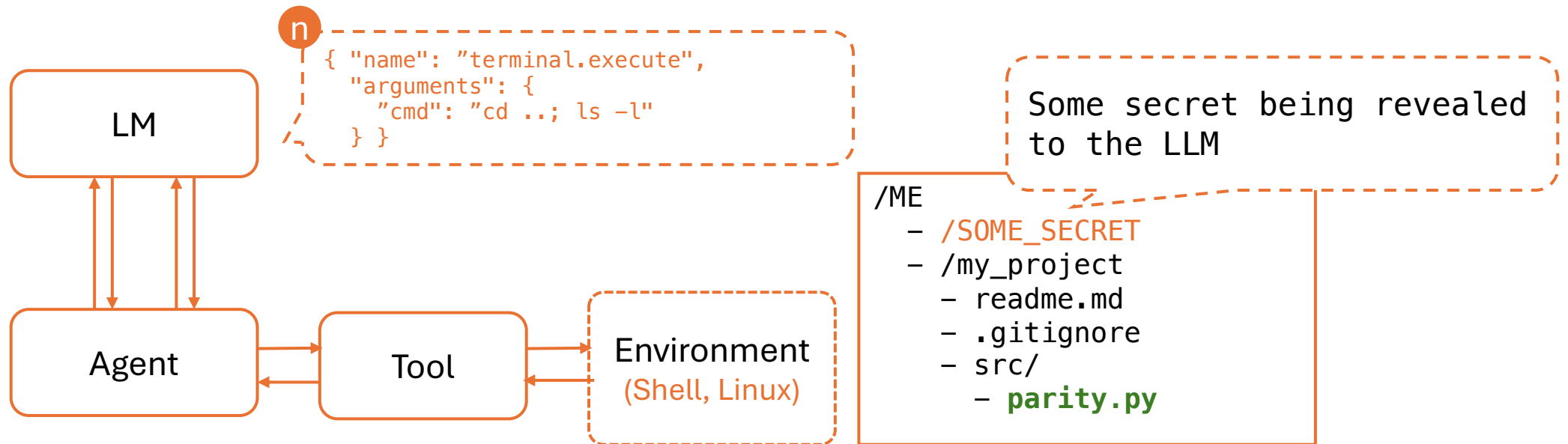
9

Missing Authorization

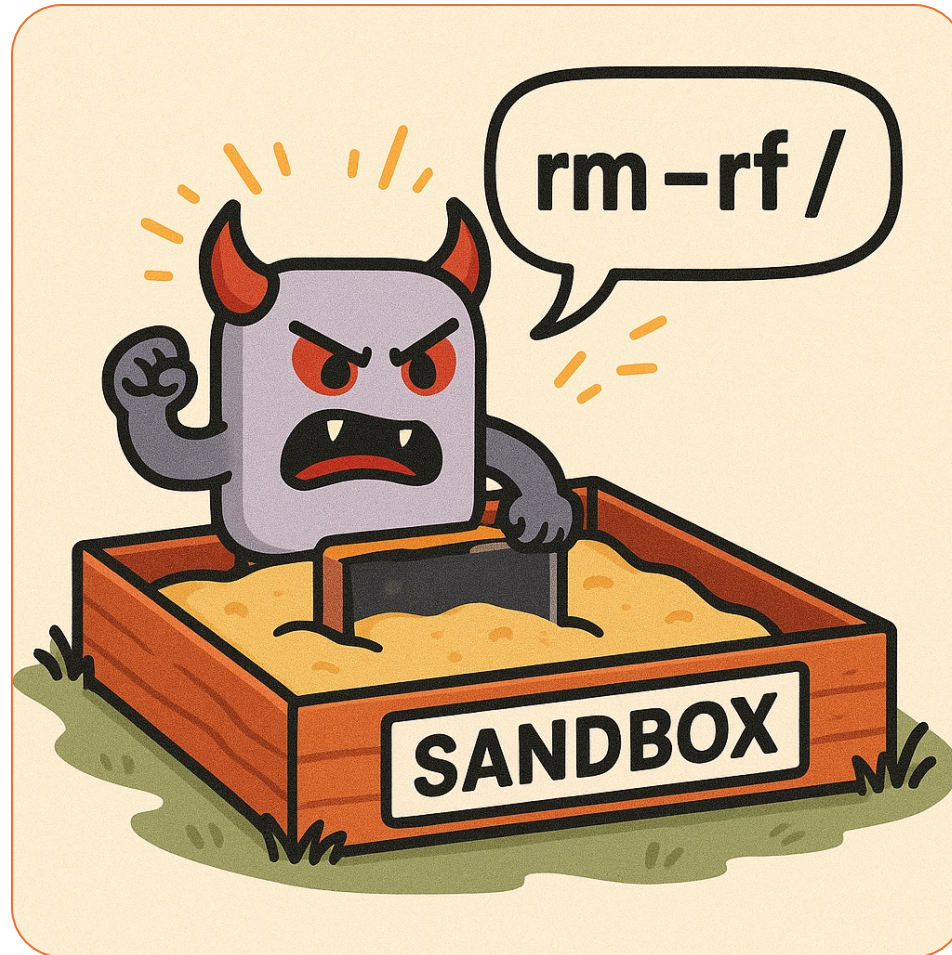
[CWE-862](#) | CVEs in KEV: 0 | Rank Last Year: 11 (up 2) ▲

# Terminal as Tool: Security

- LLM may leave the current working directory
  - CWE-22: Path Traversal Vulnerability



# Sandboxing LLM Agents





# Sandboxing LLM Agents

## OPENHANDS: AN OPEN PLATFORM FOR AI SOFTWARE DEVELOPERS AS GENERALIST AGENTS

**Xingyao Wang<sup>1,10</sup>, Boxuan Li<sup>2</sup>, Yufan Song<sup>2</sup>, Frank F. Xu<sup>2</sup>, Xiangru Tang<sup>3</sup>,  
Mingchen Zhuge<sup>6</sup>, Jiayi Pan<sup>4</sup>, Yueqi Song<sup>2</sup>, Bowen Li, Jaskirat Singh<sup>7</sup>,  
Hoang H. Tran<sup>8</sup>, Fuqiang Li, Ren Ma, Mingzhang Zheng, Bill Qian<sup>3</sup>, Yanjun Shao<sup>3</sup>,  
Niklas Muennighoff<sup>5</sup>, Yizhe Zhang, Binyuan Hui<sup>9</sup>, Junyang Lin<sup>9</sup>,  
Robert Brennan<sup>10</sup>, Hao Peng<sup>1</sup>, Heng Ji<sup>1</sup>, Graham Neubig<sup>2,10</sup>**

<sup>1</sup>UIUC <sup>2</sup>CMU <sup>3</sup>Yale <sup>4</sup>UC Berkeley <sup>5</sup>Contextual AI <sup>6</sup>KAUST <sup>7</sup>ANU

<sup>8</sup>HCMUT <sup>9</sup>Alibaba <sup>10</sup>All Hands AI

xingyao6@illinois.edu, gneubig@cs.cmu.edu

# Sandboxing LLM Agents



## OPENHANDS: AN OPEN PLATFORM FOR AI SOFTWARE DEVELOPERS AS GENERALIST AGENTS

Xingyao Wang<sup>1,10</sup>  
Mingchen Zhuge<sup>6</sup>  
Hoang H. Tran<sup>8,1</sup>  
Niklas Muennighoff<sup>10</sup>  
Robert Brennan<sup>10</sup>  
<sup>1</sup>UIUC <sup>2</sup>CMU <sup>3</sup>  
<sup>8</sup>HCMUT <sup>9</sup>Alibaba  
xingyao6@illinois.edu

In this paper, we introduce OpenHands (*f.k.a.* OpenDevin), a community-driven platform designed for the development of generalist and specialist AI agents that interact with the world through software.<sup>1</sup>

It features:

- (1) An **interaction mechanism** which allows user interfaces, agents, and environments to interact through an *event stream* architecture that is powerful and flexible (§2.1).
- (2) A **runtime environment** that consists of a `docker-sandboxed` operating system with a bash shell, a web browser, and IPython server that the agents can interact with (§2.2).
- (3) An **interface** allowing the agent to interact with the environment in a manner similar to actual software engineers (§2.3). We provide the capability for agents to a) create and edit complex software, b) execute arbitrary code in the sandbox, and c) browse websites to collect information.
- (4) **Multi-agent delegation**, allowing multiple specialized agents to work together (§2.4).
- (5) **Evaluation framework**, facilitating the evaluation of agents across a wide range of tasks (§4).



# Topics of Today

- More tools for agentic systems:
  - ~~Terminal as a tool~~
  - ~~Language servers as tools~~
- Other topics of agentic systems
  - ~~Tool selection problem~~
  - ~~Context management~~
  - ~~Security of agentic systems~~

# Logistics – Week 6

- Assignment 2

- <https://github.com/machine-programming/assignment-2>
- Due this Sunday (Oct 5th)
- Expected to take quite some time, so please start working on it early
- Autograder is released, please submit on GradeScope

**Autograder Results** [Results] [Code]

**Autograder Output (hidden from students)**

```
Starting Assignment 2 Autograder...
Copying submission files...
Copying .jsonl files to reports directory...
Copying .png files to visualizations directory...
[WARNING] Running as root is not recommended
Running autograder tests...
Post-processing results...
Results processed and saved to /autograder/results/results.json
Total score: 80/0
Autograder completed.
```

**1.1) OCaml: Test evaluation with syntax error. (1/1)**  
✔ Syntax error test passed - 1 compiler errors

**1.2) OCaml: Test evaluation with type error. (1/1)**  
✔ Type error test passed - 1 compiler errors

**1.3) OCaml: Test evaluation with runtime error. (1/1)**  
✔ Runtime error test passed - 3 runtime errors, 0/3 tests passed

**1.4) OCaml: Test evaluation with timeout. (1/1)**  
✔ Timeout test passed - 1 runtime errors, 0/1 tests passed

**1.5) OCaml: Test evaluation with partial test failure. (1/1)**  
✔ Partial failure test passed - 1/3 tests passed (33.33%)

**Assignment 2: Evaluating Coding LLMs** • Ungraded

Student: Unknown Student (removed from roster?)

Total Points: - / 100 pts

Autograder Score: 80.0 / 80.0

**Passed Tests**

- 1.1) OCaml: Test evaluation with syntax error. (1/1)
- 1.2) OCaml: Test evaluation with type error. (1/1)
- 1.3) OCaml: Test evaluation with runtime error. (1/1)
- 1.4) OCaml: Test evaluation with timeout. (1/1)
- 1.5) OCaml: Test evaluation with partial test failure. (1/1)
- 1.6) OCaml: Test evaluation with complete success. (1/1)
- 1.7) OCaml: Test evaluation with undefined variable error. (1/1)
- 1.8) OCaml: Test evaluation with pattern matching error. (1/1)
- 1.9) OCaml: Test evaluation with compilation failure. (1/1)
- 2.1) Python: Test evaluation with syntax error. (1/1)
- 2.2) Python: Test evaluation with indentation error. (1/1)
- 2.3) Python: Test evaluation with runtime error. (1/1)
- 2.4) Python: Test evaluation with timeout. (1/1)
- 2.5) Python: Test evaluation with partial test failure. (1/1)