

JTPRO: A Joint Tool–Prompt Reflective Optimization Framework for Language Agents

Anonymous ACL submission

Abstract

Large language model (LLM) agents augmented with external tools often struggle as number of tools grow large and become domain-specific. In such settings, ambiguous tool descriptions and under-specified agent instructions frequently lead to tool mis-selection and incorrect slot/value instantiation. We hypothesize that this is due to two root causes: generic, one-size-fits-all prompts that ignore tool-specific nuances, and underspecified tool schemas that lack clear guidance on when and how to use each tool and how to format its parameters. We introduce **Joint Tool–Prompt Reflective Optimization (JTPRO)**, a framework that iteratively uses rollout-driven reflection to co-optimize global instructions and per-tool schema/argument descriptions. This is based on tool-confusion and slot/formatting errors. JTPRO is designed to preserve only tool-local cues needed for correct disambiguation and slot filling. We evaluate JTPRO across multi-tool benchmarks, which account for different number of tools using three metrics: Tool Selection Accuracy (TSA), Slot Filling Accuracy (SFA), and Overall Success Rate (OSR) (correct tool + correct slots + correct values). JTPRO consistently outperforms strong baselines, including CoT-style agents, and prompt optimizers such as GEPA by 5%–20% (relative) on OSR. Ablations show that joint optimization of instructions and tool schemas is more effective and robust than optimizing either component in isolation.

1 Introduction

Tool-augmented large language model (LLM) (Vaswani et al., 2017) agents extend their capabilities by invoking external tools for specialized operations and up-to-date information (Wang et al., 2024). As tool inventories grow in domain-specific deployments, agents must (i) select the correct tool among many conflicting options, (ii) instantiate



Figure 1: **Impact of slot-filling accuracy.** *Slot filling drives end-to-end success:* On the *Enterprise Tool-Inventory Dataset (ETID)* with complex schemas, we report TSA, SFA, and OSR; green overlays show absolute gains from JTPRO over baselines, highlighting that argument correctness is critical for OSR.

correct arguments from natural language requests; both suffer when tool/slot descriptions are ambiguous or underspecified (Qin et al., 2023). **Figure 2** quantifies this scaling failure on ToolACE (Liu et al., 2025): (a) tool selection accuracy drops as the tool universe expands, (b) a basic retrieval filter (top-20) only partially mitigates the decline. Crucially, end-to-end success is often bottlenecked by *slot/value instantiation*: on ETID¹ (Enterprise Tool Inventory Dataset - newly introduced synthetic dataset that we created for this study), **Figure 1** shows that improving slot filling produces large gains in overall success. Attempts to encode exhaustive tool and slot rules in lengthy global prompts are brittle—agents often fail to reliably follow extensive instructions, and maintaining cross-tool consistency becomes infeasible (Levy et al., 2024). **Figure 3** highlights two recurring failure modes that motivate JTPRO and show why *joint* optimization of global instructions and tool/argument schemas is necessary. In (a), two tools with over-

¹www.sites.google.com/view/jtpro-etid/home

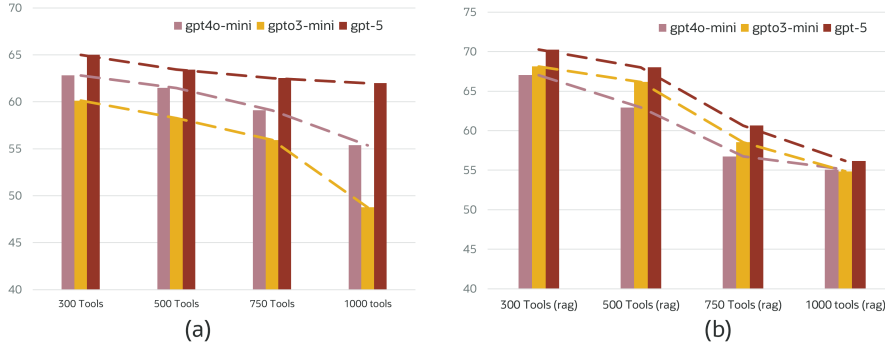


Figure 2: **Tool scaling failures and slot-filling impact.** (a) *All tools in context*: On ToolACE with an augmented inventory, tool selection accuracy drops as the tool set grows (300 to 1000), even for larger-context frontier models. (b) *Top-k retrieval*: A basic RAG with reranker stage (top-20) does not remove the drop, indicating residual tool disambiguation/argument issues.

lapping descriptions cause baseline mis-selection; JTPRO adds a brief tool-local decision rule to disambiguate them. In (b), the baseline picks the right tool but hallucinates/misnames arguments; JTPRO strengthens global tool-calling rules (while keeping tool-local parameter specs) to improve slot/value instantiation. Overall, these examples show that robust tool use at scale requires updating both the global policy layer and the implicated tool/slot descriptions.

Prior work improves tool use via largely separate levers: model tuning, tuning-free prompting/documentation, retrieval-based tool filtering, and prompt/context optimization. Tuning-free prompting (CoT (Wei et al., 2023), ReAct (Yao et al., 2023a)) and documentation refinement (DRAFT (Qu et al., 2025)) avoid weight updates but typically treat global instructions and tool schemas as static; retrieval-based selection reduces overload and iterative variants refine retrievers with agent feedback (Xu et al., 2024), yet retrieval alone does not fix downstream argument/format errors when slot semantics remain unclear. Prompt optimization and context evolution methods MIPRO (Opsahl-Ong et al., 2024); GEPA (Agrawal et al., 2025); AVATAR (Wu et al., 2024c); Dynamic Cheatsheet (Suzgun et al., 2025); ACE (Zhang et al., 2025) improve instruction-level behavior, but do not *jointly* adapt global decision rules and per-tool argument schemas at scale; similarly, Wu et al. (2025) refine prompts and tool descriptions but target efficiency rather than call-level tool/slot/value correctness under large domain tool stacks.

Our core contributions are as follows:

- **Joint optimization of tool/slot-schema and global instructions (JTPRO).** We formulate *joint*

optimization of (i) the global instruction prompt P and (ii) per-tool schema/argument descriptions $\{T_i\}$, targeting end-to-end invocation correctness (tool + slots + values) *without* model fine-tuning. This is critical as tool-use failures are inherently *coupled*: global policies depend on tool-local distinctions, and accurate slot/value instantiation relies on global conventions. Isolated optimization of P or $\{T_i\}$ is insufficient to address these interdependent failure modes.

- **Reflection-driven, localized edits with controlled growth.** Inspired by reflection-augmented prompt engineering (Agrawal et al., 2025), JTPRO diagnoses systematic rollout failures (tool confusion, missing constraints, format errors) and issues targeted edits to both P and relevant tool/slot descriptions. To prevent context bloat, we *globalize* recurring cross-tool slot semantics (date/time, bounds, currency/units) into P while keeping tool-specific exceptions locally without merging / aliasing for usability in real-world production systems.

- **Empirical evaluation under realistic constraints.** We benchmark JTPRO in both single- and multi-tool environments with variable argument structures, reporting Tool Selection Accuracy, Slot Filling Accuracy (conditional on tool correctness), and Overall Success Rate. JTPRO demonstrates clear gains over strong baselines (such as baseline CoT, GEPA, and MIPRO) and further enhances retrieval-based pipelines by improving both retrieval and downstream slot filling.

2 Related Work

Tool-use learning spans (i) tuning-based adaptation, (ii) tuning-free prompting and documentation refinement, (iii) retrieval-based tool selection, and

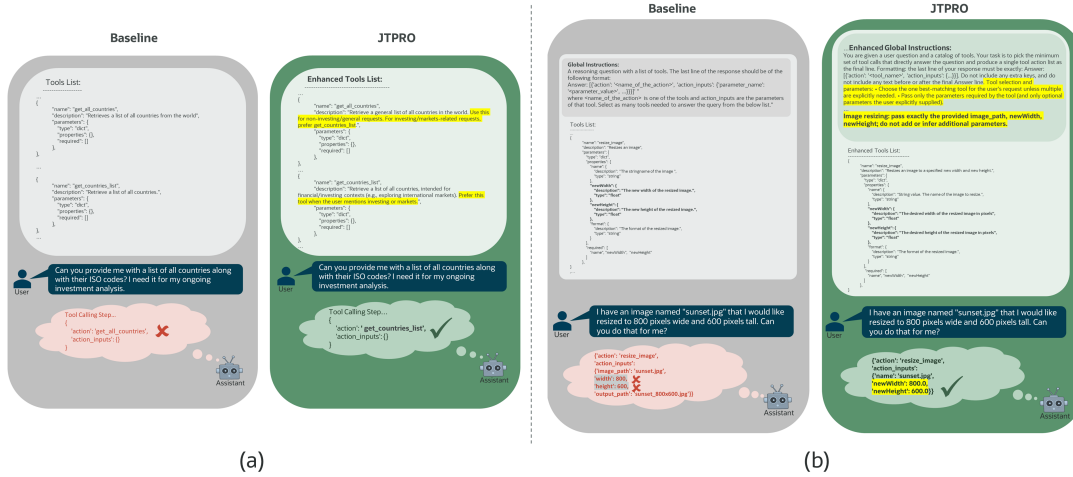


Figure 3: **Motivating context refinements with JTPRO.** (a) *Tool disambiguation*: Baseline does under-specify two similar tools, causing mis-selection; JTPRO adds brief per-tool decision rules (highlighted) to enable correct choice. (b) *Slot filling*: The baseline mis-specifies or hallucinates arguments; JTPRO tightens global tool-calling rules (highlighted) to enforce required fields and forbid extra/inferred slots.

(iv) prompt/context optimization. Most methods improve *either* the global prompt *or* tool specifications, but rarely their *joint* co-adaptation under large, evolving tool inventories.

Tuning-based tool learning. Model-tuning methods learn tool use by updating parameters or adding trainable modules from tool traces or preference/reward signals, including supervised fine-tuning (Qin et al., 2024), (Liu et al., 2025), contrastive objectives (Wu et al., 2024a), reinforcement learning (Feng et al., 2025; Qian et al., 2025), and tool-token embedding extensions (Alazraki and Rei, 2025). While effective, these methods require retraining as the underlying tools/schemas evolve.

Tuning-free prompting and documentation refinement. Prompting approaches like CoT and ReAct and agentic planners like RestGPT (Song et al., 2023) and HuggingGPT (Shen et al., 2023) elicit multi-step reasoning without weight updates, but usually treat instructions and tool schemas as static. DRAFT (Qu et al., 2025) improves per-tool documentation via trial-and-error, yet does not optimize global instruction policies or multi-tool interactions mediated by shared prompt rules.

Retriever-based tool selection. Retriever-based pipelines filter candidates via lexical/dense retrieval and specialized rerankers e.g., CRAFT (Yuan et al., 2024), ToolRerank (Zheng et al., 2024), COLT (Qu et al., 2024), improving scalability but not resolving argument/format errors when slot semantics are unclear. Iterative retrieval refinement with agent feedback (Xu et al., 2024) reduces retriever-agent

mismatch, but typically leaves the agent’s instruction layer largely unchanged.

Tool-using agents, tool construction, and prompt optimization. Toolformer (Schick et al., 2023), ReAct (Yao et al., 2023b), and ReWOO (Xu et al., 2023) integrate tool calls into reasoning traces; DSPy (Khattab et al., 2024) and AutoPDL (Spiess et al., 2025) support declarative tool programs but assume static prompts/schemas. Other work constructs tools (TOOLMAKER (Wölflein et al., 2025)), optimizes tool-use prompts (AvaTaR (Wu et al., 2024c)), calibrates tool use (CITI (Hao et al., 2025), PROBEAL (Liu et al., 2024)), or improves tool policies via SFT/RL (Sullivan et al., 2025). Separately, self-refinement and prompt optimization (Madaan et al., 2023; Shin et al., 2020; Lester et al., 2021; Pryzant et al., 2023; Yuksekgonul et al., 2025) and evolutionary search EvoPrompt (Guo et al., 2024) automate instruction improvement; MIPRO (Opsahl-Ong et al., 2024) optimizes module prompts and demonstrations, while GEPA (Agrawal et al., 2025) uses reflection over trajectories with Pareto selection, and AVATAR (Wu et al., 2024c) applies contrastive feedback. Dynamic Cheatsheet (Suzgun et al., 2025) and ACE (Zhang et al., 2025) motivate maintaining an evolving, curated context, but focus on strategy/memory rather than tool/argument schema co-adaptation.

Distinction. In contrast to prior work that optimizes prompts or tool documentation separately, JTPRO jointly updates global instructions P and per-tool *tool/argument* schema descriptions $\{T_i\}$

using rollout-driven reflection, targeting call-level correctness (tool, slots, values) without model fine-tuning. JTPRO also reduces redundancy by abstracting shared slot conventions globally while preserving tool-specific details locally, leading to improved results in retrieval-based pipelines.

3 Problem Statement

We consider an LLM agent with access to a set of N external tools (APIs/functions) $\{T_1, \dots, T_N\}$. Each tool T_i is specified by a schema/documentation entry describing its functionality and expected parameters (slots). Given a user query Q , the agent must produce an answer A , potentially by issuing one or more tool calls with structured arguments. The agent is guided by a global instruction prompt P and the collection of tool schemas $\{T_i\}_{i=1}^N$.

For a query Q , the LLM is invoked with context

$$C(P, T, Q) = P \mid T_1 \mid \dots \mid T_N \mid Q, \quad (1)$$

and produces a tool-call trace $\hat{\tau} = \hat{\tau}(P, T, Q)$.

Our objective is to optimize the textual content of P and $\{T_i\}$ to maximize tool-use performance *without* model fine-tuning. Because tool identities and interfaces are typically fixed in production, we do *not* merge or alias tools. Instead, we allow edits to P and each T_i , and we *globalize* recurring slot conventions (e.g., date/time formats, inclusive/exclusive bounds, currency/units) by lifting duplicated per-tool guidance into P .

We evaluate **call-level correctness**: correct tool selection and correct slot/value instantiation, summarized by **Tool Selection Accuracy**, **Slot Filling Accuracy** (conditional on correct tool), and **Overall Success Rate** (correct tool + correct slots + correct values). This emphasis matches deployments where executing tools and validating response-level correctness may be infeasible due to security, access control, rate limits, or non-deterministic backends.

Given a dataset $\mathcal{D} = \{(Q_j, \tau_j)\}_{j=1}^M$ with gold traces τ , we optimize only two variables—the global instructions P and tool descriptions T —to maximize expected call-level correctness:

$$(P^*, T^*) = \arg \min_{P, T} \mathbb{E}_{(Q, \tau) \sim \mathcal{D}} [\mathcal{L}(\hat{\tau}(P, T, Q), \tau)]. \quad (2)$$

The loss function can be defined using tool se-

lection, slot filling, and overall success:

$$\begin{aligned} \mathcal{L}(\hat{\tau}, \tau) &= \lambda_{\text{TSA}} (1 - \mathbb{I}[\hat{t} = t]) \\ &+ \lambda_{\text{SFA}} \mathbb{I}[\hat{t} = t] (1 - \text{Rec}(\hat{a}, a)) \\ &+ \lambda_{\text{OSR}} (1 - \mathbb{I}[\hat{t} = t \wedge \hat{a} = a]), \end{aligned} \quad (3)$$

where \hat{t} and t are the predicted and gold tool identifiers, $\mathbb{I}[\cdot]$ is the indicator function, \hat{a} and a are the predicted and gold argument structures, $\text{Rec}(\hat{a}, a)$ are slot/value recall conditional on $\hat{t} = t$, and $\lambda_{\text{TSA}}, \lambda_{\text{SFA}}, \lambda_{\text{OSR}}$ are nonnegative loss weights.

4 Technique

We present **Joint Tool–Prompt Reflective Optimization (JTPRO)**, a weight-free, context-level optimizer that iteratively updates (i) global agent instructions P and (ii) per-tool schemas $\{T_i\}_{i=1}^N$ from labeled tool-call traces. Algorithm 1 summarizes the loop.

Setup and objective For each query q , the agent runs under $C(q) = P \mid T_1 \mid \dots \mid T_N \mid q$ and produces a predicted trace $\hat{\tau}$. Given gold traces τ^* , JTPRO edits P and $\{T_i\}$ to improve TSA, SFA (conditional on correct tool), and OSR (correct tool + correct slots + correct values).

Candidate selection (Pareto) JTPRO maintains a pool \mathcal{C} of candidate contexts and uses GEPA-style Pareto selection: retain candidates that achieve the best score on at least one training instance, prune strictly dominated candidates, then sample a starting candidate with probability biased toward those that win on more instances.

Rollouts, diagnostics, and localized edits On a minibatch $\mathcal{B} \subset \mathcal{D}_{tr}$, we compute rollout metrics and extract structured failure signals \mathcal{F} via $\text{DIAGNOSE}(\hat{\tau}, \tau^*)$ (e.g., tool confusions, missing required slots, formatting/value violations). A reflector proposes targeted edits $(\Delta P, \{\Delta T_i\}) \leftarrow \text{PROPOSEEDITS}(\mathcal{F}, P^o, \{T_i^o\})$, which are applied to produce a draft context P^d and $\{T_i^d\}$. Edits are localized to the implicated global rules and tool/slot descriptions.

Merge-with-best for incremental tool adaptation JTPRO tracks a validation-best context $C^* = (P^*, \{T_i^*\})$. After editing, we merge P^d with P^* using $\text{MERGEWITHBEST}(P^d, P^*)$ to form P' , implementing a “growing playbook” that preserves cross-cutting rules while adding new, rollout-driven

guidance. This accumulation also supports incremental toolset expansion: when new T_i are appended, stable global conventions remain intact and new tool-triggered rules are integrated without re-optimizing from scratch.

Algorithm 1 JTPRO: Reflective Schema-Instruction Co-Optimization with Slot-Semantics Globalization

Input: initial global instructions $P^{(0)}$, initial tool schemas $\{D_i^{(0)}\}_{i=1}^N$, labeled training set $\mathcal{D}_{tr} = \{(q, \tau^*)\}$, labeled validation set \mathcal{D}_{val} , max iterations I , batch size B
Output: optimized global instructions P^* , optimized tool schemas $\{D_i^*\}_{i=1}^N$

```

Initialize  $P \leftarrow P^{(0)}$ ,  $T_i \leftarrow T_i^{(0)}$  for all  $i \in \{1, \dots, N\}$ 
Initialize best context  $C^* \leftarrow (P, \{T_i\})$  and best validation score  $s^* \leftarrow -\infty$ 
Initialize pool  $\mathcal{C} \leftarrow \{(P^*, T^*)\}$  // candidate contexts
for  $t = 1$  to  $I$  do // main optimization loop
  Sample a minibatch  $\mathcal{B} \subset \mathcal{D}_{tr}$  of size  $B$ 
   $(P^o, T^o) \leftarrow \text{ParetoSelect}(\mathcal{C})$ 
  Initialize aggregated feedback  $\mathcal{F} \leftarrow \emptyset$ 
  foreach  $(q, \tau^*) \in \mathcal{B}$  do
    Construct context  $C(q) \leftarrow P^o | T_1^o | \dots | T_N^o | q$ 
    Run agent to obtain predicted trace  $\hat{\tau} \leftarrow \text{AGENT}(C(q))$ 
    Compute rollout metrics (TSA, SFA, OSR)  $\leftarrow \text{EVAL}(\hat{\tau}, \tau^*)$ 
    Extract error signals  $f \leftarrow \text{DIAGNOSE}(\hat{\tau}, \tau^*)$ 
     $\mathcal{F} \leftarrow \mathcal{F} \cup \{f\}$ 
  // Propose localized edits to both global instructions and tool context
   $(\Delta P, \{\Delta T_i\}) \leftarrow \text{PROPOSEEDITS}(\mathcal{F}, P^o, \{T_i^o\})$ 
   $P^d \leftarrow \text{APPLY}(P^o, \Delta P)$ 
   $T_i^d \leftarrow \text{APPLY}(T_i^o, \Delta T_i) \forall i$ 
  // Merge localized edits with global best
   $P' \leftarrow \text{MERGEWITHBEST}(P^d, P^*)$ 
  // Globalize repetitive slot semantics
   $(P'', \{T_i''\}) \leftarrow \text{GLOBALIZESLOTS}(P', \{T_i^d\})$ 
  // Accept/reject based on held-out performance
   $s'' \leftarrow \text{SCORE}((P'', \{T_i''\}), \mathcal{B})$ 
  if  $s'' > \text{SCORE}((P^o, \{T_i^o\}), \mathcal{B})$  then
    // Run eval on the entire validation set
     $s''_{val} \leftarrow \text{SCORE}((P'', \{T_i''\}), \mathcal{D}_{val})$ 
    if  $s''_{val} \geq s^o_{val}$  then
       $\mathcal{C} \leftarrow \text{ADDTOPool}(\mathcal{C}, (P'', T''), K)$ 
      // Update global best instructions if improved
      if  $s''_{val} > s^*$  then
         $C^* \leftarrow (P'', T'')$ 
         $s^* \leftarrow s''_{val}$ 
  return  $C^*$  as  $(P^*, \{T_i^*\}_{i=1}^N)$ 

```

Globalizing repetitive slot semantics To reduce duplicated schema text, JTPRO applies GLOBALIZESLOTS($P', \{T_i^d\}$) \mapsto ($P'', \{T_i''\}$), lifting recurring cross-tool slot conventions into P'' and removing redundant per-tool restatements while keeping tool-specific exceptions and disambiguation rules in T_i'' . Figure 7 motivates this step: in ETID, a

small set of slot families (e.g., identifiers and date/time) recur across many tools (up to 77/124), producing substantial verbatim repetition in per-tool schemas.

Acceptance and pool update We score ($P'', \{D_i''\}$) on \mathcal{B} and, if improved, evaluate on \mathcal{D}_{val} . Improved candidates are added to \mathcal{C} (bounded size K), and if a candidate is best on validation we update C^* accordingly.

Summary JTPRO combines Pareto-selected candidate search, reflection-driven localized edits to P and $\{T_i\}$, and globalization of shared slot semantics to improve both tool selection and argument correctness in large tool inventories.

5 Datasets and Evaluation

5.1 Datasets

We evaluate JTPRO on 3 complementary benchmarks that stress different failure modes in tool-using agents: (i) complex, domain-specific slot filling with a moderate tool inventory, (ii) tool selection under toolset scaling and (iii) a *multi-tool calling* setting where a single query may require invoking multiple tools in parallel and correctly instantiating arguments at each step.

Enterprise Tool-Inventory Dataset (ETID). ETID is a domain-specific tool-calling dataset targeting *argument correctness* under complex schemas. It contains 124 tools with 3.4 parameters on average (max 12) and ~ 13 labeled examples per tool (min 10). We evaluate both an *all-tools* setting and *value-stream* subsets. For data efficiency, we use intent-aligned regimes Train- N ex where each tool contributes N train and N validation examples (total $124 \times N$), reporting Train-1ex/2ex/4ex. The test set is fixed at 404 queries.

ToolACE (tool scaling). ToolACE evaluates performance degradation as the tool universe expands. We use fixed splits (Train = 199, Validation = 76, Test = 121) and augment the tool inventory to create ToolACE-300/500/750/1000 variants.

SEAL-Tools (parallel multi-tool calling). SEAL-Tools (Wu et al., 2024b) benchmarks *parallel* multi-tool calling across diverse domains. We use a curated multiple-overlap subset with 1,138 tools and 2,743 arguments, split into Train = 600, Validation = 100, Test = 100. Each query requires 3.2 parallel tool calls on average (typically

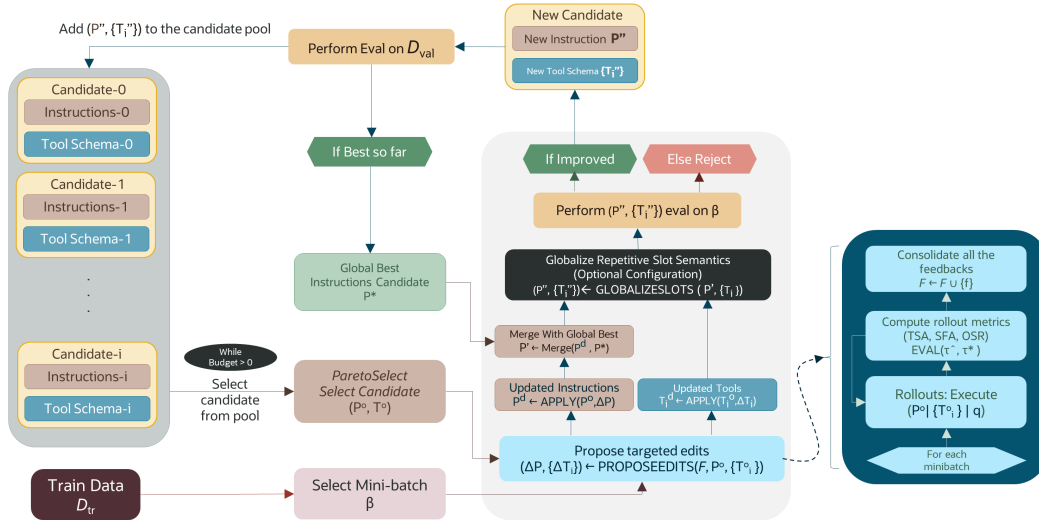


Figure 4: **JTPRO optimization loop (block-diagram view)**. JTPRO maintains a pool of candidate contexts (global instructions P and tool schemas $\{T_i\}$) and repeatedly (i) selects a candidate via Pareto-based sampling, (ii) runs minibatch rollouts on \mathcal{D}_{tr} to compute tool-use metrics (TSA, SFA, OSR) and aggregate error feedback, and (iii) proposes localized edits to both P and the implicated tool schemas. The edited instructions are merged with the current global-best P^* , followed by optional globalization of repetitive slot semantics to avoid duplicated cross-tool parameter rules. Candidates that improve minibatch performance are validated on \mathcal{D}_{val} ; improved candidates are added back to the pool, and the global best ($P^*, \{T_i^*\}$) is updated when a new highest validation score is observed.

3), with 5.8 arguments filled per query, stressing joint multi-tool selection and argument filling.

We use a curated multiple-overlap subset containing 1,138 tools with 2,743 arguments. The split is Train = 600, Validation = 100, Test = 100 examples. Each query requires an average of 3.2 parallel tool calls (77% require exactly 3 tools) with 5.8 arguments filled per query. Tool coverage overlap ensures training tools appear in evaluation splits. This setting isolates the challenge of *joint* tool selection and slot filling at scale, where models must correctly identify multiple tools and fill all arguments for each.

5.2 Evaluation Metrics

Following prior tool-use evaluations, we measure call-level correctness rather than answer accuracy. Specifically, we report:

- **Tool Selection Accuracy (TSA)**: fraction of queries for which the agent chose the correct tool(s) required (including choosing none if no tool needed).
- **Slot Filling Accuracy (SFA)**: recall of correct slot/value assignments *conditional on correct tool selection*.
- **Overall Success Rate (OSR)**: (correct tool +

| Dataset | #Tools | Total Args | | Required Args | |
|--------------|--------|------------|-----|---------------|-----|
| | | Avg | Max | Avg | Max |
| ETID | 124 | 3.4 | 12 | 0.81 | 6 |
| ToolACE-300 | 336 | 2.05 | 14 | 1.20 | 6 |
| ToolACE-500 | 536 | 2.14 | 17 | 1.20 | 7 |
| ToolACE-750 | 786 | 2.17 | 23 | 1.21 | 7 |
| ToolACE-1000 | 1036 | 2.10 | 23 | 1.21 | 7 |
| SEAL-Tools | 1138 | 2.41 | 8 | 1.60 | 5 |

Table 1: Dataset statistics. “#Tools” denotes the size of the tool universe available at inference time. “Total Args” counts all schema parameters per tool, and “Required Args” counts mandatory parameters.

correct slots + correct values).

This evaluation reflects practical deployments where executing the true tool backend may be infeasible (e.g., security constraints, access controls, rate limits, or non-deterministic systems), so correctness must be assessed at the tool-call level.

6 Results and Analysis

6.1 ToolACE: Scaling the Tool Universe

Table 2 reports Tool Selection Accuracy (TSA), Slot Filling Accuracy (SFA; conditional on correct tool), and Overall Success Rate (OSR; correct tool

| Model | #Tools | TSA (%) | | | SFA (%) ↑ TSA | | | OSR (%) | | |
|------------|--------|---------|-------|-------|-----------------|-------|--------------|---------|-------|--------------|
| | | Base | GEPA | JTPRO | Base | GEPA | JTPRO | Base | GEPA | JTPRO |
| gpt4o-mini | 500 | 63.832 | 73.33 | 75.25 | 86.996 | 85.27 | 88.12 | 60.00 | 61.98 | 69.42 |
| gpt4o-mini | 1000 | 61.115 | 73.39 | 75.13 | 86.67 | 83.36 | 83.59 | 58.18 | 60.33 | 63.64 |
| gpto3-mini | 500 | 70.78 | 73.33 | 76.19 | 84.994 | 85.27 | 88.52 | 59.454 | 61.98 | 65.29 |
| gpto3-mini | 1000 | 58.916 | 70.11 | 71.48 | 85.036 | 86.59 | 87.46 | 51.272 | 58.68 | 64.46 |
| gpt-5 | 500 | 73.02 | 77.17 | 82.28 | 84.785 | 85.75 | 90.00 | 62.73 | 66.12 | 74.38 |
| gpt-5 | 1000 | 67.658 | 75.13 | 78.72 | 87.35 | 86.40 | 89.26 | 62.366 | 67.77 | 73.55 |

Table 2: **ToolACE results under tool-universe scaling (500 vs. 1000 tools).** JTPRO achieves the strongest end-to-end performance (OSR) by jointly improving tool selection (TSA) and argument correctness (SFA), with the largest gains appearing in the 1000-tool regime where tool confusions are most frequent.

+ correct slots + correct values) for ToolACE with 500 and 1000 tools.

As the tool inventory grows, baseline performance drops primarily in TSA, which cascades to lower OSR even when SFA remains high. GEPA improves TSA in most settings, but gains in OSR are limited because failures often stem from tool-specific disambiguation and argument constraints that global instruction refinement alone cannot resolve.

JTPRO consistently achieves the highest TSA and OSR across all models and tool counts. Gains are especially pronounced in the 1000-tool setting (e.g., +13.2 OSR points for gpto3-mini over baseline). While SFA is already strong, JTPRO further boosts end-to-end success by reducing tool confusions and encoding missing slot/value conventions. These results show that, on ToolACE, OSR improvements are primarily driven by better tool selection, emphasizing that accurate TSA is critical for downstream argument correctness.

6.2 ETID: Complex Slot Filling with Moderate Tool Counts

We evaluate on the *Enterprise Tool-Inventory Dataset (ETID)*, which features complex multi-argument schemas (avg. 3.4 parameters/intent; max 12) and measures correctness at the *call level* (tool + slots + values). Table 3 reports results under low-supervision regimes (Train-1/2/4 examples per intent; fixed test set of 404 queries).

Two trends stand out. First, **slot/value correctness is the main bottleneck**. Baseline TSA is high (85–94%), yet OSR remains much lower, showing that SFA errors dominate once the correct tool is chosen. JTPRO addresses this directly, improving SFA and boosting OSR—e.g., for gpt4o-mini,

OSR rises from 44.8→60.15 (+15.35) in Train-1ex and 46.53→66.83 (+20.30) in Train-4ex, despite similar TSA.

Second, **JTPRO delivers robust gains across models and training regimes**. For gpto3-mini, OSR improves over both baseline and GEPA in all regimes, with larger gains as supervision increases (Train-4ex: 67.33→82.67). For gpt-5, GEPA raises TSA, but JTPRO achieves the highest OSR by combining strong tool selection with higher SFA (Train-2ex: SFA 92.77, OSR 85.64). Overall, ETID shows that optimizing tool selection alone is insufficient; joint refinement of instructions and tool/slot descriptions is necessary to convert high TSA into end-to-end success.

6.3 SEAL-Tools: Multi-Tool Calling

On SEAL-Tools, JTPRO consistently improves end-to-end accuracy by enhancing slot filling while keeping tool selection high. For gpt-41-mini, TSA remains stable at 82–83%, but SFA rises from 31.3% to 36.6% (+16.8%), boosting OSR from 26.0% to 30.0% (+15.4%). A similar pattern holds for o4-mini (TSA 78–79%, SFA 26.6→30.8%, OSR 21.0→24.0%) and gpt-4o (TSA 81–82%, SFA 28.4→33.5%, OSR 23.0→27.5%).

These gains arise from JTPRO’s tool description refinements, including explicit parameter guidance (e.g., “*pass country_code as 2-letter lower-case ISO code*”) and disambiguation cues (e.g., “*PRE-FERRED for general-purpose web discovery*”). In complex multi-tool schemas with 1,138 tools and an average of 3.2 parallel calls per query, these joint instruction and tool description optimizations translate high baseline tool selection into substantially higher end-to-end success.

| Model | Train | TSA (%) | | | SFA (%) ↑ TSA | | | OSR (%) | | |
|------------|-----------|---------|-------|-------|-----------------|--------|--------------|---------|--------|--------------|
| | | Base | GEPA | JTPRO | Base | GEPA | JTPRO | Base | GEPA | JTPRO |
| gpt4o-mini | Train-1ex | 85.91 | 86.23 | 83.14 | 69.81 | 78.71 | 82.72 | 44.80 | 50.19 | 60.15 |
| gpt4o-mini | Train-2ex | 86.36 | 85.32 | 88.27 | 70.90 | 77.12 | 83.37 | 45.79 | 52.7 | 65.10 |
| gpt4o-mini | Train-4ex | 86.96 | 88.38 | 87.18 | 70.09 | 74.93 | 85.16 | 46.53 | 54.34 | 66.83 |
| gpto3-mini | Train-1ex | 94.00 | 95.08 | 95.78 | 80.48 | 83.73 | 89.77 | 68.81 | 75.00 | 79.46 |
| gpto3-mini | Train-2ex | 94.15 | 94.40 | 95.40 | 80.20 | 85.30 | 90.01 | 67.33 | 73.02 | 79.70 |
| gpto3-mini | Train-4ex | 94.41 | 94.90 | 95.95 | 79.60 | 87.455 | 90.98 | 67.325 | 77.725 | 82.67 |
| gpt-5 | Train-1ex | 94.06 | 97.27 | 95.76 | 82.45 | 90.21 | 92.15 | 68.81 | 80.20 | 84.65 |
| gpt-5 | Train-2ex | 94.03 | 97.91 | 95.80 | 83.90 | 89.89 | 92.77 | 71.53 | 79.18 | 85.64 |
| gpt-5 | Train-4ex | 94.16 | 98.47 | 96.81 | 83.29 | 88.80 | 92.30 | 70.54 | 80.69 | 85.15 |

Table 3: **ETID results under low-supervision training regimes.** TSA: tool selection accuracy. SFA: slot/value correctness conditional on correct tool. OSR: end-to-end call-level correctness (correct tool + correct slots + correct values). JTPRO yields the most consistent OSR improvements, indicating that improved argument semantics are crucial for complex enterprise schemas.

| Model | TSA (%) | | SFA (%) | | OSR (%) | |
|-------------|---------|-------|---------|-------------|---------|-------------|
| | Base | JTPRO | Base | JTPRO | Base | JTPRO |
| gpt-41-mini | 83 | 82 | 31.3 | 36.6 | 26.0 | 30.0 |
| o4-mini | 79 | 78 | 26.6 | 30.8 | 21.0 | 24.0 |
| gpt-4o | 81 | 82 | 28.4 | 33.5 | 23.0 | 27.5 |

Table 4: **SEAL-Tools results (multi-tool calling).** TSA: tool selection accuracy. SFA: slot filling accuracy conditional on correct tools. OSR: end-to-end success rate (correct tools + correct arguments). JTPRO improves SFA and OSR across all models while maintaining stable TSA.

6.4 Tool Description Disambiguation

JTPRO resolves ambiguity between semantically similar tools by refining tool descriptions. On ToolACE-500, 11% of descriptions (55/500) were updated with explicit disambiguation cues. For example, search was modified to “*NOT for general web article discovery—prefer ‘web_search’*”, while web_search became “*PREFERRED for general-purpose discovery*”.

We quantify this effect via intra-group cosine similarity across 37 groups of confusable tools (e.g., get_ip_*, get_user_*). The group with the largest improvement reduced similarity from 0.668 to 0.502 (−16.6%), demonstrating clearer differentiation. Full details are in Appendix D.9.

7 Conclusion

We presented **Joint Tool–Prompt Optimization (JTPRO)**, a weight-free context optimization framework that jointly refines global agent instructions and per-tool schema/argument descriptions from rollout-driven feedback. JTPRO targets the two dominant failure modes in large, domain-specific tool inventories: tool mis-selection and argument mis-instantiation. The method uses reflective diagnostics to produce localized edits, maintains a candidate pool with Pareto-style selection to preserve diverse effective behaviors, and prevents context bloat by *globalizing* recurring slot semantics into the instruction layer while retaining tool-specific disambiguation cues in local schemas. Across ToolACE tool-scaling experiments and ETID enterprise slot-filling tasks, JTPRO improves tool selection, slot filling, and overall success relative to strong baselines including CoT-style agents and prompt optimizers (e.g., GEPA), highlighting that accurate argument semantics are necessary to translate high tool selection accuracy into end-to-end tool-use success. These results support joint schema–instruction co-adaptation as a practical route to maintaining tool-using agents under evolving tool inventories without model fine-tuning.

Limitations

Our study has limitations that motivate future work. First, our experimental settings cover (i) single-tool, single-slot/value cases and (ii) multi-tool *parallel* calling with single-slot/value instantiation; we do

not evaluate *sequential* multi-tool workflows that require multi-step dependencies, intermediate state, or long-horizon planning (e.g., tool chains where earlier outputs condition later calls). Extending JTPRO to such settings will require modeling step-wise credit assignment across tool sequences and validating robustness under longer rollouts.

Second, while ETID captures complex multi-argument schemas, our current evaluation does not systematically stress *deeply nested* argument structures (e.g., multi-layer JSON objects, lists-of-objects with constraints, or schema-dependent composition rules) at scale; future benchmarks should include nested-slot correctness and structure-aware metrics beyond scalar slot/value matching.

Third, our evaluation focuses on call-level correctness (tool, slots, values) rather than executing tools and verifying response-level correctness; when tool execution is available, future experiments should extend JTPRO to end-to-end evaluation that includes tool responses and downstream post-processing logic (e.g., response parsing, aggregation, and business-rule enforcement), since these components can introduce additional failure modes beyond argument correctness.

Fourth, ETID is currently not publicly released; we are actively working with our legal team to enable publication of ETID in a compliant form, which would facilitate broader reproducibility and benchmarking by the community.

Finally, our empirical study is limited to 3 benchmarks; future work should broaden coverage to additional public and proprietary tool-use datasets spanning more domains, tool granularity, and interaction styles, to better characterize generalization under diverse tool inventories and schema conventions.

References

Lakshya A Agrawal, Shangyin Tan, Dilara Soyulu, Noah Ziemis, Rishi Khare, Krista Opsahl-Ong, Arnav Singhvi, Herumb Shandilya, Michael J Ryan, Meng Jiang, Christopher Potts, Koushik Sen, Alexandros G. Dimakis, Ion Stoica, Dan Klein, Matei Zaharia, and Omar Khattab. 2025. [Gepa: Reflective prompt evolution can outperform reinforcement learning](#). *Preprint*, arXiv:2507.19457.

Lisa Alazraki and Marek Rei. 2025. [Meta-reasoning improves tool use in large language models](#). In *Findings of the Association for Computational Linguistics: NAACL 2025*, page 7885–7897. Association for Computational Linguistics.

Jiazhan Feng, Shijue Huang, Xingwei Qu, Ge Zhang, Yujia Qin, Baoquan Zhong, Chengquan Jiang, Jinxin Chi, and Wanjun Zhong. 2025. [Retool: Reinforcement learning for strategic tool use in llms](#). *Preprint*, arXiv:2504.11536.

Qingyan Guo, Rui Wang, Junliang Guo, Bei Li, Kaitao Song, Xu Tan, Guoqing Liu, Jiang Bian, and Yujiu Yang. 2024. [Connecting large language models with evolutionary algorithms yields powerful prompt optimizers](#). In *The Twelfth International Conference on Learning Representations*.

Yupu Hao, Pengfei Cao, Zhuoran Jin, Huanxuan Liao, Yubo Chen, Kang Liu, and Jun Zhao. 2025. [Citi: enhancing tool utilizing ability in large language models without sacrificing general performance](#). In *Proceedings of the Thirty-Ninth AAAI Conference on Artificial Intelligence and Thirty-Seventh Conference on Innovative Applications of Artificial Intelligence and Fifteenth Symposium on Educational Advances in Artificial Intelligence, AAAI’25/IAAI’25/EAAI’25*. AAAI Press.

Omar Khattab, Arnav Singhvi, Paridhi Maheshwari, Zhiyuan Zhang, Keshav Santhanam, Sri Vardhamanan, Saiful Haq, Ashutosh Sharma, Thomas T. Joshi, Hanna Moazam, Heather Miller, Matei Zaharia, and Christopher Potts. 2024. [Dspy: Compiling declarative language model calls into self-improving pipelines](#).

Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. [The power of scale for parameter-efficient prompt tuning](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 3045–3059, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Mosh Levy, Alon Jacoby, and Yoav Goldberg. 2024. [Same task, more tokens: the impact of input length on the reasoning performance of large language models](#).

Hao Liu, Zi-Yi Dou, Yixin Wang, Nanyun Peng, and Yisong Yue. 2024. [Uncertainty calibration for tool-using language agents](#). In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 16781–16805, Miami, Florida, USA. Association for Computational Linguistics.

Weiwen Liu, Xu Huang, Xingshan Zeng, Xinlong Hao, Shuai Yu, Dexun Li, Shuai Wang, Weinan Gan, Zhengying Liu, Yuanqing Yu, Zezhong Wang, Yuxian Wang, Wu Ning, Yutai Hou, Bin Wang, Chuhan Wu, Xinzhi Wang, Yong Liu, Yasheng Wang, and 8 others. 2025. [Toolace: Winning the points of llm function calling](#). *Preprint*, arXiv:2409.00920.

Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, Shashank Gupta, Bodhisattwa Prasad Majumder, Katherine Hermann, Sean Welleck, Amir Yazdanbakhsh, and Peter Clark. 2023. [Self-refine: iterative](#)

| | | |
|-----|---|-----|
| 606 | refinement with self-feedback. In <i>Proceedings of the 37th International Conference on Neural Information Processing Systems, NIPS '23</i> , Red Hook, NY, USA. Curran Associates Inc. | |
| 607 | | |
| 608 | | |
| 609 | | |
| 610 | Krista Opsahl-Ong, Michael J Ryan, Josh Purtell, David Broman, Christopher Potts, Matei Zaharia, and Omar Khattab. 2024. Optimizing instructions and demonstrations for multi-stage language model programs . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 9340–9366, Miami, Florida, USA. Association for Computational Linguistics. | |
| 611 | | |
| 612 | | |
| 613 | | |
| 614 | | |
| 615 | | |
| 616 | | |
| 617 | | |
| 618 | Reid Pryzant, Dan Iter, Jerry Li, Yin Lee, Chenguang Zhu, and Michael Zeng. 2023. Automatic prompt optimization with “gradient descent” and beam search . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 7957–7968, Singapore. Association for Computational Linguistics. | |
| 619 | | |
| 620 | | |
| 621 | | |
| 622 | | |
| 623 | | |
| 624 | | |
| 625 | Cheng Qian, Emre Can Acikgoz, Qi He, Hongru Wang, Xiusi Chen, Dilek Hakkani-Tür, Gokhan Tur, and Heng Ji. 2025. Toolrl: Reward is all tool learning needs . <i>Preprint</i> , arXiv:2504.13958. | |
| 626 | | |
| 627 | | |
| 628 | | |
| 629 | Yujia Qin, Shengding Hu, Yankai Lin, Weize Chen, Ning Ding, Ganqu Cui, Zheni Zeng, Yufei Huang, Chaojun Xiao, Chi Han, Yi Ren Fung, Yusheng Su, Huadong Wang, Cheng Qian, Runchu Tian, Kunlun Zhu, Shihao Liang, Xingyu Shen, Bokai Xu, and 22 others. 2024. Tool learning with foundation models . <i>Preprint</i> , arXiv:2304.08354. | |
| 630 | | |
| 631 | | |
| 632 | | |
| 633 | | |
| 634 | | |
| 635 | | |
| 636 | Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, Sihan Zhao, Lauren Hong, Runchu Tian, Ruobing Xie, Jie Zhou, Mark Gerstein, Dahai Li, Zhiyuan Liu, and Maosong Sun. 2023. Toollm: Facilitating large language models to master 16000+ real-world apis . | |
| 637 | | |
| 638 | | |
| 639 | | |
| 640 | | |
| 641 | | |
| 642 | | |
| 643 | Changle Qu, Sunhao Dai, Xiaoqi Wei, Hengyi Cai, Shuaiqiang Wang, Dawei Yin, Jun Xu, and Ji-Rong Wen. 2024. Towards completeness-oriented tool retrieval for large language models . In <i>Proceedings of the 33rd ACM International Conference on Information and Knowledge Management, CIKM '24</i> , page 1930–1940. ACM. | |
| 644 | | |
| 645 | | |
| 646 | | |
| 647 | | |
| 648 | | |
| 649 | | |
| 650 | Changle Qu, Sunhao Dai, Xiaoqi Wei, Hengyi Cai, Shuaiqiang Wang, Dawei Yin, Jun Xu, and Ji-Rong Wen. 2025. From exploration to mastery: Enabling llms to master tools via self-driven interactions . | |
| 651 | | |
| 652 | | |
| 653 | | |
| 654 | Timo Schick, Jane Dwivedi-Yu, Roberto Dessi, Roberta Raileanu, Maria Lomeli, Eric Hambro, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. 2023. Toolformer: Language models can teach themselves to use tools . In <i>Thirty-seventh Conference on Neural Information Processing Systems</i> . | |
| 655 | | |
| 656 | | |
| 657 | | |
| 658 | | |
| 659 | | |
| 660 | Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li, Weiming Lu, and Yueting Zhuang. 2023. Hugging-gpt: Solving ai tasks with chatgpt and its friends in hugging face . <i>Preprint</i> , arXiv:2303.17580. | |
| 661 | | |
| 662 | | |
| 663 | | |
| | Taylor Shin, Yasaman Razeghi, Robert L. Logan IV, Eric Wallace, and Sameer Singh. 2020. AutoPrompt: Eliciting Knowledge from Language Models with Automatically Generated Prompts . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 4222–4235, Online. Association for Computational Linguistics. | 664 |
| | | 665 |
| | | 666 |
| | | 667 |
| | | 668 |
| | | 669 |
| | | 670 |
| | Yifan Song, Weimin Xiong, Dawei Zhu, Wenhao Wu, Han Qian, Mingbo Song, Hailiang Huang, Cheng Li, Ke Wang, Rong Yao, Ye Tian, and Sujian Li. 2023. Restgpt: Connecting large language models with real-world restful apis . <i>Preprint</i> , arXiv:2306.06624. | 671 |
| | | 672 |
| | | 673 |
| | | 674 |
| | | 675 |
| | Claudio Spiess, Mandana Vaziri, Louis Mandel, and Martin Hirzel. 2025. AutoPDL: Automatic Prompt Optimization for LLM Agents . <i>arXiv e-prints</i> , arXiv:2504.04365. | 676 |
| | | 677 |
| | | 678 |
| | | 679 |
| | Michael Sullivan, Mareike Hartmann, and Alexander Koller. 2025. Procedural environment generation for tool-use agents . In <i>Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing</i> , pages 18555–18573, Suzhou, China. Association for Computational Linguistics. | 680 |
| | | 681 |
| | | 682 |
| | | 683 |
| | | 684 |
| | | 685 |
| | Mirac Suzgun, Mert Yuksekogun, Federico Bianchi, Dan Jurafsky, and James Zou. 2025. Dynamic cheat-sheet: Test-time learning with adaptive memory . <i>Preprint</i> , arXiv:2504.07952. | 686 |
| | | 687 |
| | | 688 |
| | | 689 |
| | Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In <i>Proceedings of the 31st International Conference on Neural Information Processing Systems, NIPS'17</i> , page 6000–6010, Red Hook, NY, USA. Curran Associates Inc. | 690 |
| | | 691 |
| | | 692 |
| | | 693 |
| | | 694 |
| | | 695 |
| | | 696 |
| | Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang, Xu Chen, Yankai Lin, Wayne Xin Zhao, Zhewei Wei, and Jirong Wen. 2024. A survey on large language model based autonomous agents . <i>Frontiers of Computer Science</i> , 18(6). | 697 |
| | | 698 |
| | | 699 |
| | | 700 |
| | | 701 |
| | | 702 |
| | Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2023. Chain-of-thought prompting elicits reasoning in large language models . | 703 |
| | | 704 |
| | | 705 |
| | | 706 |
| | Georg Wölflein, Dyke Ferber, Daniel Truhn, Ognjen Arandjelovic, and Jakob Nikolas Kather. 2025. LLM agents making agent tools . In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 26092–26130, Vienna, Austria. Association for Computational Linguistics. | 707 |
| | | 708 |
| | | 709 |
| | | 710 |
| | | 711 |
| | | 712 |
| | | 713 |
| | Bin Wu, Edgar Meij, and Emine Yilmaz. 2025. A joint optimization framework for enhancing efficiency of tool utilization in LLM agents . In <i>Findings of the Association for Computational Linguistics: ACL 2025</i> , pages 22361–22373, Vienna, Austria. Association for Computational Linguistics. | 714 |
| | | 715 |
| | | 716 |
| | | 717 |
| | | 718 |
| | | 719 |

720 Jiayi Wu, Renyu Zhu, Nuo Chen, Qiushi Sun, Xi-
721 ang Li, and Ming Gao. 2024a. [Structure-aware](#)
722 [fine-tuning for code pre-trained models](#). *Preprint*,
723 arXiv:2404.07471.

724 Mengsong Wu, Tong Zhu, Han Han, Chuanyuan Tan,
725 Xiang Zhang, and Wenliang Chen. 2024b. [Seal-](#)
726 [tools: Self-instruct tool learning dataset for agent](#)
727 [tuning and detailed benchmark](#). *arXiv preprint*
728 *arXiv:2405.08355*.

729 Shirley Wu, Shiyu Zhao, Qian Huang, Kexin Huang,
730 Michihiro Yasunaga, Kaidi Cao, Vassilis N. Ioanni-
731 dis, Karthik Subbian, Jure Leskovec, and James Zou.
732 2024c. [Avatar: Optimizing LLM agents for tool usage](#)
733 [via contrastive reasoning](#). In *The Thirty-eighth*
734 *Annual Conference on Neural Information Process-*
735 *ing Systems*.

736 Bin Feng Xu, Zhiyuan Peng, Bowen Lei, Subhabrata
737 Mukherjee, Yuchen Liu, and Dongkuan Xu. 2023.
738 [Rewoo: Decoupling reasoning from observations](#)
739 [for efficient augmented language models](#). *Preprint*,
740 arXiv:2305.18323.

741 Qiancheng Xu, Yongqi Li, Heming Xia, and Wenjie Li.
742 2024. [Enhancing tool retrieval with iterative feed-](#)
743 [back from large language models](#).

744 Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak
745 Shafraan, Karthik Narasimhan, and Yuan Cao. 2023a.
746 [React: Synergizing reasoning and acting in language](#)
747 [models](#).

748 Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak
749 Shafraan, Karthik Narasimhan, and Yuan Cao. 2023b.
750 [ReAct: Synergizing reasoning and acting in language](#)
751 [models](#). In *International Conference on Learning*
752 *Representations (ICLR)*.

753 Lifan Yuan, Yangyi Chen, Xingyao Wang, Yi R. Fung,
754 Hao Peng, and Heng Ji. 2024. [Craft: Customiz-](#)
755 [ing llms by creating and retrieving from specialized](#)
756 [toolsets](#). *Preprint*, arXiv:2309.17428.

757 Mert Yuksekgonul, Federico Bianchi, Joseph Boen,
758 Sheng Liu, Pan Lu, Zhi Huang, Carlos Guestrin,
759 and James Zou. 2025. [Optimizing generative ai by](#)
760 [backpropagating language model feedback](#). *Nature*,
761 639:609–616.

762 Qizheng Zhang, Changran Hu, Shubhangi Upasani,
763 Boyuan Ma, Fenglu Hong, Vamsidhar Kamanuru,
764 Jay Rainton, Chen Wu, Mengmeng Ji, Hanchen Li,
765 Urmish Thakker, James Zou, and Kunle Olukotun.
766 2025. [Agentic context engineering: Evolving con-](#)
767 [texts for self-improving language models](#). *Preprint*,
768 arXiv:2510.04618.

769 Yuanhang Zheng, Peng Li, Wei Liu, Yang Liu, Jian
770 Luan, and Bin Wang. 2024. [Toolrerank: Adap-](#)
771 [tive and hierarchy-aware reranking for tool retrieval](#).
772 *Preprint*, arXiv:2403.06551.

A Method prompts and details 773

A.1 Update Global Instructions and Tool Revisions Prompt 774 775

Prompt Template: Propose Updated Global Instructions and Tool Revisions (JT-PRO Reflector)

Goal: produce clean, minimal updates to the global instructions and only the tools that require revision, based on the feedback trace.

Task. Update the global instructions and tool descriptions using the feedback on the current context.

Current Global Instructions

<curr_instructions>

Current Tool Definitions (Full List)

<tools_list_to_update>

Objective. Produce:

- `global_instructions`: updated system-level guidance.
- `example_specific_instructions`: batch-specific guidance to append to prior examples/hints.
- `tool_revisions`: only the tools you modified (not the entire tool list).

Feedback Trace

<dataset_with_feedback>

For each example, the feedback indicates whether the model:

- failed to call a tool that should have been called, or
- called the wrong tool instead of the correct one, or
- selected the correct tool but produced incorrect `action_inputs` (missing/wrong parameters), or
- selected the correct parameters but assigned incorrect slot names/values (formatting/value errors).

Instructions Revision Rules (Important)

- The instructions may already contain prior revisions and examples.
- Always preserve previously incorporated `example_specific_instructions` and example hints; do not alter them.
- Modify `example_specific_instructions` existing `example_specific_instructions` only if there is a direct conflict with the current feedback.
- Add new guidance as bullet points appended to the existing list under `example_specific_instructions`.

Tool Revision Rules (Important)

- If an answer is marked wrong, check two cases:

- **Case 1 (Model error):** If a tool-selection error occurred (a tool should have been chosen but was not, or was chosen but should not have been), you *must* revise the relevant tool description(s) to make correct usage clearer. Otherwise, you may leave the tool unchanged.
 - **Case 2 (Documentation/data issue):** Tool documentation may be ambiguous or incomplete, and ground-truth traces may contain incorrect tool arguments or values. If you detect such issues, revise the tool documentation to remove ambiguity so future runs avoid the same failure.
- Return only the tools you modified from the provided tool list.

Output Format (Strict)

- Return a single JSON object immediately after the literal text Answer:
- Do not add any extra text before or after the JSON.

Required JSON Schema

```
{
  "global_instructions": "UPDATED
  GLOBAL INSTRUCTIONS HERE",
  "example_specific_instructions":
  "UPDATED INSTRUCTIONS AS PER THE
  CURRENT BATCH",
  "tool_revisions": [
    {
      "name": "<tool1_name>",
      "description":
      "<tool1_description>",
      "parameters": {
        "type": "dict",
        "properties": {
          "<property1>": {
            "description": "updated
            tool1_property1_description",
            "type": "<property1_type; same as
            original>",
            "<property2>": {
              "description": "updated
              tool1_property2_description",
              "type": "<property2_type; same as
              original>"
            },
            "required": ["<property1;
            required parameters; identical to
            original>"]
          },
          "required": null
        }
      },
      "required": null
    }
  ]
}
```

Final constraint: Answer: must be followed by *only* the JSON object.

Figure 5: JTPRO reflector prompt for proposing new global instructions and tool candidates. The reflector updates system-level guidance and selectively revises only the implicated tools/slots based on rollout feedback, while preserving prior batch-specific examples and enforcing a strict JSON-only output format.

Prompt Template: Merge Draft Instructions with the Global Best (Merge-WithBest)

Role. You are the *instruction merger* in agent. Your job is to combine a draft update with the current best global instructions into a single improved instruction prompt.

Inputs

- **Global best instructions** (P^*):
<best_global_instructions>
- **Draft instructions from current rollout** (P^d):
<draft_global_instructions>
- **(Optional) Newly added tools since P^* :**
<new_tools_summary>
(Names + 1–2 lines per tool describing the new capability.)

Objective. Produce merged instructions P' that:

- preserve stable, broadly useful guidance from P^* (the “growing playbook”),
- incorporate *new* and *validated* guidance from P^d *additively*,
- remain concise and non-redundant (avoid restating the same rule twice),
- support incremental toolset growth: keep cross-cutting rules stable while adding any new decision rules introduced by newly appended tools.

Merge Rules (Strict)

- **Do not overwrite:** never delete a rule from P^* unless P^d provides a clearly conflicting correction.
- **Prefer generalization:** if P^d adds a rule that generalizes an existing one, keep the generalized version and remove the narrower duplicate.
- **Resolve conflicts explicitly:** if a draft rule contradicts an existing rule, keep the version that is more precise and operational (clear triggers, clear expected behavior), and remove the other.
- **Tool-growth compatibility:** if a draft rule is specific to a newly added tool, include it only if it can be stated as a general decision rule (when-to-use / how-to-fill), otherwise keep it minimal and non-invasive.
- **No tool merging:** do not rename, alias, or merge tools; only adjust global instruction text.

Output Format (Strict)

- Return *only* the merged global instructions P' as plain text.
- No JSON. No commentary. No additional sections.

Figure 6: **MergeWithBest prompt.** The merger composes rollout-specific draft instructions P^d with the current global best P^* to form P' , preserving stable cross-cutting guidance while integrating validated new rules. This “growing playbook” mechanism supports incremental toolset expansion by keeping existing conventions stable and appending new decision rules when new tools are introduced.

Prompt Template: Slot-Semantics Globalization (Two-Level Context Editing)

Role. You are a context editor for a tool-using LLM agent. You may revise (i) *Global Instructions* P and (ii) per-tool schemas $\{T_i\}$ (tool and argument descriptions).

Objective. Reduce repeated slot/argument guidance across tools while preserving tool-specific distinctions needed for correct tool selection and slot filling.

Step 1: Scan for repeated slot semantics. Read each tool schema T_i and its argument descriptions carefully. Identify *recurring* slot conventions that appear across many tools, such as: date/time windows, identifier formatting, numeric bounds (inclusive/exclusive), units/currency normalization, boolean/defaulting rules, pagination parameters, and sorting conventions.

Step 2: Globalize shared rules. For each repeated convention, write a *single*, canonical rule in the Global Instructions P that: (a) states the default interpretation and formatting requirements, and (b) specifies when to apply default values versus using user-provided constraints. The global rule should be phrased generically so it applies to any tool that contains the relevant slot(s).

Step 3: Keep exceptions local. Do *not* merge, alias, or rename tools. For each tool schema T_i :

- Remove redundant restatements of globalized rules and replace them with a short pointer (e.g., “See Global Instructions: [Rule Name]”).
- If a tool requires different semantics (e.g., a different date format, special rounding, a stricter constraint), keep that information *locally* in T_i and explicitly

mark it as an **override** of the global rule.

Constraints.

- Do not change tool interfaces: do **not** add/remove arguments or invent fields.
- Prefer minimal, high-impact edits: globalize only clearly repetitive conventions; keep tool-unique decision rules and edge cases local.

Output.

- An updated Global Instructions block to append to P (named rules + concise definitions).
- Updated schemas for only the tools you modified (short pointers + explicit overrides).

780

B Enterprise Tool-Inventory Dataset (ETID): Release Status and Disclosure Constraints

781

782

783

Dataset overview. We introduce the *Enterprise Tool-Inventory Dataset (ETID)* to stress-test tool-use under complex, multi-argument schemas and realistic enterprise-style tool catalogs. ETID is designed to contain *no personally identifiable information (PII)*. Explicit PII-avoidance constraints were applied during generation and further *verification* was done through human review to confirm that examples do not contain sensitive personal content or proprietary identifiers.

784

785

786

787

788

789

790

791

792

793

Legal review and current disclosure limitations. At the time of submission, ETID is undergoing an internal legal review. As a result, we are not yet able to release the dataset or disclose certain concrete artifacts beyond the aggregate statistics and evaluation results already reported in the paper. This constraint also limits the motivating examples we can provide for the *Globalizing Repetitive Slot Semantics* step: while we empirically demonstrate the effect and characterize the repetition patterns (e.g., via parameter-frequency analysis), we omit verbatim examples that could inadvertently expose dataset-specific schema fragments.

794

795

796

797

798

799

800

801

802

803

804

805

806

Planned release. We are actively working with the legal team toward a version of ETID that can be shared publicly. Subject to review outcomes,

807

808

809

we intend to release an appropriately redacted and documented version (or a functionally equivalent public variant) that supports reproducibility.

C Experimental Setup

Evaluation protocol and reporting. To reduce variance from stochastic decoding and optimization dynamics, all results are aggregated over multiple independent runs. Unless otherwise stated, each experiment is repeated **5–10 times** and we report the **average** performance across runs. For fair comparison, we run **JTPRO and GEPA under matched optimization budgets**: both methods use the same maximum number of rollouts and identical optimization settings (including the same minibatch size and the same reflector configuration). Across all experiments, we use **gpt-o3-mini** as the reflector model and set the **LLM temperature parameter to 1**.

D Additional Figure Discussion

D.1 Figure 7: Repetitive slot semantics across tools

Figure 7 motivates GLOBALIZESLOTS. Panel (b) shows a heavy-tailed frequency distribution: a small number of parameter families (notably identifiers and date/time) recur across a large portion of the tool inventory (up to 77/124 tools). Panel (a) illustrates that these recurring families are often described with near-verbatim text (e.g., ISO formatting, inclusive bounds, defaulting behavior), which increases context length without improving tool-specific disambiguation. JTPRO therefore *lifts* shared slot conventions into the global instruction layer (reducing repeated schema text) while preserving tool-local exceptions and decision rules in T_i'' to maintain accurate disambiguation and slot filling.

D.2 Figure 8: ToolACE scaling results

Figure 8 evaluates robustness under tool-universe growth (500 vs. 1000 tools). The dominant failure mode under scaling is reduced TSA: as inventories expand, overlapping tool descriptions and increased distractors cause more routing errors, which then cascade into lower OSR. GEPA partially mitigates this via global instruction refinement, but it does not directly repair tool-local ambiguity or slot semantics. JTPRO delivers the most consistent OSR gains because it **jointly** revises global policies *and* the specific tool/slot descrip-

tions implicated by observed failures, improving both selection and downstream argument correctness.

D.3 Figure 9: ETID performance across supervision levels

Figure 9 studies data efficiency on ETID under Train-1ex/2ex/4ex regimes. Across models, baselines often achieve relatively strong TSA but substantially lower OSR, indicating that **slot/value instantiation** is the primary bottleneck under complex schemas. JTPRO improves SFA (conditional on TSA) and therefore consistently lifts OSR across supervision levels, reflecting that many ETID failures stem from underspecified or inconsistent argument semantics that can be corrected through targeted tool/slot documentation edits plus strengthened global tool-calling rules.

D.4 Figure 10: Per-example slot-filling improvement rate

Figure 10 reports the fraction of test instances for which JTPRO improves per-query slot/value correctness over the baseline. The gains are larger on ETID (complex schemas) than on ToolACE, aligning with the hypothesis that real-world OSR is often bottlenecked by argument instantiation even after correct tool selection. Importantly, the improvements occur on a non-trivial fraction of held-out examples across all evaluated models, suggesting that joint context refinement yields robust, example-level corrections rather than isolated wins.

D.5 Figure 11: ToolACE-500 example-wise corrections (GPT-5)

Figure 11 provides an example-wise view of slot/value corrections on ToolACE-500 for GPT-5. Each bar corresponds to a test instance, contrasting baseline vs. JTPRO slot correctness. Overall, JTPRO improves slot correctness on **26/121** examples (**21.48%**). This plot highlights that improvements are distributed across the test set (rather than concentrated in a single cluster), consistent with JTPRO correcting recurring slot conventions and tool-specific documentation ambiguities that manifest in diverse queries.

D.6 Figure 12: ETID example-wise corrections (GPT-5)

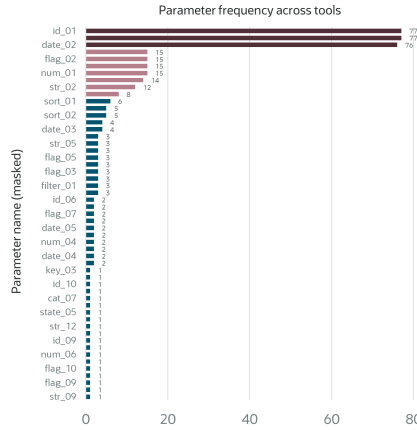
Figure 12 shows the analogous example-wise comparison for ETID (GPT-5). JTPRO improves slot correctness on **94/403** examples (**23.33%**), rein-

```

"startDate": {
  "type": "string", "format": "date-time",
  "description": "Specifies the inclusive
start date and time (ISO 8601, e.g., 2024-
01-01T00:00:00Z) for the calculation window.
Provide this when an explicit time range
is requested; otherwise default to the
current fiscal year start (e.g., 2024-07-
01T00:00:00Z)."}
"endDate": {
  "type": "string", "format": "date-time",
  "description": "Specifies the inclusive
end date and time (ISO 8601, e.g., 2024-
12-31T23:59:59Z) for the calculation window.
Provide this when an explicit time range
is requested; otherwise default to the
current fiscal year end (e.g., 2024-12-
31T23:59:59Z)."}

```

(a)



(b)

Figure 7: **Parameter frequency across tools reveals highly repetitive slot semantics.** The distribution is heavy-tailed: a small number of slot families (e.g., date/time and identifiers) recur in a large fraction of tools (up to 77/124), motivating globalization of shared slot conventions to reduce duplicated schema text.

forcing that complex multi-argument schemas benefit strongly from (i) tightening global tool-calling policies (e.g., required-field completeness, no hallucinated keys) and (ii) clarifying per-tool parameter semantics. Together with Figures 9 and 10, this example-level view supports the claim that SFA is a major driver of end-to-end OSR gains on ETID.

D.7 Figures and Plots

D.8 Figure 13: Convergence behavior

Figure 13 plots validation OSR over JTPRO iterations for three base models, with \star denoting the final test OSR obtained using the best validation-selected context. The curves show rapid early gains followed by saturation, consistent with the reflector first correcting high-impact, systematic errors (e.g., frequent tool confusions, missing required slots, formatting/defaulting mistakes) and later iterations focusing on smaller refinements. The separation between validation trajectories and the final test markers indicates that improvements transfer to held-out queries rather than merely optimizing minibatch idiosyncrasies.

D.9 Tool Description Disambiguation Analysis

We study how joint optimization improves tool selection by analyzing semantic changes in tool descriptions on the ToolAce-500 benchmark.

D.9.1 Description Enrichment

The optimizer modifies 55 out of 500 tool descriptions (11%), increasing the average description length from 86.1 to 100.1 characters (+16.3%). These edits primarily target *disambiguation*, ex-

PLICITLY differentiating tools with overlapping semantics.

Example: search vs. web_search

- **search (before):** “Perform Google search and get results.”
- **search (after):** “Perform Google search with advanced locale controls (gl/hl), country restrictions (cr), and time filters (tbs). *NOT for general web article or paper discovery—prefer web_search for generic queries.*”
- **web_search (after):** “Search the web for relevant pages. *PREFERRED for general-purpose web, article, and paper discovery. Do not confuse with similarly named search tools.*”

D.9.2 Confusable Tool Groups

We identify tools sharing common name prefixes (e.g., `get_user_*`, `get_all_*`) as potentially confusable. This yields 37 groups comprising 109 tools, representing high-risk ambiguity regions where models frequently select semantically similar but incorrect tools.

D.9.3 Embedding-Based Disambiguation Metric

To quantify disambiguation quality, we compute pairwise cosine similarity between tool descriptions within each confusable group using sentence embeddings (all-MiniLM-L6-v2). Lower intra-group similarity indicates stronger semantic separation.

The `get_ip_*` group exhibits the largest improvement, with similarity reduced from 0.668 to 0.502.

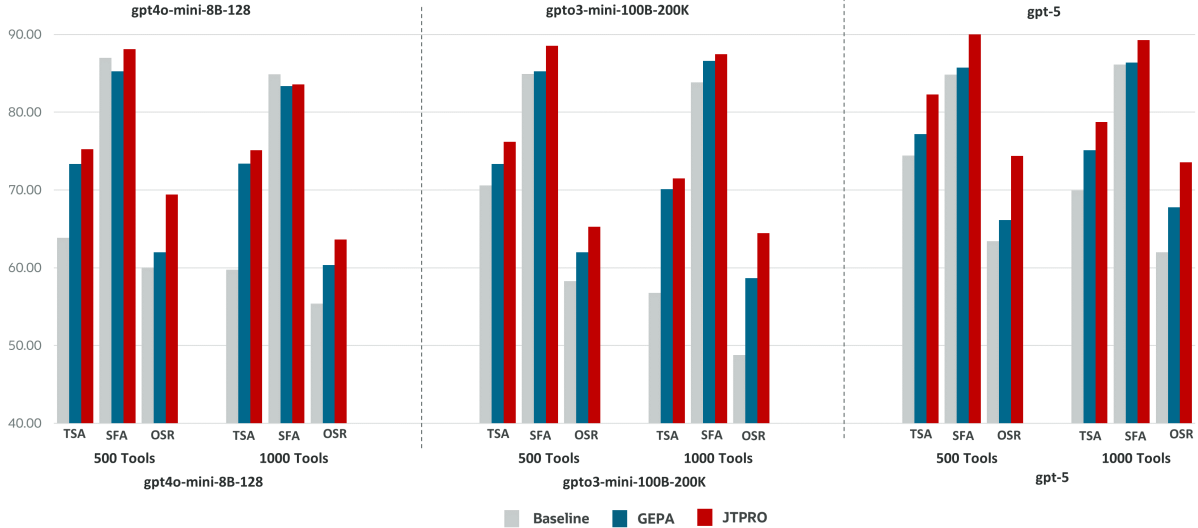


Figure 8: **ToolACE scaling results across models and metrics.** For each model, we report TSA, SFA (conditional on correct tool), and OSR at 500 and 1000 tools. Tool-universe growth primarily reduces TSA for the baseline, which cascades to lower OSR; GEPA partially mitigates this via global instruction refinement, while JTPRO provides the most consistent improvements in OSR by jointly refining global instructions and tool/slot descriptions.

| Tool Group | Before | After | Δ |
|---------------------|--------|-------|----------|
| get_ip_* (2) | 0.668 | 0.502 | -0.166 |
| get_page_* (2) | 0.849 | 0.736 | -0.113 |
| get_languages_* (2) | 0.676 | 0.584 | -0.092 |
| get_trending_* (3) | 0.437 | 0.377 | -0.060 |
| search_by_* (2) | 0.281 | 0.265 | -0.016 |

Table 5: Intra-group cosine similarity (lower indicates better disambiguation) for the most improved confusable tool groups. Parentheses denote group size.

| Metric | Value |
|---------------------------|-------------|
| Total tools analyzed | 500 |
| Modified descriptions | 55 (11.0%) |
| Avg. length (before) | 86.1 chars |
| Avg. length (after) | 100.1 chars |
| Relative increase | +16.3% |
| Confusable groups | 37 |
| Confusable tools | 109 |
| Groups improved | 6 (16.2%) |
| Avg. similarity reduction | 0.012 |

Table 6: Summary of tool description disambiguation on ToolAce-500.

This change reflects the addition of explicit preference guidance, e.g., “Preferred for general IP geolocation requests. Use this instead of get_geolocation_by_ip unless extended fields are required.”

D.9.4 Disambiguation Patterns Learned

Across the 55 modified tools, we observe four recurring disambiguation strategies:

- Parameter format guidance** (28 tools): e.g., “Pass country_code as a 2-letter lowercase ISO code.”
- Explicit preference signals** (9 tools): e.g., “PREFERRED for...”
- Negative constraints** (5 tools): e.g., “NOT for general web article discovery.”
- Cross-tool references** (3 tools): e.g., “Use this instead of get_geolocation_by_ip.”

D.9.5 Summary Statistics

Overall, joint optimization learns to resolve tool ambiguity through targeted description edits, complementing instruction-level optimization and improving tool selection robustness.

D.9.6 Visualizations

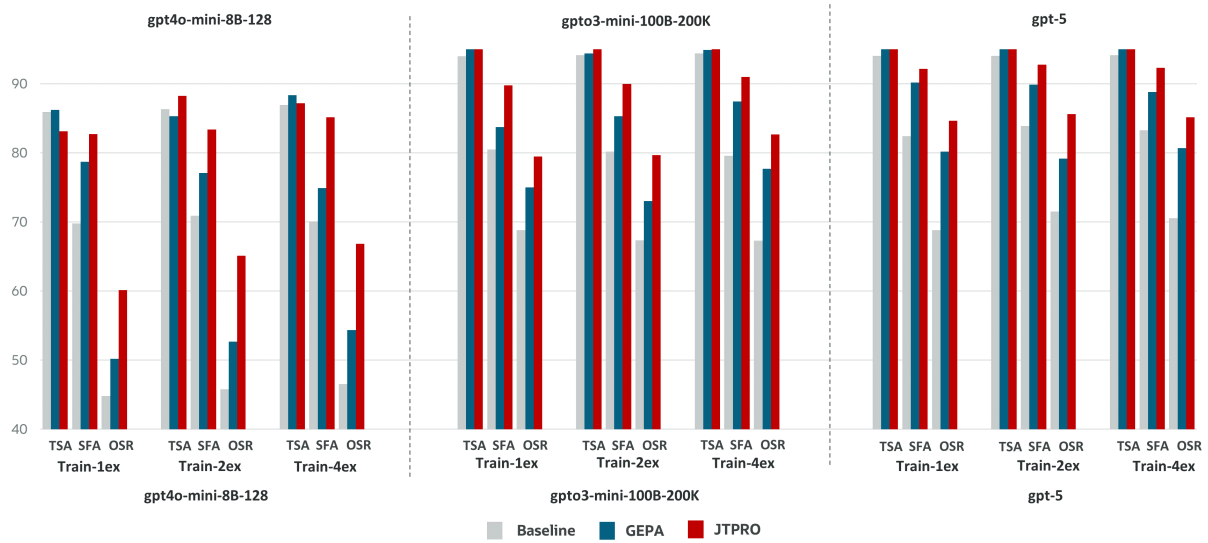


Figure 9: **ETID performance across supervision levels.** Grouped bars show TSA, SFA (conditional on TSA), and OSR for three models under Train-1ex/2ex/4ex regimes. Baselines achieve high TSA but substantially lower OSR, revealing slot/value errors as the dominant failure mode; JTPRO improves SFA and therefore OSR consistently across regimes, while GEPA primarily improves TSA for larger models.

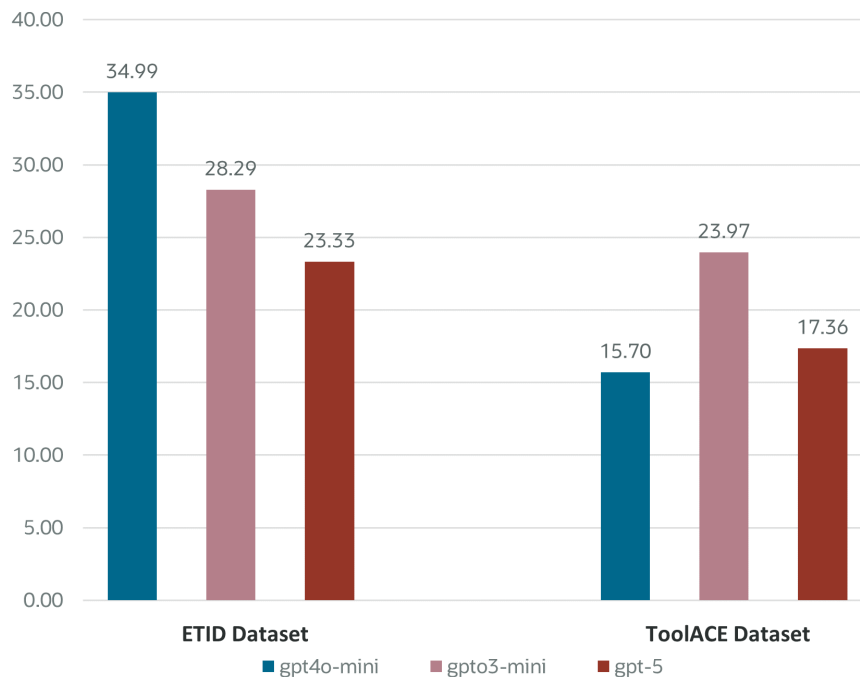


Figure 10: **Per-example slot-filling improvements from JTPRO.** For each base model and dataset, we report the *average percentage of test instances* on which slot filling is more accurate after JTPRO optimization than the corresponding baseline (i.e., per-query slot/value correctness improves). Gains are larger on the complex slot-filling ETID benchmark (e.g., 34.99% for gpt4o-mini) and remain substantial on ToolACE (e.g., 23.97% for gpt-o3-mini), indicating that JTPRO improves argument instantiation on a non-trivial fraction of held-out queries across models.

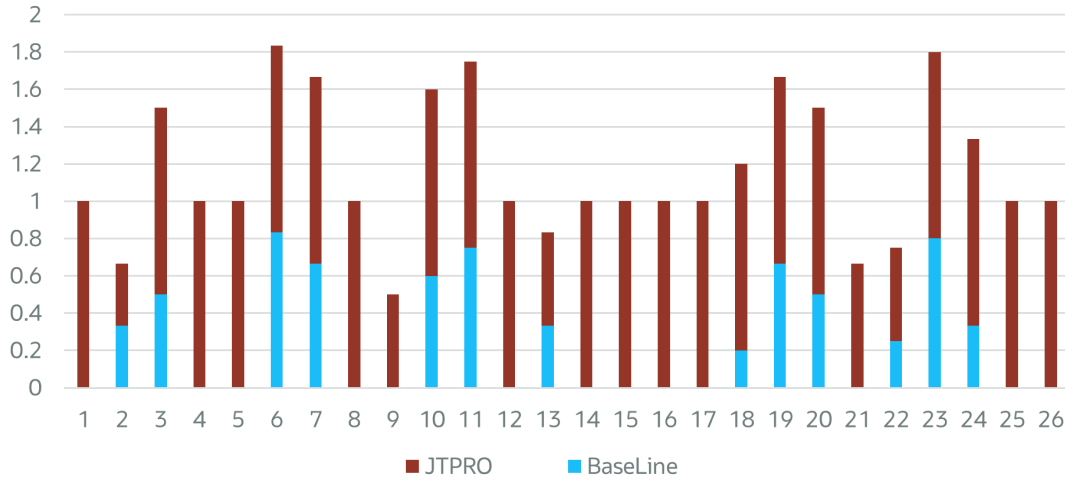


Figure 11: **Per-example slot/value corrections after JTPRO (ToolACE-500, GPT-5)**. Example-wise comparison of slot-filling outcomes on the ToolACE test set with 500 tools for GPT-5: each bar corresponds to a test instance (x-axis indices), highlighting instances where JTPRO fixes previously incorrect slot/value instantiations relative to the baseline. Overall, JTPRO improves slot correctness on **26 out of 121** test examples (**21.48%**).

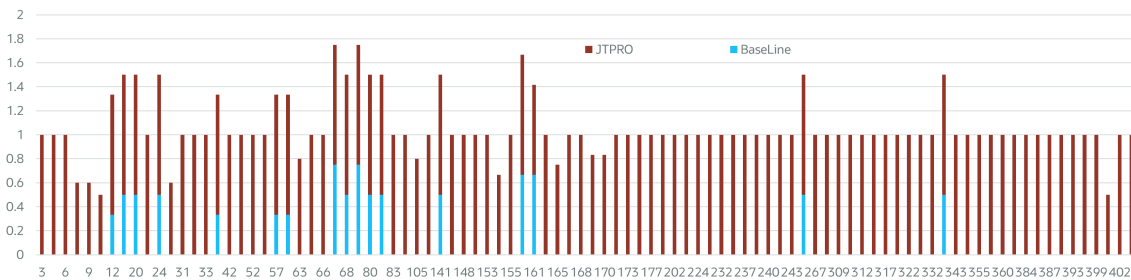


Figure 12: **Per-example slot/value corrections after JTPRO (ETID, GPT-5)**. Example-wise comparison of slot-filling outcomes on the ETID test set for GPT-5: each bar corresponds to a test instance (x-axis indices), highlighting instances where JTPRO fixes previously incorrect slot/value instantiations relative to the baseline. Overall, JTPRO improves slot correctness on **94 out of 403** test examples (**23.33%**).

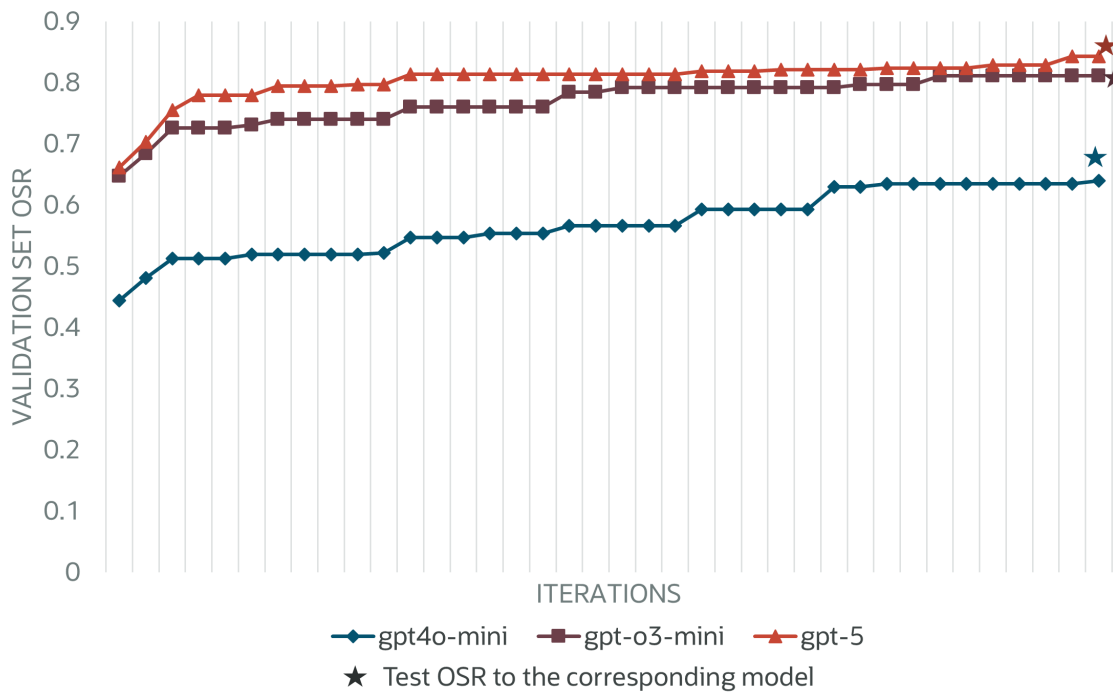


Figure 13: **JTPRO convergence on validation OSR.** Validation-set Overall Success Rate (OSR) as a function of JTPRO optimization iterations for three base models (gpt4o-mini, gpt-o3-mini, gpt-5). The \star markers denote the final OSR measured on the corresponding held-out test set using the best validation-selected context. The curves show rapid early improvements followed by diminishing returns, indicating that most gains are captured within the first few refinement cycles while later iterations primarily yield incremental consolidation. Across all three models, OSR increases sharply in early iterations and then saturates, consistent with the reflector quickly correcting high-impact failure modes (e.g., tool confusions and missing slot constraints) before converging to smaller, fine-grained edits. The final \star markers report test OSR achieved by the validation-selected best context, showing that the improvements learned during optimization transfer to held-out queries rather than overfitting to the optimization batches.

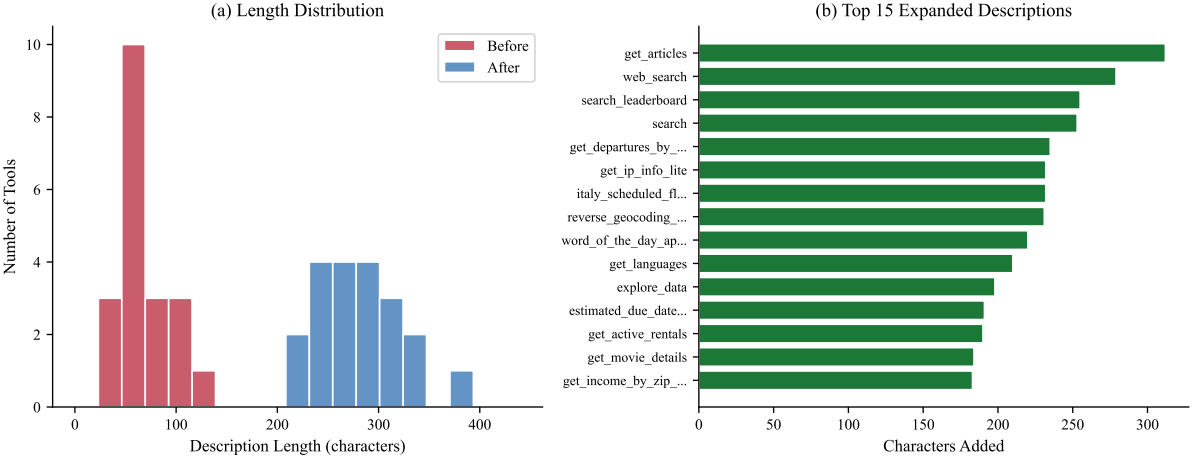


Figure 14: Tool description length analysis. (a) Distribution of description lengths before and after optimization. (b) Per-tool length changes for the 55 modified tools.

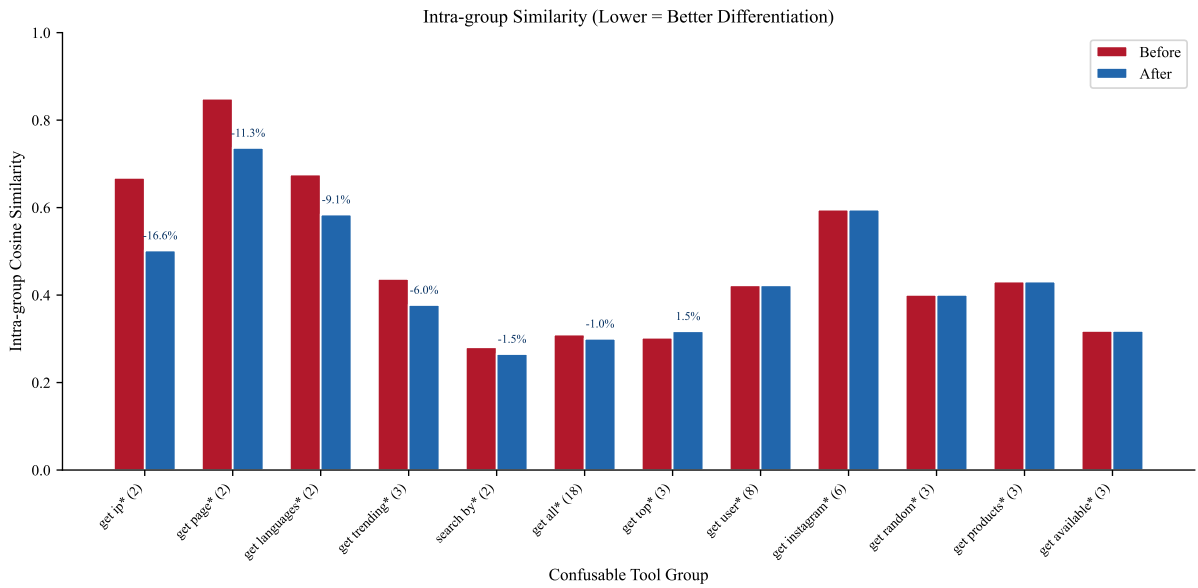


Figure 15: Intra-group cosine similarity for the top 15 confusable tool groups. Lower values indicate stronger semantic differentiation.

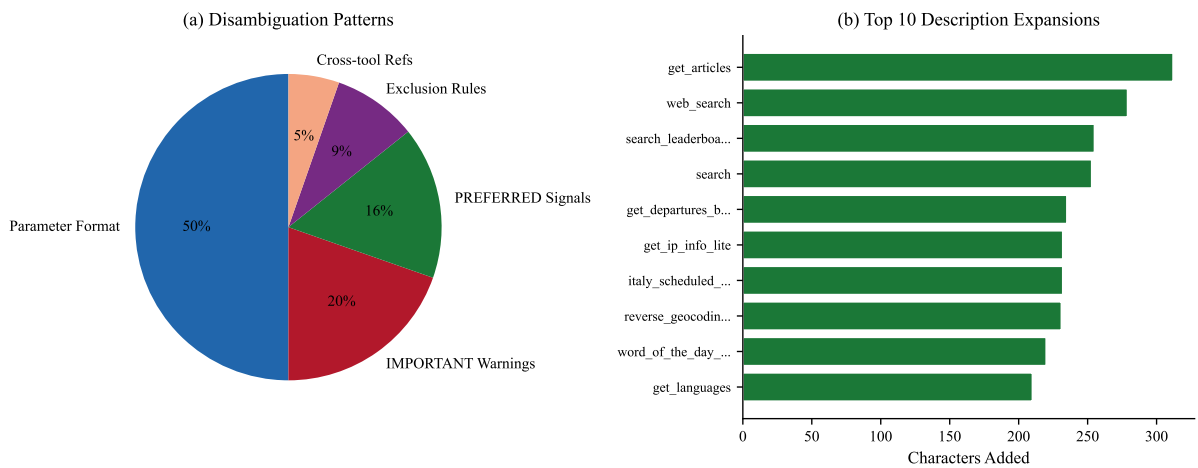


Figure 16: Disambiguation strategies learned by the optimizer. (a) Distribution of pattern types. (b) Tools with the largest description expansions.