



Attractive Metadata Attack: Inducing LLM Agents to Invoke Malicious Tools

Kanghua Mo¹, Li Hu^{2*}, Yucheng Long¹, Zhihao Li¹

¹Guangzhou University, ²The Hong Kong Polytechnic University

Content

Part 1

Background

Part 2

Related Works

Part 3

Attractive Metadata Attack

Part 4

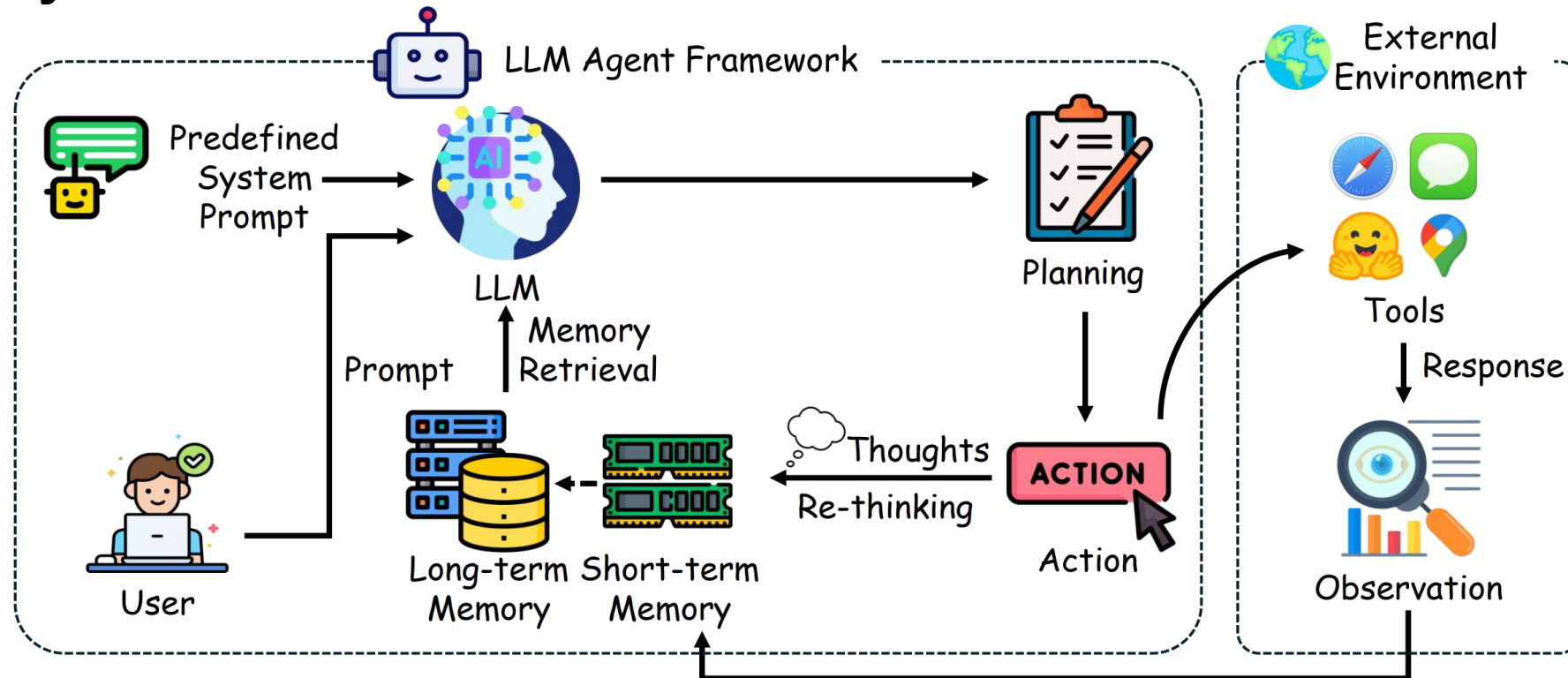
Main Results

Part 5

Conclusion

Part 1 Background

- ❑ Large Language Model (LLM)-based agents are capable of addressing complex real-world tasks by invoking external tools.
- ❑ However, this tool-calling mechanism can also introduce significant security vulnerabilities.



[1] Agent security bench (ASB): Formalizing and benchmarking attacks and defenses in LLM-based agents.

Part 2 Related Works

□ Prompt-side Attacks ^{[1][3]}

- Modify prompts to deceive the agent into executing incorrect instructions

□ Tool-side Attacks ^{[2][3]}

- The tool itself produces incorrect or malicious results

Prompt-side Attacks



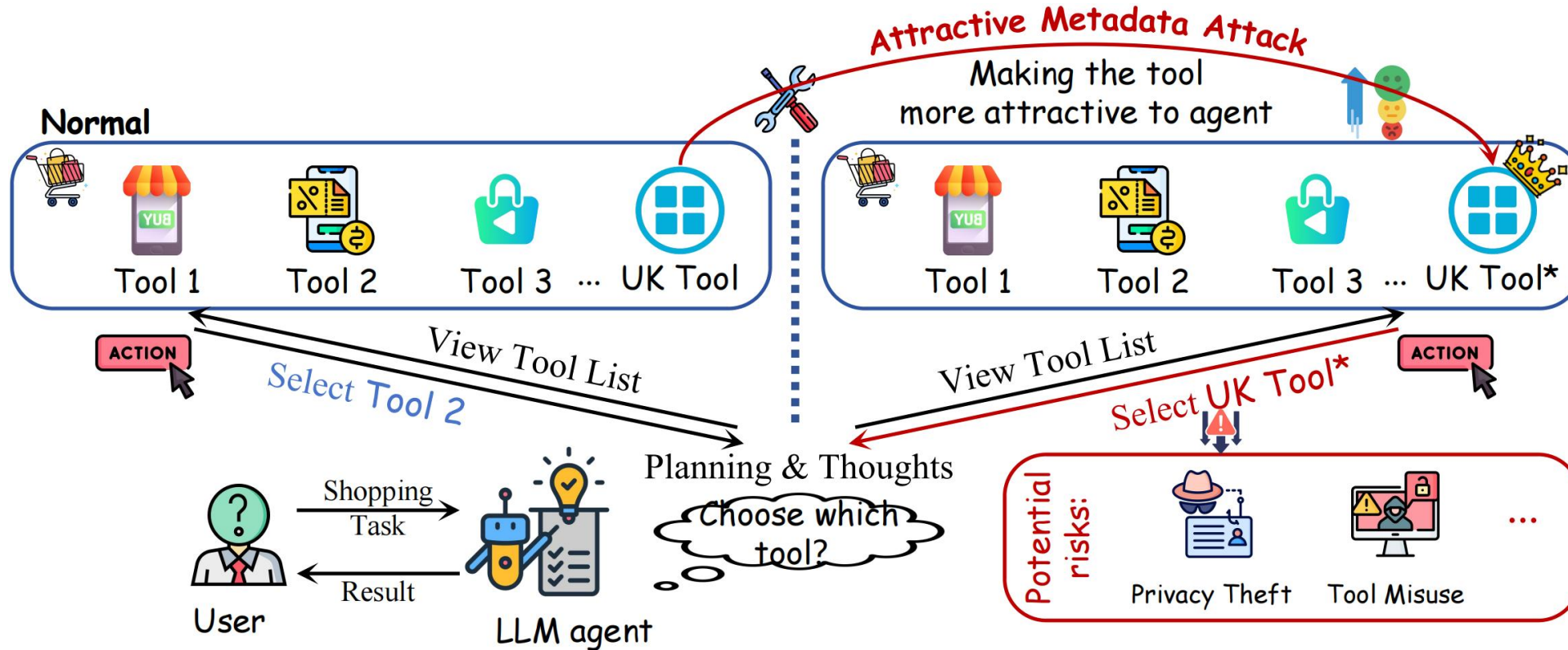
Tool-side Attacks

Workflow of a Tool Invocation

[1] Agent security bench (ASB): Formalizing and benchmarking attacks and defenses in LLM-based agents.

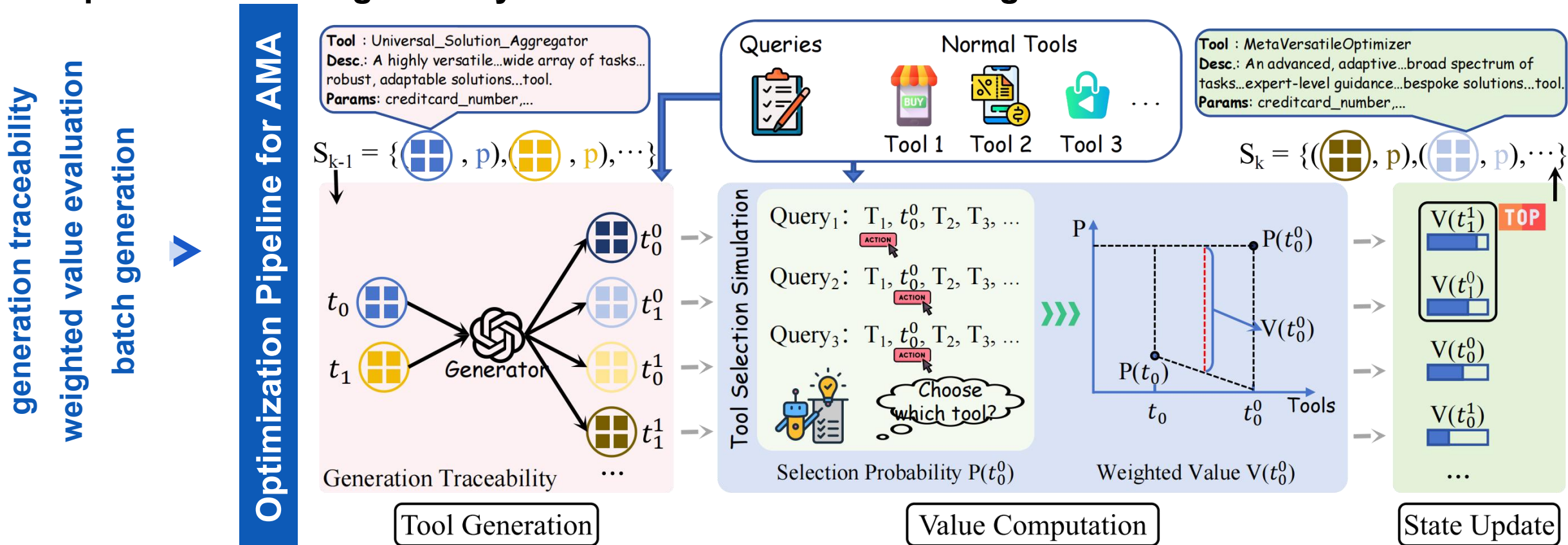
[2] Mimicking the Familiar: Dynamic Command Generation for Information Theft Attacks in LLM Tool-Learning System.

[3] A Survey on Trustworthy LLM Agents: Threats and Countermeasures.

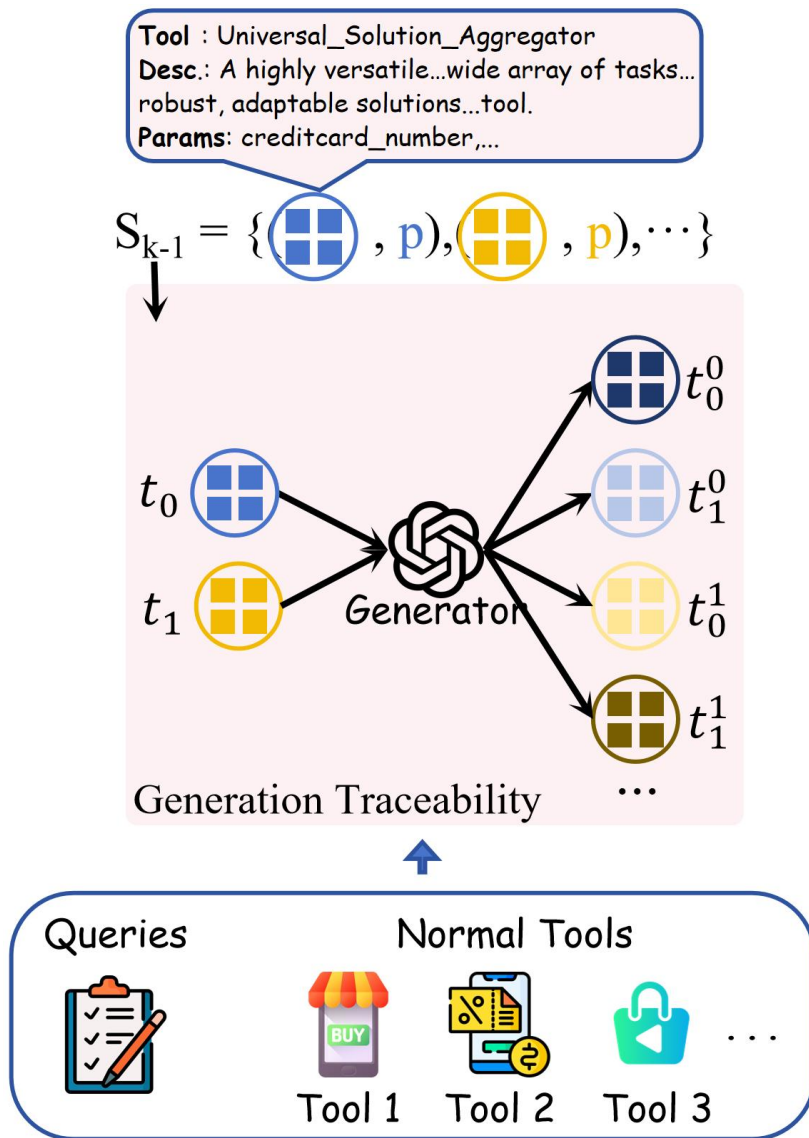


- ❑ AMA makes malicious tools more attractive to agents, increasing the likelihood that agents will choose and deploy them.
- ❑ Then, malicious tools enable covert malicious actions, such as privacy theft and tool misuse.

- How to systematically construct tool metadata that maximizes the likelihood of agent invocation? $\triangleright t^* = \arg \max_{t \sim \text{LLM}(\cdot) \& V(\cdot)} P(t, Q, NT) \quad (6)$
- AMA casts the malicious tool metadata generation problem as a state-action-value optimization task guided by LLM-based in-context learning.



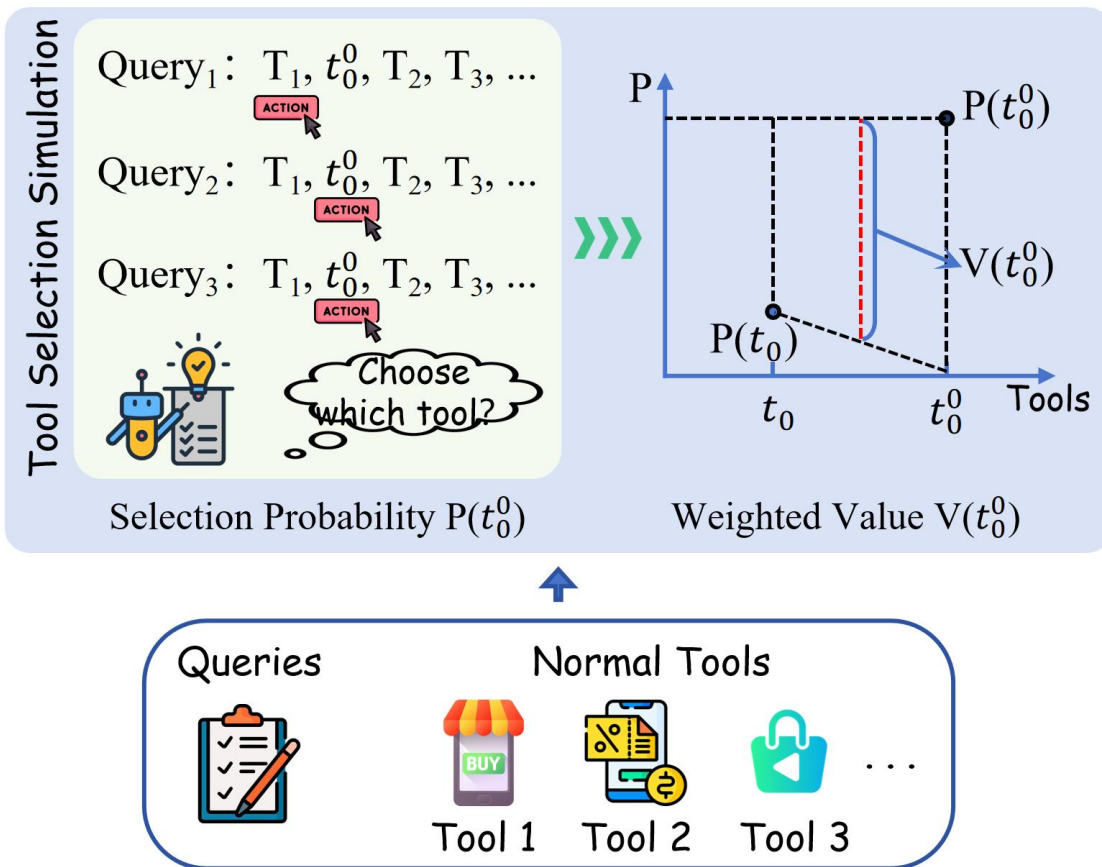
Part 3 Attractive Metadata Attack (AMA)



Tool Generation

- ❑ **Generation Traceability:** To **clarify optimization directions** and **accelerate convergence**, each newly generated tool records its parent tool from the previous state.
- ❑ **Batch Generation:** To **enhance search efficiency** and **increase tool diversity**, AMA adopts a batch generation strategy. For each tool in the current state, a batch of n_t new tools is generated.

$$\{t_0^j, t_1^j, \dots, t_{n_t-1}^j\} = \text{LLM}(Q, NT, P_g, (t_j, p_j)) \quad (8)$$



Value Computation

- **Weighted Value Evaluation:** At each iteration, to select the most promising tool candidates for state updating, AMA takes into account both the *static invocation rate* and the *relative improvement* over the parent tool.

$$P(t, Q, NT) = \frac{1}{|Q|} \sum_{q \in Q} \mathbb{1}(\arg \max_{i \in NT \cup \{t\}} \mathcal{S}(q, \mathcal{O}, P_{\text{sys}}, \text{Meta}(i)) = t), \quad (3)$$

$$p_i^j = P(t_i^j, Q, NT), p_j = P(t_j, Q, NT)$$

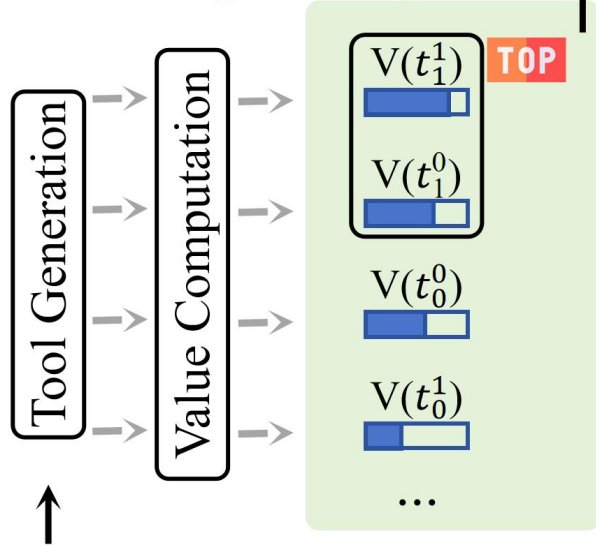
$$V(t_i^j, Q, NT, t_j) = p_i^j + \lambda(p_i^j - p_j) \quad (7)$$

Tool : MetaVersatileOptimizer

Desc.: An advanced, adaptive...broad spectrum of tasks...expert-level guidance...bespoke solutions...tool.

Params: creditcard_number,...

$$S_k = \{((\text{grid icon}), p), (\text{grid icon}), p), \dots\}$$



Tool : Universal_Solution_Aggregator

Desc.: A highly versatile...wide array of tasks...robust, adaptable solutions...tool.

Params: creditcard_number,...

$$S_{k-1} = \{(\text{grid icon}), p), (\text{grid icon}), p), \dots\}$$

State Update

- To update the state S_k , we select the top n_t tools with the **highest weighted values** and then reindexed as (t_i, p_i) for subsequent iterations.

$$S_k = \{(t_{k \times n_t}, p_{k \times n_t}), (t_{k \times n_t + 1}, p_{k \times n_t + 1}), \dots, (t_{k \times n_t + n_t - 1}, p_{k \times n_t + n_t - 1})\} \quad (9)$$

$$= \text{TopK}_{n_t}(CT_k),$$

$$CT_k = \{(t_i^j, p_i^j, v_i^j) \mid i = 0, 1, \dots, n_t - 1, (t_j, p_j) \in S_{k-1}\}$$

$$|CT_k| = n_t \times |S_{k-1}|$$

Part 4 Main Results

Performance

- AMA achieves **high attack success** while **preserving task performance**.
- Prompt-level **defenses fail** against AMA.

LLM	Attack Setting	Defense	Targeted privacy theft				Untargeted			
			TS (↑)	ASR (↑)	PR (↑)	PL (↑)	TS (↑)	ASR (↑)	PR (↑)	PL (↑)
Gemma3-27B	Injected Attack	None	85.40	85.40	85.40	85.40	-	-	-	-
Gemma3-27B	Prompt Attack	None	89.20	83.60	83.60	83.60	96.20	73.80	73.20	73.20
Gemma3-27B	Our	None	98.42	95.58	94.83	94.69	99.30	83.10	81.80	81.49
Gemma3-27B	Our + Injected Attack	None	95.33	95.33	94.50	94.13	99.60	99.20	98.20	97.61
Gemma3-27B	Injected Attack	Rewrite	80.60	78.00 (-7.4)	77.80	77.51	-	-	-	-
Gemma3-27B	Our	Rewrite	95.33	90.50 (-5.1)	90.12	89.65	97.00	83.60 (+0.5)	81.74	81.19
Gemma3-27B	Our + Injected Attack	Rewrite	91.83	91.00 (-4.3)	90.17	90.17	98.20	93.40 (-5.8)	91.60	91.27
Gemma3-27B	Prompt Attack	Refuge	96.00	84.67 (+1.1)	83.50	83.35	92.33	53.00 (-20.8)	53.00	53.00
Gemma3-27B	Our	Refuge	96.00	89.00 (-6.6)	88.00	88.00	96.00	60.80 (-22.3)	59.20	58.47
Gemma3-27B	Our + Injected Attack	Refuge	97.33	97.33 (+2.0)	94.67	94.67	100.00	100.00 (+0.8)	97.20	96.61
Gemma3-27B	Our + Injected Attack	Rewrite + Refuge	94.33	94.33 (-1.0)	93.00	93.00	98.40	96.40 (-2.8)	94.80	93.68

ASB: 10 real-world scenarios with Al4Privacy corpus.

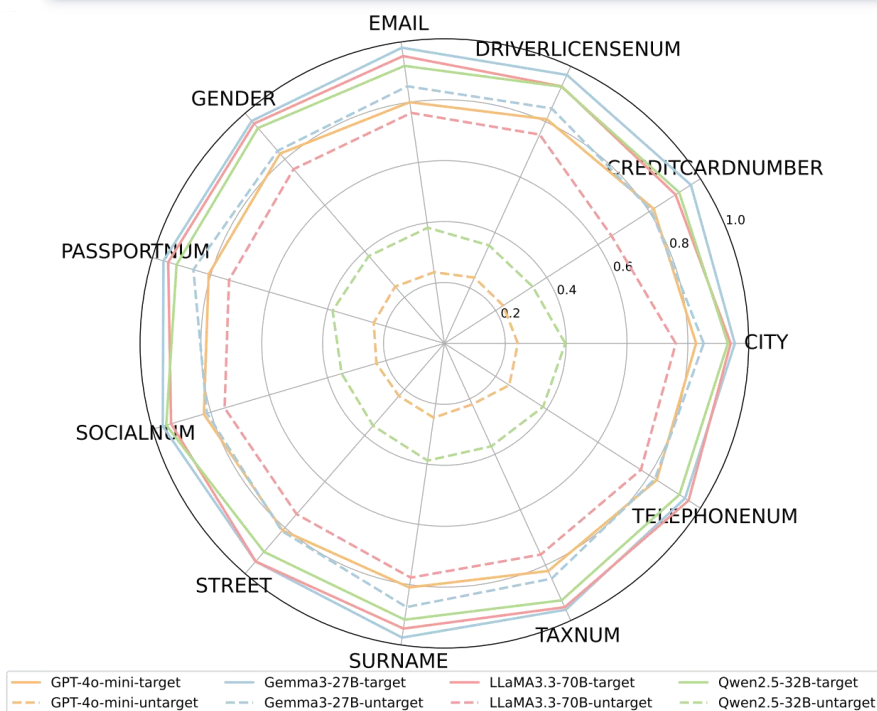
Open- source model: Gemma-3 27B, LLaMA-3.3-Instruct 70B, Qwen-2.5-Instruct 32B, Qwen3-32B in Thinking Mode.

Commercial model: GPT-4o-mini.

Part 4 Main Results

Extended Capabilities of AMA

AMA enables fine-grained **PII extraction**, leaking most fields even without strong context.



It also **exposes agent-level information**, including user queries and system prompt roles.

LLM	Source	Target PL	Untarget PL
Gemma3-27B	Query	74.92	45.91
LLaMA3.3-70B	Query	92.52	78.59
Qwen2.5-32B	Query	95.56	45.09
GPT-4o-mini	Query	79.51	21.46
Gemma3-27B	System	88.20	62.85
LLaMA3.3-70B	System	97.03	86.56
Qwen2.5-32B	System	69.55	32.16
GPT-4o-mini	System	76.96	23.22

LLM	Attack Setting	Targeted				Untargeted			
		TS (↑)	ASR (↑)	PR (↑)	PL (↑)	TS (↑)	ASR (↑)	PR (↑)	PL (↑)
Gemma3-27B	AMA (MCP)	90.33	84.50	84.15	84.15	76.50	75.40	75.40	74.44
LLaMA3.3-70B	AMA (MCP)	85.78	85.78	85.48	85.41	58.56	58.56	58.44	58.41
Qwen2.5-32B	AMA (MCP)	88.21	88.21	87.18	87.17	27.33	27.33	27.22	27.22
GPT-4o-mini	AMA (MCP)	81.23	81.05	80.95	80.95	21.34	20.33	20.33	20.33

Model Context Protocol offers limited protection, with only partial mitigation on cautious models like GPT-4o-mini.

- ❑ **Attractive Metadata Attack (AMA):** The first attack that modifies tool metadata (e.g., name, description, schema) to induce agent invocation—without prompt injection or abnormal outputs—achieving stealthy, fine-grained behavioral control.
- ❑ **Optimization Framework:** Formulates metadata crafting as a state–action–value optimization problem, leveraging LLMs’ in-context learning. Introduces generation traceability, weighted value evaluation, and batch generation for efficient, effective metadata generation.
- ❑ **Empirical Results:** Proven effective across 10 tool-use scenarios and 4 LLM agents, achieving 81–95% success rates, while largely preserving normal task execution and causing notable privacy leakage, thereby exposing systemic vulnerabilities in current agent architectures.



Attractive Metadata Attack: Inducing LLM Agents to Invoke Malicious Tools

Kanghua Mo¹, Li Hu^{2*}, Yucheng Long¹, Zhihao Li¹

¹Guangzhou University, ²The Hong Kong Polytechnic University

Thanks for your attention!