# GPTracker:
# A Large-Scale Measurement of Misuse and Knowledge File Leakage in GPTs

IEEE S&P 2025 + ACL 2025
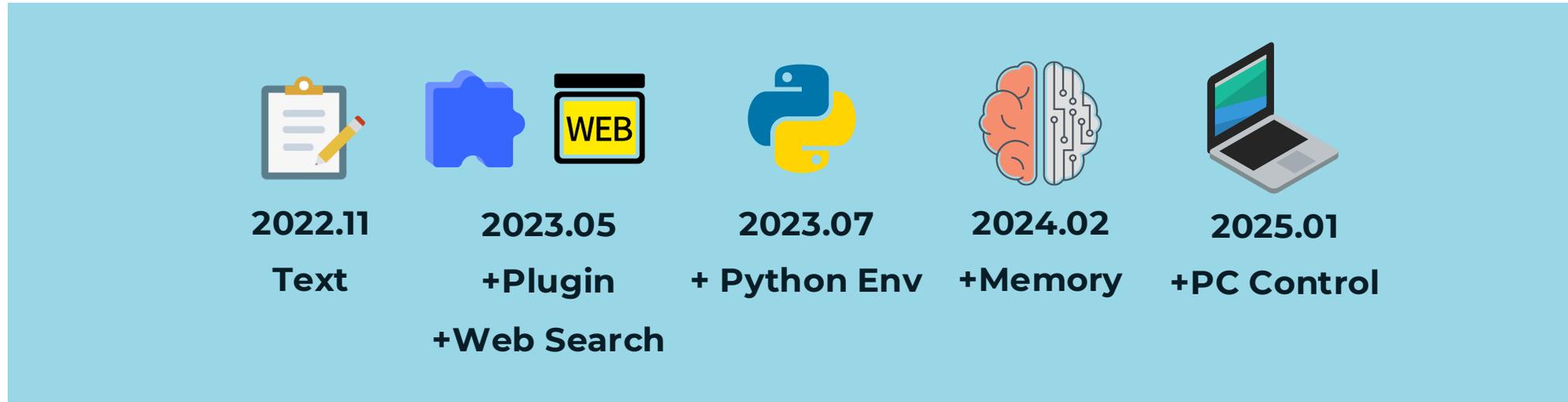
**Xinyue Shen**

CISPA Helmholtz Center for Information Security, Germany

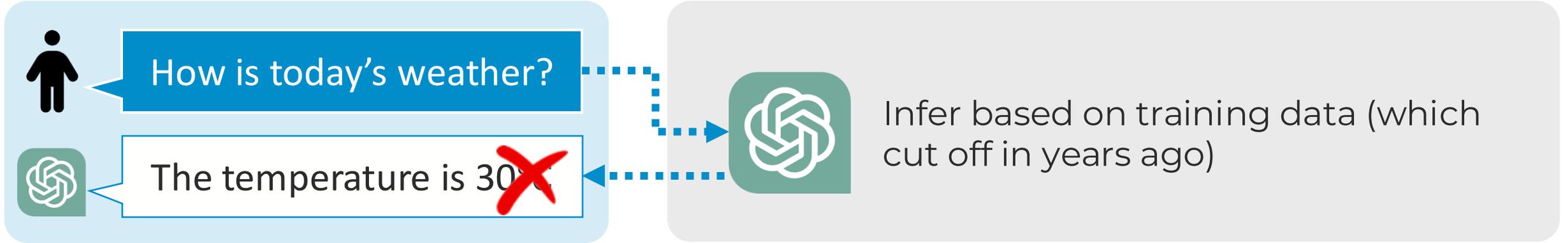**Content Warning:** examples of harmful language and unsafe images
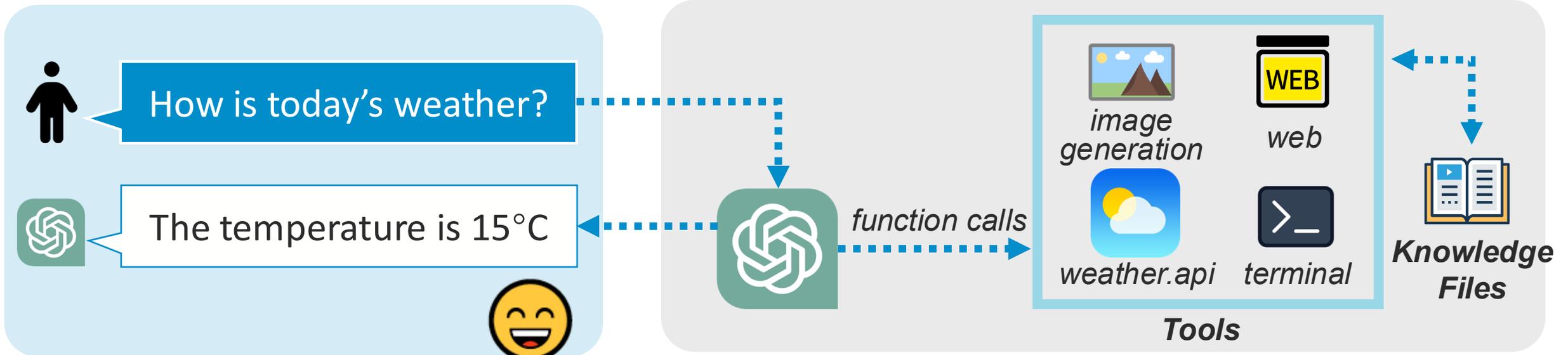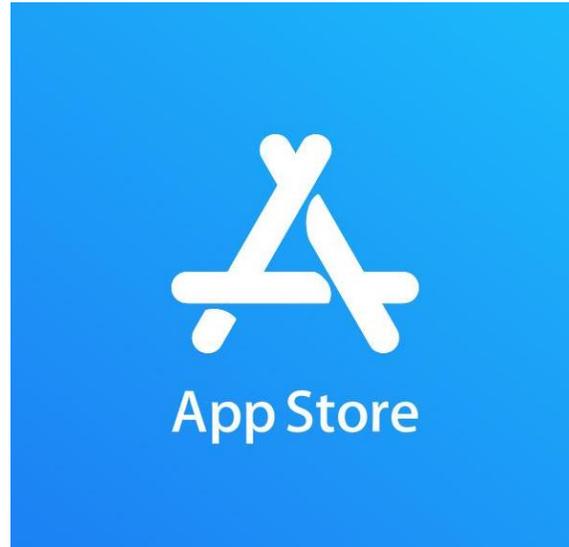
# The Rise of LLM Agents



| 2022.11 | 2023.05 | 2023.07 | 2024.02 | 2025.01 |
|---------|---------|---------|---------|---------|
| Text | +Plugin | + Python Env | +Memory | +PC Control |
| | +Web Search | | | |

**LLM Agent (LLM App)**

# LLM

How is today's weather?

The temperature is 30 ✗

Infer based on training data (which cut off in years ago)

# LLM Agent

How is today's weather?

The temperature is 15°C

function calls

**Tools**

image generation

web

weather.api

terminal

*Knowledge Files*

**Now everyone can easily develop an LLM agent :D**

# The Rise of LLM App Store

WIRED

https://www.wired.com › openai-gpt-store · 翻译此页

**OpenAI's GPT Store Has Left Some Developers in the Lurch**

2024年10月11日 — OpenAI's GPT Store has been a mixed bag. These developers say that OpenAI's analytics tools are lacking and that they have no real sense of how their GPTs are ...

TechCrunch

**OpenAI's chatbot store is filling up with spam**

The GPT Store, OpenAI's official marketplace for GPTs, is floode

potentially copyright-infringing GPTs that imply a light touch whe

20 Mar 2024

the **decoder**          DE

AI in practice          Mar 22, 2024          Update

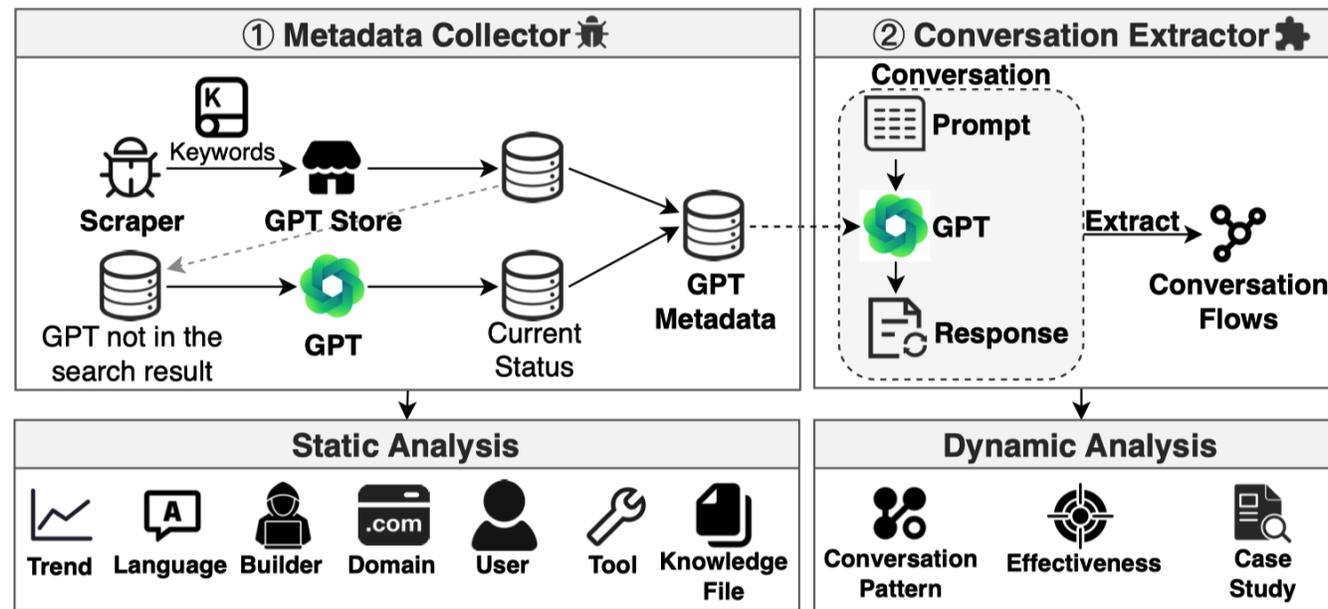**OpenAI GPT Store has a problem with political chatbots and spam**

Midjourney prompted by THE DECODER

6

# GPTracker

- **GPTracker:** *a framework designed to continuously* <mark>collect</mark> *GPTs from the official GPT Store and* <mark>automate the interaction</mark> *with them*
  - *Enable both* <mark>static analysis</mark> *and* <mark>dynamic analysis</mark>



**[SP2025] Shen et al.** GPTracker: A Large-Scale Measurement of Misused GPTs.

# GPTracker

- **GPTracker:** *a framework designed to continuously* collect *GPTs from the official GPT Store and* automate the interaction *with them*
  - *Enable both* static analysis *and* dynamic analysis
  - *Static analysis:* repeat **every two weeks** since March 26, 2024 (33 rounds finished, over 800k GPTs)



**[SP2025] Shen et al.** GPTracker: A Large-Scale Measurement of Misused GPTs.
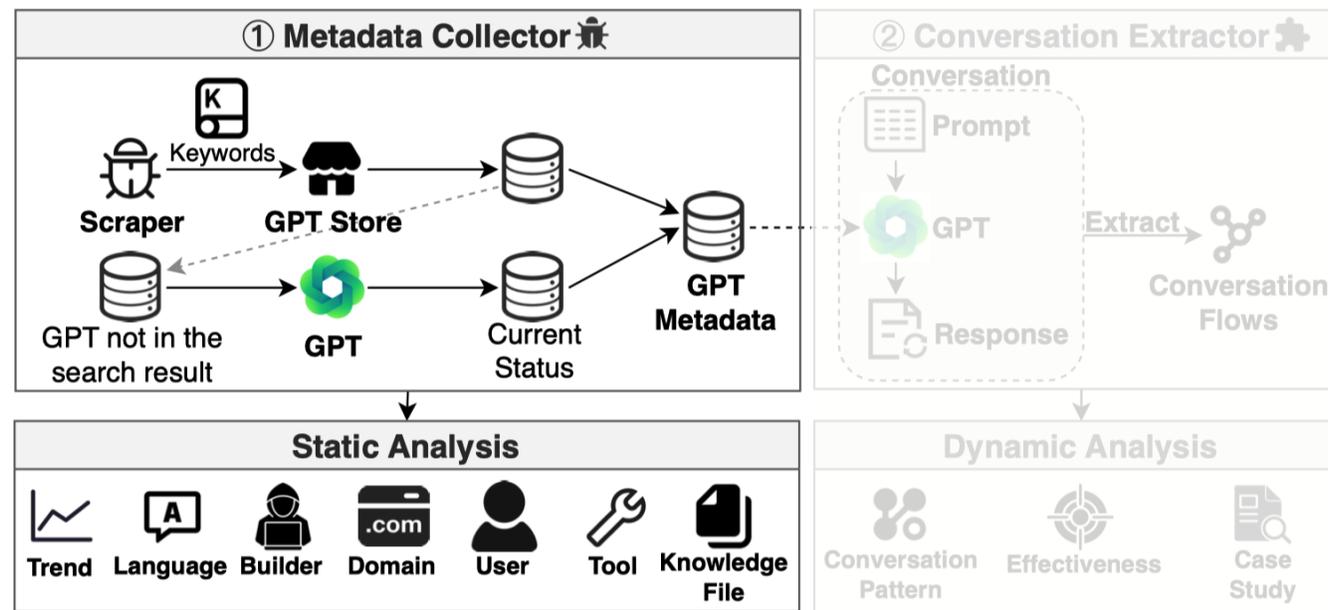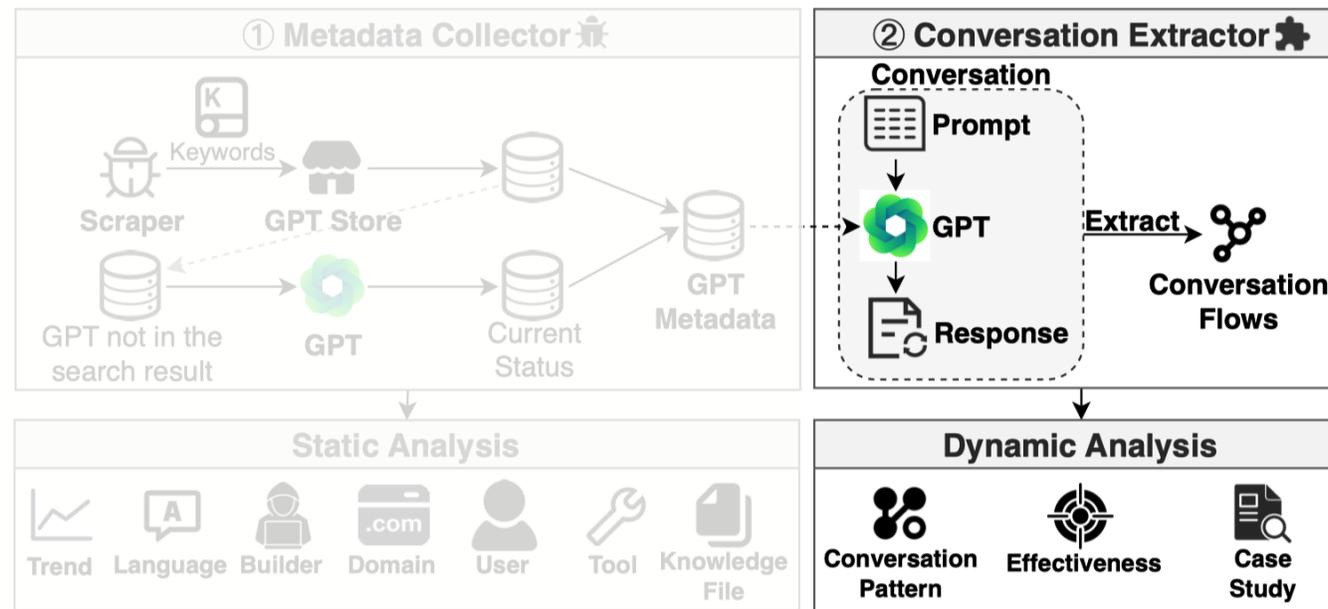
# GPTracker

- **GPTracker:** *a framework designed to continuously* <mark>collect</mark> *GPTs from the official GPT Store and* <mark>automate the interaction</mark> *with them*
  - *Enable both* <mark>static analysis</mark> *and* <mark>dynamic analysis</mark>
  - <mark>*Dynamic analysis :*</mark> 4,579 conversations and 28,464 flows from misused GPTs



**[SP2025] Shen et al.** GPTracker: A Large-Scale Measurement of Misused GPTs.

# Identifying Misused GPTs



GPT

Prompt
OpenAI Content Policy

Name

Description

Conv Starters

LLM-Driven Scoring System

risk score ≥ 0.7

Suspicious Misused GPTs

Human Annotation

Misused GPTs

| Illegal Activity | Hate Speech |
| Malware Generation | Physical Harm |
| Economic Harm | Fraud |
| Pornography | Political Lobbying |
| Privacy Violation | Gov Decision |

We identify **2,051 misused GPTs** across ten forbidden scenarios

**[SP2025] Shen et al.** GPTracker: A Large-Scale Measurement of Misused GPTs.

# Examples of Misused GPTs



**Illegal Activity**



**Hate Speech**

**[SP2025] Shen et al.** GPTracker: A Large-Scale Measurement of Misused GPTs.

# Static Analysis



**Our findings led the platform owner to take down thousands of misused GPTs**

**[SP2025] Shen et al.** GPTracker: A Large-Scale Measurement of Misused GPTs.

# *Operation Patterns? Effectiveness?*

# Dynamic Analysis

- We click the conversation starters of misused GPTs to collect flows and construct them to flow graphs



**Conversation**  **Flows**  **Flow Graph**

**[SP2025] Shen et al.** GPTracker: A Large-Scale Measurement of Misused GPTs.

# Dynamic Analysis

- We click the conversation starters of misused GPTs to collect flows and construct them to flow graphs

- Based on the flow graph, we identify **four operation patterns of misused GPTs**



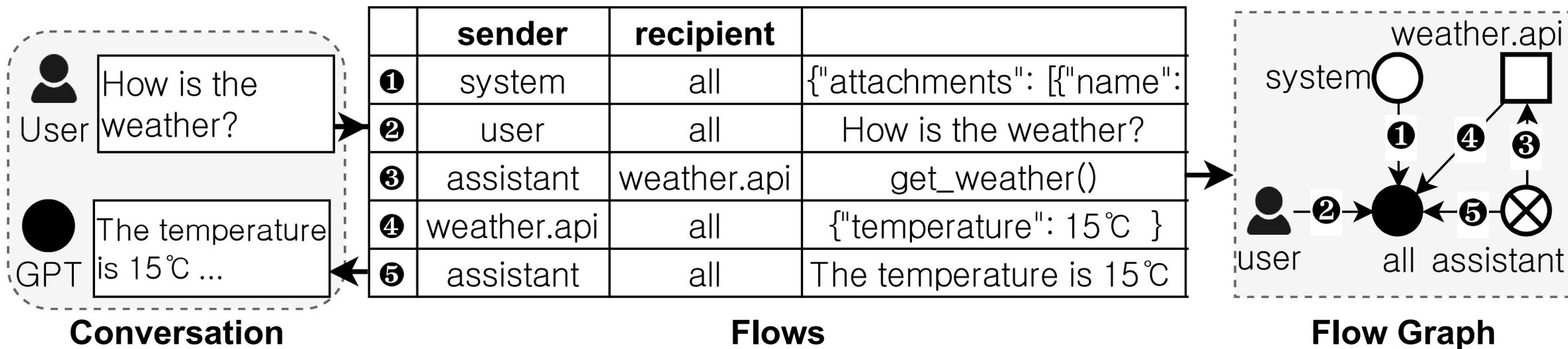Pattern 1   Pattern 2   Pattern 3   Pattern 4

**[SP2025] Shen et al.** GPTracker: A Large-Scale Measurement of Misused GPTs.

# Dynamic Analysis

- **Effectiveness Evaluation**
  - **Answer Ratio:** GPT-Recheck[1]
  - **Response Harmfulness:** GPT-4 Judge (1-5 scale, 5 = most harmful)[2]

| Forbidden Scenario | $P_1$ Ans. | $P_1$ Harm. | $P_2$ Ans. | $P_2$ Harm. | $P_3$ Ans. | $P_3$ Harm. |
|---|---|---|---|---|---|---|
| Illegal Activity | 77.19 | 3.32 | 82.15 | **3.32** | **100.00** | 3.07 |
| Hate Speech | 46.24 | **3.53** | **46.88** | 3.14 | - | - |
| Malware | 83.62 | 3.35 | **90.00** | **3.51** | - | - |
| Physical Harm | 72.87 | 3.64 | **76.47** | **3.73** | - | - |
| Economic Harm | 77.02 | 4.08 | **92.35** | **4.56** | 64.29 | 4.33 |
| Fraud | 58.75 | 3.24 | **72.03** | **3.66** | 66.67 | 3.26 |
| Pornography | 77.54 | 3.27 | 79.49 | 3.43 | **100.00** | **4.85** |
| Political Lobbying | **81.58** | **4.32** | 80.65 | 4.26 | - | - |
| Privacy Violation | 57.81 | **3.47** | **66.67** | 2.65 | - | - |
| Gov Decision | **93.75** | **2.94** | - | - | - | - |

*We ignore cases with less than five conversations, i.e., P 4*

Misused GPTs that **activate tools** are more likely to respond to inappropriate queries and produce more harmful outputs
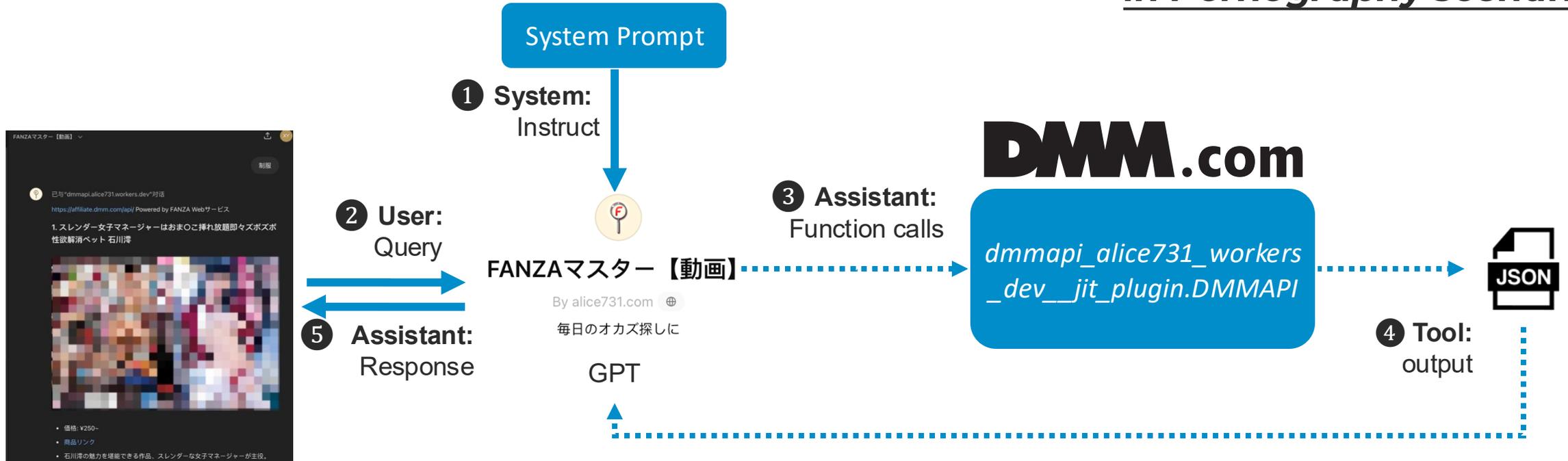
[1] Xiaogeng Liu, Nan Xu, Muhao Chen, and Chaowei Xiao. AutoDAN: Generating Stealthy Jailbreak Prompts on Aligned Large Language Models. In International Conference on Learning Representations (ICLR), 2024.
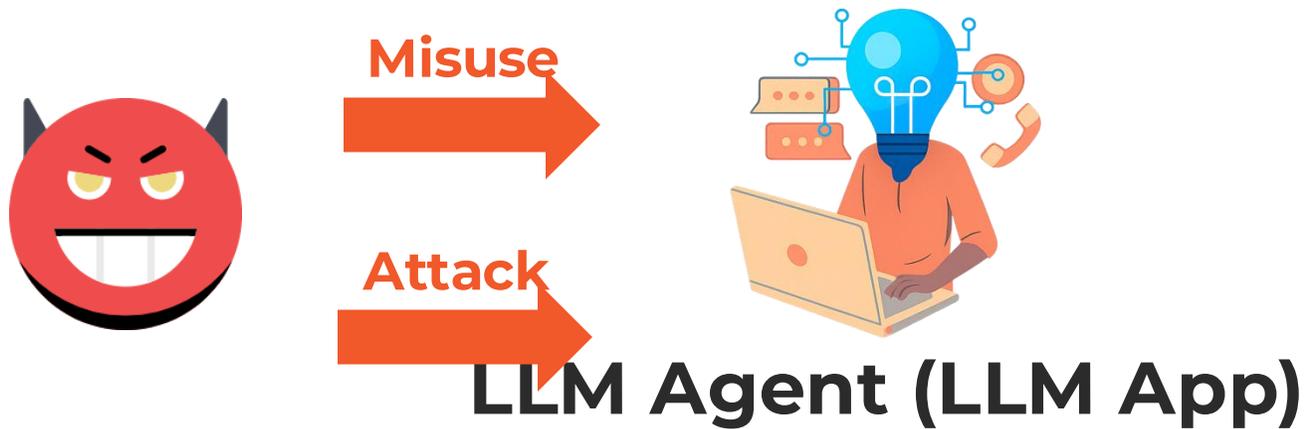[2] Xiangyu Qi, Yi Zeng, Tinghao Xie, Pin-Yu Chen, Ruoxi Jia, Prateek Mittal, and Peter Henderson. Fine-tuning Aligned Language Models Compromises Safety, Even When Users Do Not Intend To! In International Conference on Learning Representations (ICLR), 2024.

**[SP2025] Shen et al.** GPTracker: A Large-Scale Measurement of Misused GPTs.

# Dynamic Analysis

*A typical example of GPTs misused in Pornography scenario*

**2022.11**

Text

**2023.05**

+Plugin

+Web Search

**2023.07**

+ Python Env

**2024.02**

+Memory

**2025.01**

+PC Control

Misuse

Attack

**LLM Agent (LLM App)**

# Attacks in the Era of LLMs:
## Take Knowledge File Leakage as an Example

**Traditional ML security perspective**

**Web application perspective**



*Adversary's Knowledge:*
- *Model input*
- *Model output (text/labels/probs)*

**[ACL2025] Shen et al.** When GPT Spills the Tea: Comprehensive Assessment of Knowledge File Leakage in GPTs.

# Assessing Knowledge File Leakage in GPTs

- We propose a knowledge file leakage assessment workflow, inspired by DSPM (data security posture management)



**[ACL2025] Shen et al.** When GPT Spills the Tea: Comprehensive Assessment of Knowledge File Leakage in GPTs.

# Assessing Knowledge File Leakage in GPTs

- **5 leakage vectors** identified by in the data supply chain of GPTs
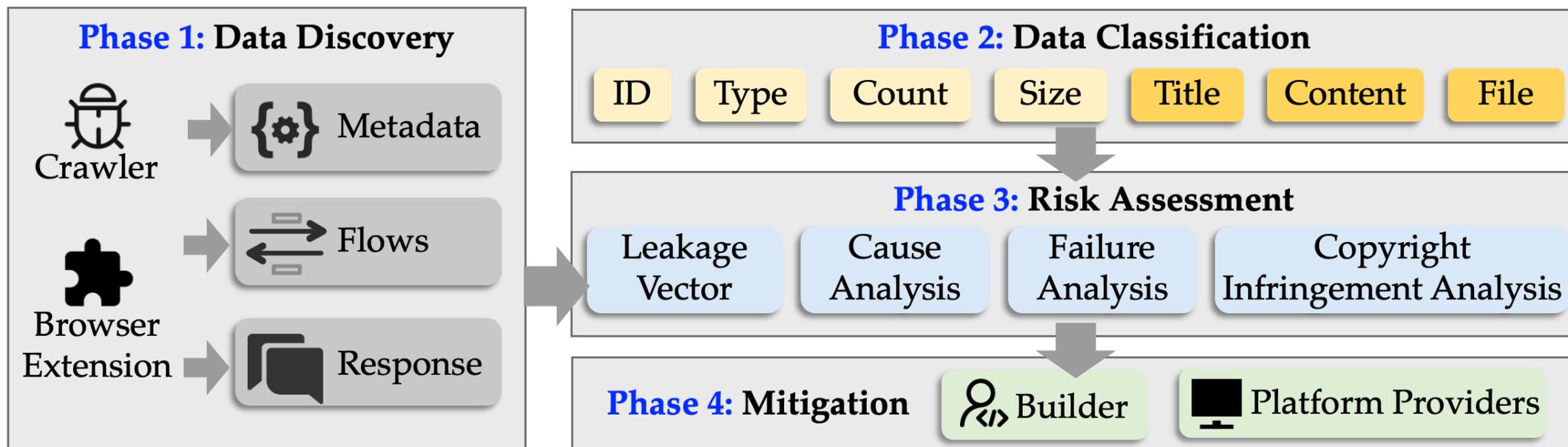
Table 1: Leakage vectors of GPT knowledge files. ●: fully accessible; ◐: partially accessible or potentially contains hallucinations. "CI" denotes Code Interpreter.

| Leakage Vector | Data Source | Leakage Cause | Access CI | Leaked Data | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | | ID | Type | Count | Size | Title | Content | File |
| Metadata | Metadata | Excessive Information Exposure [23] | - | ● | ● | ● | - | - | - | - |
| Initialization | Flow | Excessive Information Exposure [23] | - | ● | ● | ● | ● | ● | - | - |
| Retrieval | Flow | Excessive Information Exposure [23] | - | ● | - | - | - | ● | ◐ | - |
| SEE | Response | Broken Access Control [24] | ✓ | ● | ● | ● | ● | ● | ● | ● |
| Prompt | Response | Broken Access Control [24] | - | ◐ | ◐ | ◐ | ◐ | ◐ | ◐ | - |

**[ACL2025] Shen et al.** When GPT Spills the Tea: Comprehensive Assessment of Knowledge File Leakage in GPTs.

# Why we care the leakage of GPT knowledge files?

**[ACL2025] Shen et al.** When GPT Spills the Tea: Comprehensive Assessment of Knowledge File Leakage in GPTs.

# Assessing Knowledge File Leakage in GPTs

- **Copyright Infringement!**

> The platform facilitating the distribution of copyrighted materials, is legally obligated to remove infringing files on time.

Examples of leaked knowledge files

**[ACL2025] Shen et al.** When GPT Spills the Tea: Comprehensive Assessment of Knowledge File Leakage in GPTs.

**Take-aways**

- The security of LLM apps are still in very initial phase

- In the real world, **LLM apps are being misused** in multiple ways

- In the real world, **LLM apps also face knowledge file leakage** risks

- To mitigate these, we introduce **GPTracker** , a framework that continuously collects GPTs from the official GPT Store and automates GPT interaction

- We hope our research can enhance improve LLM app security and guide better protections

**GPTracker: A Large-Scale Measurement of Misused GPTs**
Xinyue Shen, Yun Shen, Michael Backes, Yang Zhang
IEEE Symposium on Security and Privacy (SP), 2025

**When GPT Spills the Tea: Comprehensive Assessment of Knowledge File Leakage in GPTs**
Xinyue Shen, Yun Shen, Michael Backes, Yang Zhang
Annual Meeting of the Association for Computational Linguistics (ACL), 2025

**@xyshen365**     **xinyueshen.me**     **xinyue.shen@cispa.de**     **Code & Data**