# 🐕 Samoyed: Towards Generalized LLM Agents via Fine-Tuning on 50000+ Interaction Trajectories

**Anonymous ACL submission**

## Abstract

Fine-tuning on agent-environment interaction trajectory data is of high potential for surfacing generalized agent capabilities in open-source large language models (LLMs). In this work, we curate SuperAgent, by far the largest trajectory tuning data collection featuring more than 50k diverse high-quality interaction trajectories with GPT-annotated rationale. SuperAgent comprises 16 tasks covering five distinct agent skill dimensions. Furthermore, we present Samoyed[1], a series of open-source LLMs fine-tuned on SuperAgent. Our comparative experiments show that Samoyed outperforms strong baseline LLMs on both held-in and held-out tasks, demonstrating the effectiveness of scaling the interaction trajectory data to acquire generalized agent capabilities. Additional studies also reveal some key observations regarding trajectory tuning and agent skill generalization.

## 1 Introduction

An agent is an entity that possesses the ability to exercise volition, make choices, take actions, and, most importantly, perceive its environment (Jennings et al., 1998). In the realm of cognitive science, previous literature has suggested that interaction with environment derives an agent's generalized intelligence, and intelligent behavior emerges from a synergistic blend of simpler behaviors, including reasoning, programming, and game playing (Brooks, 1991). Modern large language models such as GPT-3.5 (OpenAI, 2022) and GPT-4 (OpenAI, 2023) have demonstrated strong capabilities in instruction following, reasoning, and planning, which encourage many attempts to build autonomous agent systems utilizing LLMs as core controllers (Richards, 2023; Nakajima, 2023).

However, comprehensive evaluations have shown that most open-sourced LLMs fall short in agent capabilities when compared with GPTs (Liu et al., 2023; Wang et al., 2023).

Previous research has shown that learning from golden interaction trajectories (we name it as **Trajectory Tuning**) could enhance the capabilities of weaker agents (Brooks, 1991; Hussein et al., 2017). However, many such studies heavily focus on specialized agents for specific tasks, and the enhancement of generalized agent capabilities starts to draw attention only recently. Existing attempts are exampled by Chen et al. (2023) and Yin et al. (2023), who build agent trajectory data from teacher agents (*e.g.*, GPT-4) and fine-tune open-source LLMs to improve specific agent abilities like reasoning. Taking a step further, Zeng et al. (2023) adopt a multi-task tuning approach called AgentTuning. While training on a small trajectory dataset comprising six tasks with 1.8k trajectories, AgentTuning struggles to enhance the generalized agent capability, particularly for the 7B and 13B LLMs[2].

Inspired by the data scaling effect in instruction tuning (Chung et al., 2022; Wang et al., 2022b), we construct a much larger trajectory tuning dataset to further explore the efficacy of incorporating interaction trajectory data on agent ability generalization. To the best of our knowledge, SuperAgent is the largest agent interaction trajectory dataset so far, which features 16 distinct tasks covering five agent skill dimensions and contains 50k+ trajectories with high-quality chain-of-thought (CoT) rationale for each action step. We adopt a hybrid data construction pipeline combined with heuristic action searching and LLM-based rationale generation. Compared with previous work that collects successful trajectories of GPTs to build the training data (Chen et al., 2023; Zeng et al., 2023), Super-

---

[2] The experiments in Zeng et al. (2023) show that models trained on only interaction trajectories overfit and suffer performance degradation on unseen tasks.

AGENT has extraordinary quality and suffers less from the *difficulty bias* problem[3].

We further develop SAMOYED, a set of models that possess enhanced agent capabilities, by trajectory tuning Llama-2 (Touvron et al., 2023) on SUPERAGENT. Evaluation on both held-in and unseen held-out tasks suggests that by fine-tuning on extensive multi-task trajectories, our models exhibit remarkable agent intelligence in comparison with untuned ones. Specifically, SAMOYED outperforms GPT-3.5-Turbo on average on held-in tasks, which can be attributed to the in-domain trajectory tuning. Furthermore, our models also demonstrate superior performance on held-out tasks, manifesting the efficacy of large-scale trajectory tuning in acquiring generalized agent capabilities.

To trace the emergence of agent capabilities generalization, we follow the initial evaluation with a systematic analysis across various dimensions. First, we validate the effectiveness of massive trajectory tuning on different base models. Next, we conduct an ablation study by combining generalist instruction data and code data to examine the benefits of hybrid training, which reveals further enhancements in the agent capabilities. Lastly, we plot the scaling trends of tasks and the number of trajectories, while also studying the ability transfer across different agent skills.

Our contributions are summarized as follows:

- We build SUPERAGENT, a high-quality agent interaction trajectory dataset with 50k+ trajectories covering 16 tasks across five skill dimensions.

- We train SAMOYED, the most powerful open-source LLM set specialized at agent tasks at 7B/13B scales, by trajectory tuning. Extensive experiments show SAMOYED can achieve transferable agent intelligence on unseen tasks.

- We further conduct comprehensive experiments and in-depth analysis on agent intelligence acquisition, including the relations with instruction following and code capability, scaling law of interaction trajectories, and generalization of different agent skills.

- To facilitate future research on developing open-source LLM agents, we release SUPERAGENT, SAMOYED, and the unified evaluation framework.

---

[3]See Section 3 for details.

| | SUPERAGENT (this work) | FireAct (Chen et al., 2023) | AgentInstruct (Zeng et al., 2023) |
|---|---|---|---|
| Task Num. | **16** | 3 | 6 |
| Instruct. Num. | **51287** | 1344 | 1866 |
| Avg. Turns | 3.9 | - | **5.2** |
| Reasoning | ✓ | ✓ | ✗ |
| Math | ✓ | ✗ | ✗ |
| Programming | ✓ | ✗ | ✓ |
| Web | ✓ | ✗ | ✓ |
| Embodied AI | ✓ | ✗ | ✓ |

Table 1: A comparison of SUPERAGENT with other datasets for agent trajectory tuning.

## 2 Related Work

### 2.1 Instruction Tuning

Instruction tuning is a simple yet powerful approach to align LLMs with human preferences (Zhang et al., 2023). Previous studies have primarily focus on improving general-purpose instruction following capabilities of LLMs. FLAN series (Wei et al., 2021; Chung et al., 2022), T0 (Sanh et al., 2021), and NaturalInstruction (Wang et al., 2022b) scale up the instruction datasets to activate the generalized instruction following capabilities of LLMs. More recently, utilizing synthetic instruction following data distilled from GPTs to align open-source LLMs has also been proposed (Taori et al., 2023; Chiang et al., 2023). Furthermore, multiple works have shown the promise of instruction tuning in enhancing the specialized abilities of LLMs, such as math (Yu et al., 2023; Yue et al., 2023), reasoning (Lee et al., 2023), and agent tasks (Chen et al., 2023; Zeng et al., 2023).

### 2.2 LLM-based Agent

Modern LLMs have demonstrated various emergent abilities that encourage researchers to build agent systems based on LLMs. ReAct (Yao et al., 2022b) combines CoT reasoning with agent actions to accomplish tasks such as QA. Auto-GPT (Richards, 2023) and BabyAGI (Nakajima, 2023) harness LLMs as the core controllers to constitute powerful agent frameworks capable of solving real-world complex problems. While advanced proprietary models exampled by GPT-3.5/4 have shown strong performances on agent tasks, their open-source counterparts still lag far behind (Liu et al., 2023; Wang et al., 2023). In response, recent studies including FireAct (Chen et al., 2023), AgentTuning (Zeng et al., 2023) and Lumos (Yin et al., 2023) collect agent trajectory data from teacher agents (*e.g.*, GPT-4) and fine-tune open-
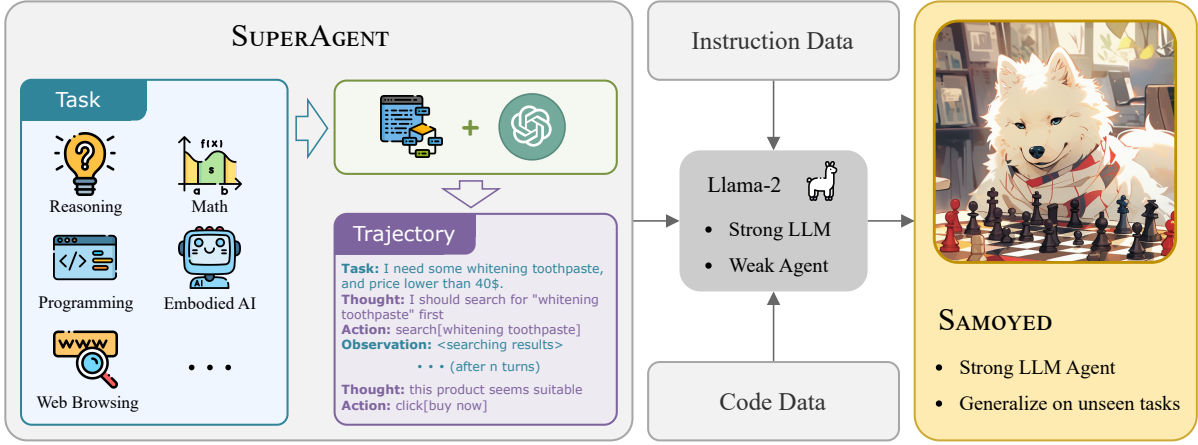
Figure 1: Overview of the construction process of SUPERAGENT and the training procedure of SAMOYED

source LLMs (*e.g.*, Llama series) with the data. However, limited by the number of tasks and expert trajectories, existing research has not yet exhaustively explored whether open-source LLMs can acquire generalized agent abilities, a gap that this study aims to bridge.

## 3 Method

### 3.1 Agent Task Formulation

Given an agent task described by the instruction $u$, an LLM agent generates an action $a_1$ based on its policy. Next, an environment receives the action, transfers to a new latent state, and provides an observation $o_i$ in natural language format. Subsequently, the agent generates another action for the next step, $a_{i+1}$, and repeats this circle of interaction with the task environment until either the task is completed or the maximum number of steps is reached. This "conversation" between the agent with the environment is denoted as the interaction trajectory $(u, a_1, o_1, ..., a_n)$. Finally, a final reward $r \in [0, 1]$ is returned depending on the task completion status.

In this work, we employ ReAct (Yao et al., 2022b) as the agent tasking framework, which outputs Chain-of-Thought (CoT) rationale before the action. CoT method (Wei et al., 2022) is an effective approach to enhance the inferential capabilities of LLMs by a step-by-step reasoning process.

### 3.2 Task and Instruction Collection

A generalized agent needs to possess a wide range of capabilities across various dimensions. As shown in Table 2, we curate 16 publicly available agent datasets to lay the foundation of SUPERAGENT and categorize specific tasks into five

skill dimensions: reasoning, math, programming, web navigation, and embodied tasks. Additionally, some tasks aggregated in SUPERAGENT involve the usage of external tools, such as search engine, calculator, and code interpreter, as the ability to effectively operate tools is also a crucial aspect for generalized agents. From the perspective of action space, tasks in SUPERAGENT can be classified into two types: those with a *continuous* action space (including natural language and code) and those with a predefined *discrete* action space. Our dataset also covers a broad range of interaction turns, ranging from 1 to 30. Note that some tasks are originally evaluated in a single-turn QA style, such as HotpotQA (Yang et al., 2018) and MATH (Hendrycks et al., 2021). Following Wang et al. (2023), we modify these datasets to accommodate multi-turn interaction environments with tool usage.

Since most of the original benchmarks have a training set, we use them to construct our dataset. To balance data sources, we down-sample some tasks which have a huge training set. See Appendix A for detailed descriptions of each dataset.

### 3.3 Difficulty Bias in Trajectory Collection

Previous works (Chen et al., 2023; Zeng et al., 2023) have employed GPT-4 as teacher agents to interact with the environment and collect successful interaction trajectories. To ensure the quality of generated data, a trajectory filtering mechanism is used to remove the cases where GPT failed. However, this **GPT-exploration** pipeline automates the trajectory construction at some significant cost. Scaling up this process to a larger trajectory amount is challenging due to the low success rate of GPT-4. For instance, AgentInstruct (Zeng et al., 2023)

| Skill Dim. | Task | Action Space | Tool | #Inst. | Avg. Turns | Traj. Annotation |
|---|---|---|---|---|---|---|
| Reasoning | HotpotQA (Yang et al., 2018) | Continuous | Search | 4273 | 3.1 | Explore |
| | StrategyQA (Geva et al., 2021) | Continuous | Search | 1267 | 3.6 | Explore |
| | TriviaQA (Joshi et al., 2017) | Continuous | Search | 4134 | 2.5 | Explore |
| Math | GSM8K (Cobbe et al., 2021) | Continuous | Calculator | 7471 | 4.5 | Reformat |
| | MathQA (Amini et al., 2019) | Continuous | Python | 4000 | 2.0 | Explore |
| | MATH (Hendrycks et al., 2021) | Continuous | Python, Wiki | 2312 | 2.5 | Explore |
| Programming | IC-SQL (Yang et al., 2023) | Continuous | MySQL | 4540 | 4.8 | Explore+Answer Force |
| | APPS (Hendrycks et al., 2021) | Continuous | Python | 4408 | 1.0 | Reformat |
| | HumanEval (Chen et al., 2021) | Continuous | Python | 134 | 2.7 | Explore+Answer Force |
| | MBPP (Austin et al., 2021) | Continuous | Python | 608 | 2.2 | Explore+Answer Force |
| Web | Mind2Web (Deng et al., 2023) | Discrete | - | 7770 | 1.0 | Reformat |
| | WebArena (Zhou et al., 2023) | Discrete | - | 657 | 1.0 | Reformat |
| | WebShop (Yao et al., 2022a) | Discrete | - | 5315 | 3.4 | Explore & Reformat |
| Embodied | ALFWorld (Shridhar et al., 2020b) | Discrete | - | 3554 | 10.1 | Reformat |
| | RoomR (Weihs et al., 2021) | Discrete | - | 300 | 30.2 | Search+Reformat |
| | IQA (Gordon et al., 2018) | Discrete | - | 1627 | 28.4 | Search+Reformat |
| | Total (SUPERAGENT) | - | - | 51287 | 3.9 | - |

Table 2: Overview of SUPERAGENT dataset. It compiles 16 agent tasks covering 5 skill dimensions, formulating the largest interaction trajectory dataset. "Inst." and "Traj." refer to instruction and interaction trajectory.

discards more than 90% generated trajectories due to GPT failures.

More importantly, we argue that GPT-exploration pipelines inevitably introduce *difficulty bias* to the final training data. Essentially, a trajectory filtering strategy can be regarded as grouping the instances based on whether GPT is capable of solving them. In other words, discarding failed trajectories leads to a skewed distribution of "difficulty", resulting in a training set with much easier instances than those in the test set. This violation of the *i.i.d.* assumption may hurt the generalization ability of the trained agents. We conduct a diagnostic experiment on the difficulty bias and the results support our arguments (see Appendix B).

### 3.4 Interaction Trajectory Annotation

In response to the difficulty bias we observe and to build an unbiased large-scale dataset, we develop a trajectory annotation process incorporating both heuristic and LLM-based methods. Specifically, tailored to the specific nature of different tasks, our approach involves several techniques to obtain high-quality action sequences along with their corresponding rationales accordingly.

**Answer Forcing** For tasks with a continuous natural language or code action space, we apply an answer forcing re-annotation strategy as an extension to GPT-exploration pipeline, to mitigate the difficulty bias introduced by simply filtering. Instead of discarding the failed trajectories, we prompt GPT with the failed trajectory and the gold final answer to generate a correct interaction trajectory.

**Heuristic Action Search** Meanwhile, we use heuristic depth-first-search algorithm to find the golden action sequence for tasks with a discrete action space such as Rearrange (Weihs et al., 2021) and IQA (Gordon et al., 2018),

**Trajectory Reformat** For tasks with gold action sequences, either provided in original datasets or generated by heuristic action search, we directly prompt GPT to reformat them and generate the corresponding CoT rationale of each action step.

For tasks with a huge number of instructions and GPTs have a high success rate, such as StrategyQA (Geva et al., 2021) and WebShop (Yao et al., 2022a), we directly use the GPT-exploration pipeline as Zeng et al. (2023).

Due to the large number of instructions and our budget constraints, we employ GPT-3.5-Turbo as the primary LLM in the interaction trajectory annotation process. The overview of SUPERAGENT is shown in Table 2. See the Appendix A for more details about the annotation process of each task. A human evaluation assessing the quality of our dataset can be found in Appendix C.

### 3.5 Train SAMOYED with SUPERAGENT

To initialize the training of SAMOYED, we formulate agent interaction trajectories in

| Task | Skill Dim. | #Inst. | Metric |
|------|-----------|--------|--------|
| *Held-in Tasks* | | | |
| HotpotQA (Yang et al., 2018) | Reasoning | 100 | Exact Match |
| StrategyQA (Geva et al., 2021) | Reasoning | 100 | Exact Match |
| GSM8K (Cobbe et al., 2021) | Math | 100 | Exact Match |
| MATH (Hendrycks et al., 2021) | Math | 100 | Exact Match |
| IC-SQL (Yang et al., 2023) | Programming | 100 | Avg. Reward |
| MBPP (Austin et al., 2021) | Programming | 100 | Success Rate |
| Mind2Web (Deng et al., 2023) | Web | 1173 | Step SR |
| WebShop (Yao et al., 2022a) | Web | 200 | Avg. Reward |
| ALFWorld (Shridhar et al., 2020b) | Embodied | 134 | Success Rate |
| *Held-out Tasks* | | | |
| Bamboogle (Press et al., 2022) | Reasoning | 126 | Exact Match |
| SVAMP (Patel et al., 2021) | Math | 100 | Exact Match |
| IC-Bash (Yang et al., 2023) | Programming | 200 | Avg. Reward |
| MiniWoB++ (Kim et al., 2023) | Web | 460 | Success Rate |
| ScienceWorld (Wang et al., 2022a) | Embodied | 270 | Avg. Reward |

Table 3: The held-in and held-out tasks used to evaluate the agent capabilities of different LLMs.

SUPERAGENT into a chatbot-style schema $(u, a_1, o_1..., a_i, o_i, ..., a_n)$, where $u$ is the task instruction, $o_i$ and $a_i$ denote the observation from the task environment and the corresponding action with rationale generated by the agent in the $i$-th round. During the training process, we feed the entire interaction trajectory into a decoder-only LLM, where only the auto-regressive loss on tokens of ground-truth responses $Y = \{a_1, ..., a_n\}$ is counted. We mask all tokens belonging to the instruction and observations from the environment to prevent them from loss computation. Concretely, the loss function is defined as:

$$\mathcal{L} = -\sum_j \log p_\theta(t_j | t_{<j}) \times \mathbf{1}(t_j \in Y), \quad (1)$$

where $t_j$ denotes the $j$-th input token and $\mathbf{1}$ is the indicator function.

Recent studies (Yang et al., 2024; Zeng et al., 2023) suggest that hybrid training with generalist instruction data and code data may improve the generalized ability of LLM agents. Following them, we adopt a mixture of SUPERAGENT $\mathcal{D}_{\text{agent}}$, the general domain instruction dataset $\mathcal{D}_{\text{general}}$, and the code dataset $\mathcal{D}_{\text{code}}$ for fine-tuning. We perform detailed ablation experiments to explore the effectiveness of generalist and code data in Section 5.2 and 5.3.

## 4 Experiments

### 4.1 Experimental Setup

**Base LLMs and Baselines** We select Llama-2-Chat (Touvron et al., 2023) as our base model for training SAMOYED. To ensure a fair evaluation, we include Llama-2-Chat and several Llama-2 based fine-tuned models as baselines, including Vicuna-1.5 (Chiang et al., 2023), CodeLlama-Instruct (Roziere et al., 2023), and AgentLM (Zeng et al., 2023). These models are fine-tuned with chatbot data, code, and a small number of agent trajectories, respectively. Due to our limited resources, we use 7B and 13B models for our experiments, leaving the comparison at a larger scale (*e.g.*, Llama-2-70B and Lemur-70B (Xu et al., 2023b)) for the future work. We also select GPT-3.5-Turbo (OpenAI, 2022) and GPT-4 (OpenAI, 2023) as strong baselines. For all LLMs, the decoding temperature is set to 0 for the most deterministic generation.

**Training Setup and Parameters** We use AdamW optimizer with a learning rate of 5e-5 and a cosine scheduler. The models are trained for 3 epochs with 3% warm-up steps. The batch size is set to 128 and the sequence length is 2048. We choose ShareGPT[4] as the generalist instruction data, and Evol-CodeAlpaca (Luo et al., 2023) as the code data. The mixture ratio of $\mathcal{D}_{\text{agent}}$, $\mathcal{D}_{\text{general}}$, and $\mathcal{D}_{\text{code}}$ is 80%, 10%, 10%. A corresponding data contamination analysis can be found in Appendix D. All experiments are conducted on 4 NVIDIA A100 80G GPUs. We use PyTorch FSDP (Paszke et al., 2019) for efficient training.

**Held-in/out Tasks** In an effort to balance the reliability and efficiency of the evaluation, we select nine tasks from SUPERAGENT to form the held-in test set. For tasks with a huge test set, following Wang et al. (2023), we randomly sample a subset from the original test set. To evaluate the generalized agent intelligence of SAMOYED, we additionally compile five unseen held-out tasks that do not exist in SUPERAGENT but still fall into the five skill dimensions of a foundation agent. The held-in and held-out evaluation tasks used in the experiments and the corresponding metrics are listed in Table 3. For all evaluated tasks, 1-shot in-context example is provided in prompts. We use average scores on held-in/out tasks to measure the overall capability of different agents. We also report the results on AgentBench (Liu et al., 2023), another agent benchmark, in Appendix F.

### 4.2 Main Results

Table 4 shows the results of different models on held-in and held-out tasks. For held-in tasks, it is expected that SAMOYED-7B outperforms GPT-3.5-

---

[4]https://sharegpt.com/

| Model | Held-in Tasks | | | | | | Held-out Tasks | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Reason | Math | Program | Web | Embodied | Avg. | Reason | Math | Program | Web | Embodied | Avg. |
| *Closed-Source Model* | | | | | | | | | | | | |
| GPT-4 | 61.6 | 73.0 | 54.9 | 40.6 | 77.8 | 59.8 | 41.6 | 92.0 | 69.4 | 69.4 | 36.4 | 61.8 |
| GPT-3.5-Turbo | 41.0 | 41.5 | 51.2 | 42.0 | 10.5 | 40.2 | 32.0 | 79.0 | 74.8 | 66.7 | 21.2 | 54.7 |
| *7B Open-Source Model* | | | | | | | | | | | | |
| Llama-2-7B-Chat | 4.0 | 7.5 | 2.5 | 13.9 | 0.0 | 6.2 | 4.0 | 12.0 | 7.0 | 0.4 | 7.8 | 6.3 |
| Vicuna-7B | 29.0 | 2.0 | _19.0_ | 24.2 | 6.0 | 17.1 | 8.8 | 21.9 | 19.0 | 18.2 | 12.8 | 16.1 |
| CodeLlama-7B | 3.5 | 3.5 | 1.5 | 24.8 | 0.0 | 7.4 | 1.0 | 18.0 | 21.8 | **41.3** | 5.5 | 17.5 |
| AgentLM-7B | _29.5_ | _10.0_ | 12.0 | **37.2** | **63.4** | _26.7_ | _19.2_ | _28.0_ | _50.5_ | 13.5 | _13.3_ | _25.0_ |
| SAMOYED-7B | **48.0** | **30.5** | **41.6** | _36.4_ | _61.2_ | **41.6** | **32.0** | **42.0** | **59.2** | _24.2_ | **14.2** | **34.3** |
| *13B Open-Source Model* | | | | | | | | | | | | |
| Llama-2-13B-Chat | 12.5 | 10.5 | 8.2 | 11.2 | 0.0 | 9.4 | 9.6 | 20.2 | 33.0 | 17.6 | 7.3 | 17.5 |
| Vicuna-13B | 25.5 | 6.5 | _30.4_ | 34.2 | 2.2 | 21.7 | _24.8_ | 37.0 | 37.0 | 34.2 | _14.8_ | 29.6 |
| CodeLlama-13B | 13.5 | _18.5_ | 5.1 | 15.3 | 0.0 | 11.7 | 6.4 | 36.0 | 11.1 | **46.5** | 5.5 | 21.1 |
| AgentLM-13B | _38.0_ | 13.5 | 22.8 | _38.1_ | _52.2_ | _30.8_ | 20.8 | _50.0_ | _46.6_ | 21.6 | 14.6 | _30.7_ |
| SAMOYED-13B | **54.5** | **38.5** | **55.4** | **40.9** | **72.4** | **50.1** | **35.0** | **51.0** | **62.4** | _38.9_ | **18.4** | **41.1** |

Table 4: Performance comparison of SAMOYED and baseline LLMs on held-in and held-out tasks. Due to the space constraint, we group the held-in tasks according to the skill dimensions and report the average scores. The best and the second of each model group are highlighted in **bold** and underlined respectively. See Appendix E for complete results.

Turbo due to the task coverage in our training data. However, the truly remarkable improvement lies in the performance of SAMOYED on held-out unseen tasks, which demonstrates a substantial boost in agent capabilities through large-scale trajectory tuning, compared to the untuned ones. Surprisingly, SAMOYED-7B exhibits an even greater enhancement compared to SAMOYED-13B. Our models also outperform AgentLM which is tuned on 1.8k trajectories, demonstrating the effectiveness of scaling up the tuning trajectories.

The experiment yields several noteworthy model-wise observations. Vicuna exhibits strong abilities through fine-tuning on generalist instruction data, demonstrating impressive performance on both held-in and held-out tasks. Remarkably, the performance of Vicuna-13B is even comparable to that of AgentLM-13B. It is important to highlight that AgentLM's training set comprises $80\%$ generalist instruction data, suggesting that the held-out task performance of AgentLM largely comes from the enhanced capability of instruction following. Additionally, we observe that CodeLlama excels in web browsing tasks. Nonetheless, it fails to achieve the best performance in programming tasks. This may be attributed to the fact that fine-tuning on code might compromise a model's general instruction-following capabilities, rendering it less effective in iterative programming environments.
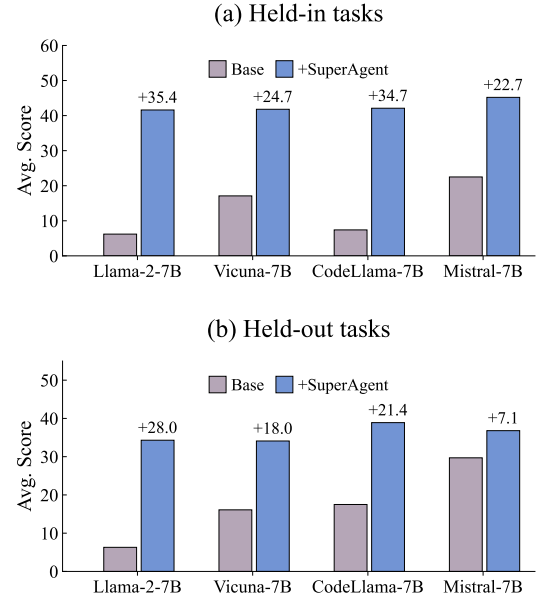


(a) Held-in tasks

(b) Held-out tasks

Figure 2: The results of different base models. "Base" denotes untrained LLMs. "+SuperAgent" denotes models after training on SUPERAGENT.

## 5 Further Analysis and Discussion

### 5.1 Different Base LLMs tuned on SUPERAGENT

To establish that the effectiveness of massive interaction trajectory training is model-unspecific, we fine-tune various base LLMs on SUPERAGENT. Two post-training versions of Llama-2 are
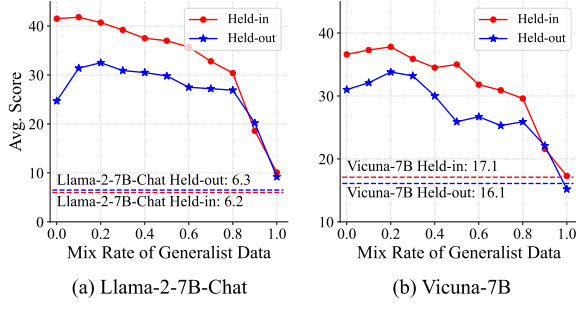
6

Figure 3: Ablation study on generalist instruction data.

| Model | Reason | Math | Program | Web | Embodied |
|---|---|---|---|---|---|
| Llama-2-7B-Chat | 4.0 | 12.0 | 7.0 | 0.4 | 7.8 |
| +SUPERAGENT | **32.0** | **42.0** | 59.2 | 24.2 | 14.2 |
| CodeLlama-7B | 1.0 | 18.0 | 21.8 | 41.3 | 5.5 |
| +SUPERAGENT | 29.6 | 40.0 | **67.7** | **42.2** | **14.8** |

Table 5: The held-out task performance of Llama-2 and CodeLlama.



Figure 4: Ablation study on code data.

selected: Vicuna-7B, which conducts instruction tuning using GPT-generated instruction data, and CodeLlama-7B which performs code training to enhance the code generation ability. Additionally, we include another popular open-source LLM with a different architecture, Mistral-7B (Jiang et al., 2023), which is engineered to excel Llama-2-13B in most benchmarks.

As illustrated in Figure 2, after large-scale trajectory tuning, all LLMs yield significant performance improvements on held-in and held-out tasks, demonstrating that LLMs can get generalized agent capability by learning from massive high-quality interaction trajectories. We also notice some interesting outcomes. Despite untuned Vicuna has a better agent performance than Llama-2, these two models exhibit similar performance after trajectory tuning. This suggests that agent ability of an LLM may be inherently acquired during the pretraining stage, while supervised fine-tuning merely surfaces these hidden capabilities. On the other hand, CodeLlama's superior performance indicates that code training can enhance agent capabilities. To delve deeper into the relationship between code data and agent intelligence, we conduct a comprehensive analysis in Section 5.3. As for Mistral, although fine-tuning on SUPERAGENT also yields improvements, the performance gain is relatively modest compared with the substantial improvement seen on Llama-2 (24% *v.s.* 442% on held-out tasks). This finding indicates that weaker LLMs may benefit more from massive trajectory tuning than their stronger counterparts.

### 5.2 The Effect of Generalist Instruction Data

When training SAMOYED, we mix 10% generalist instruction data to the training dataset. In an effort to distinguish *whether the improvement in held-out tasks is due to the agent's generalized capabilities or its enhanced ability to follow instructions*, we
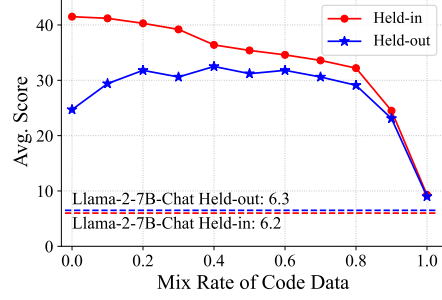
conduct an ablation study using ShareGPT dataset. We vary the mixture ratio of ShareGPT and train Llama-2-7B-Chat and Vicuna-7B for 1000 steps.

As shown in Figure 3, disagreed with Zeng et al. (2023) who find that training with only interaction trajectory data will lead to performance degradation on held-out tasks, SAMOYED trained on solely SUPERAGENT shows performance improvement on held-out tasks instead. This also manifests that generalized agent capabilities can be acquired by trajectory tuning alone. For held-in tasks, the performance degrades as the ratio of generalist data increases. This can be attributed to fewer learning trajectories during fixed training steps. On the other hand, for held-out tasks, a relatively low proportion of generalist data leads to improved performance. Nevertheless, as the amount of generalist data continues to increase, the performance on held-out tasks dramatically degrades. Based on this ablation study, we can conclude that the enhanced generalized agent capability is the primary factor behind the performance improvement on held-out tasks, and hybrid training with a small amount of generalist data will lead to further enhancement.

### 5.3 The Effect of Code Data

Code data, comprising standard syntax and logical abstraction, has the potential to enhance the planning and decision-making capabilities of LLM agents (Yang et al., 2024). In Table 5, we compare the distinctions between agents based on Llama-2-Chat and CodeLlama. Unsurprisingly, due to its ex-
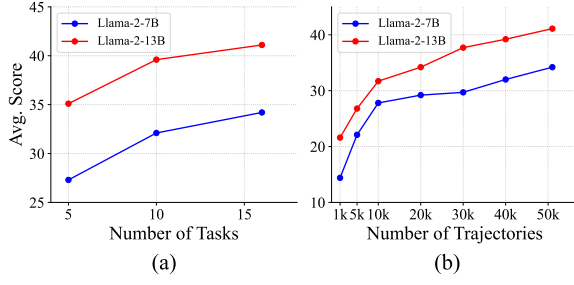
Figure 5: Scaling trends of the number of tasks and interaction trajectories.



Figure 6: Heatmap of skill-level capability transfer. We plot the relative improvements over training on generalist instruction and code data.

tensive code training, CodeLlama demonstrates excellent performance in programming tasks. Training with extensive interaction trajectories can further elevate its coding proficiency. Additionally, CodeLlama shows exceptional competence in web navigation tasks, likely attributed to the abundance of web pages present in its pretraining datasets.

We also carry out an ablation study by varying the mixture ratio of Evol-CodeAlpaca (Luo et al., 2023) and train Llama-2-7B-Chat for 1000 steps. As depicted in Figure 4, we notice an improvement in held-out task performance when a lower ratio of code data is incorporated. However, as the amount of code data continues to increase, both held-in and held-task performance begins to decline.

### 5.4 Scaling Trends of Generalization

We investigate the generalization performance of trajectory tuning with respect to two scaling factors: the number of training tasks and the number of training trajectories. Figure 5 illustrates the performance changes on held-out tasks when scaling each of these factors.

To explore the impact of task scaling, we modify the number of tasks in each skill dimension while ensuring that the skill coverage of the subsets remains consistent. We observe that increasing the number of tasks used for training results in improved performance on held-out tasks. This finding suggests that by scaling the number of distinct tasks for trajectory tuning, the model can enhance its generalized agent capabilities.

As shown in Figure 5b, comparing the performance with 1k trajectories to that with 50k+ trajectories, we observe a significant decline in the generalized agent ability with fewer trajectories. These results emphasize the importance of scaling the amount of interaction data for better performance.
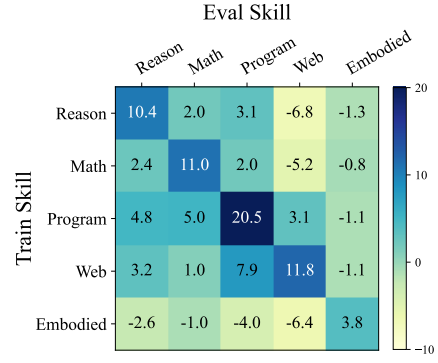
### 5.5 Skill-Level Transfer

To explore the potential transferability across different agent skills, we fine-tune Llama-2-Chat on held-in tasks corresponding to a specific agent skill and evaluate on held-out tasks. The compared baseline is fine-tuning Llama-2-Chat using a mixture of generalist instruction and code data. All models are trained for 300 steps to ensure a fair study.

As depicted in Figure 6, most skills, with the exception of embodied skill, exhibit the ability to transfer across different skill dimensions. This can be attributed to the unified agent interaction format in SUPERAGENT. The transferability of programming and web tasks further confirms the findings from Section 5.3. Notably, embodied AI skill is particularly challenging, for it receives negative impact from all other skills.

## 6 Conclusion

In this work, we explore the acquisition of generalized agent capabilities through fine-tuning open-source LLMs on massive interaction trajectories. We introduce by far the largest interaction trajectory dataset SUPERAGENT, comprising over 50k trajectories that encompass 16 tasks across five distinct agent skill dimensions. Building upon SUPERAGENT, we fine-tune Llama-2 to develop SAMOYED, an open-source LLM series specialized for agent tasks. Evaluations on both held-in and held-out tasks show that SAMOYED significantly outperforms strong baselines in terms of generalized agent capabilities. Comprehensive analysis also reveals the effectiveness of data mixture and plots the scaling law of trajectories. We hope this work to serve as a catalyst for further exploration in the development of more powerful agents.

## Limitations

We conclude the limitations of this work as follows:

- Due to the resource constraints, we only conduct experiments and analysis on 7B and 13B models. The extent to which larger models can benefit from large-scale trajectory tuning remains unknown.

- We have not fully explored the potential of equipping our SAMOYED with more sophisticated agent mechanisms, such as Reflexion (Shinn et al., 2023) and ReWOO (Xu et al., 2023a). Further investigation into these mechanisms could yield valuable insights.

- This work primarily focuses on improving the agent's performance via behavioral cloning (Pomerleau, 1991) from expert trajectories. How to further enhance the performance of the agent beyond behavioral cloning is underexplored. One potential approach is to allow the agent to self-play with the environment to facilitate improvement.

- This work is centered around building strong ReAct-style single-agent models. However, multi-agent collaboration framework has demonstrated impressive performance in handling realistic tasks. The development of strong generalized multi-agent systems based on open-source LLMs is still an under-explored area.

## Ethics Statement

This work fully complies with the ACL Ethics Policy. We declare that there are no ethical issues in this paper, to the best of our knowledge.

## References

Aida Amini, Saadia Gabriel, Peter Lin, Rik Koncel-Kedziorski, Yejin Choi, and Hannaneh Hajishirzi. 2019. Mathqa: Towards interpretable math word problem solving with operation-based formalisms. *arXiv preprint arXiv:1905.13319*.

Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. 2021. Program synthesis with large language models. *arXiv preprint arXiv:2108.07732*.

Rodney A Brooks. 1991. Intelligence without representation. *Artificial intelligence*, 47(1-3):139–159.

Baian Chen, Chang Shu, Ehsan Shareghi, Nigel Collier, Karthik Narasimhan, and Shunyu Yao. 2023. Fireact: Toward language agent fine-tuning. *arXiv preprint arXiv:2310.05915*.

Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. 2021. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*.

Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. 2023. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality.

Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. 2022. Scaling instruction-finetuned language models. *arXiv preprint arXiv:2210.11416*.

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.

Xiang Deng, Yu Gu, Boyuan Zheng, Shijie Chen, Samuel Stevens, Boshi Wang, Huan Sun, and Yu Su. 2023. Mind2web: Towards a generalist agent for the web. *arXiv preprint arXiv:2306.06070*.

Mor Geva, Daniel Khashabi, Elad Segal, Tushar Khot, Dan Roth, and Jonathan Berant. 2021. Did aristotle use a laptop? a question answering benchmark with implicit reasoning strategies. *Transactions of the Association for Computational Linguistics*, 9:346–361.

Daniel Gordon, Aniruddha Kembhavi, Mohammad Rastegari, Joseph Redmon, Dieter Fox, and Ali Farhadi. 2018. Iqa: Visual question answering in interactive environments. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4089–4098.

Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021. Measuring mathematical problem solving with the math dataset. *arXiv preprint arXiv:2103.03874*.

Ahmed Hussein, Mohamed Medhat Gaber, Eyad Elyan, and Chrisina Jayne. 2017. Imitation learning: A survey of learning methods. *ACM Computing Surveys (CSUR)*, 50(2):1–35.

Nicholas R Jennings, Katia Sycara, and Michael Wooldridge. 1998. A roadmap of agent research and development. *Autonomous agents and multi-agent systems*, 1:7–38.

Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7b. *arXiv preprint arXiv:2310.06825*.

Mandar Joshi, Eunsol Choi, Daniel S Weld, and Luke Zettlemoyer. 2017. Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension. *arXiv preprint arXiv:1705.03551*.

Geunwoo Kim, Pierre Baldi, and Stephen McAleer. 2023. Language models can solve computer tasks. *arXiv preprint arXiv:2303.17491*.

Ariel N Lee, Cole J Hunter, and Nataniel Ruiz. 2023. Platypus: Quick, cheap, and powerful refinement of llms. *arXiv preprint arXiv:2308.07317*.

Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, et al. 2022. Holistic evaluation of language models. *arXiv preprint arXiv:2211.09110*.

Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei, Hanyu Lai, Yu Gu, Hangliang Ding, Kaiwen Men, Kejuan Yang, et al. 2023. Agentbench: Evaluating llms as agents. *arXiv preprint arXiv:2308.03688*.

Ziyang Luo, Can Xu, Pu Zhao, Qingfeng Sun, Xiubo Geng, Wenxiang Hu, Chongyang Tao, Jing Ma, Qingwei Lin, and Daxin Jiang. 2023. Wizardcoder: Empowering code large language models with evol-instruct.

Yohei Nakajima. 2023. Babyagi. *Python.* *https://github.com/yoheinakajima/babyagi*.

OpenAI. 2022. Introducing chatgpt.

OpenAI. 2023. Gpt-4 technical report. *arXiv*, pages 2303–08774.

Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. 2019. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32.

Arkil Patel, Satwik Bhattamishra, and Navin Goyal. 2021. Are nlp models really able to solve simple math word problems? *arXiv preprint arXiv:2103.07191*.

Dean A Pomerleau. 1991. Efficient training of artificial neural networks for autonomous navigation. *Neural computation*, 3(1):88–97.

Ofir Press, Muru Zhang, Sewon Min, Ludwig Schmidt, Noah A Smith, and Mike Lewis. 2022. Measuring and narrowing the compositionality gap in language models. *arXiv preprint arXiv:2210.03350*.

Toran Bruce Richards. 2023. Auto-gpt: An autonomous gpt-4 experiment.

Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, et al. 2023. Code llama: Open foundation models for code. *arXiv preprint arXiv:2308.12950*.

Victor Sanh, Albert Webson, Colin Raffel, Stephen H Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Teven Le Scao, Arun Raja, et al. 2021. Multitask prompted training enables zero-shot task generalization. *arXiv preprint arXiv:2110.08207*.

Noah Shinn, Beck Labash, and Ashwin Gopinath. 2023. Reflexion: an autonomous agent with dynamic memory and self-reflection. *arXiv preprint arXiv:2303.11366*.

Mohit Shridhar, Jesse Thomason, Daniel Gordon, Yonatan Bisk, Winson Han, Roozbeh Mottaghi, Luke Zettlemoyer, and Dieter Fox. 2020a. Alfred: A benchmark for interpreting grounded instructions for everyday tasks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10740–10749.

Mohit Shridhar, Xingdi Yuan, Marc-Alexandre Côté, Yonatan Bisk, Adam Trischler, and Matthew Hausknecht. 2020b. Alfworld: Aligning text and embodied environments for interactive learning. *arXiv preprint arXiv:2010.03768*.

Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Stanford alpaca: An instruction-following llama model. `https://github.com/tatsu-lab/stanford_alpaca`.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.

Ruoyao Wang, Peter Jansen, Marc-Alexandre Côté, and Prithviraj Ammanabrolu. 2022a. Scienceworld: Is your agent smarter than a 5th grader? *arXiv preprint arXiv:2203.07540*.

Xingyao Wang, Zihan Wang, Jiateng Liu, Yangyi Chen, Lifan Yuan, Hao Peng, and Heng Ji. 2023. Mint: Evaluating llms in multi-turn interaction with tools and language feedback. *arXiv preprint arXiv:2309.10691*.

Yizhong Wang, Swaroop Mishra, Pegah Alipoormolabashi, Yeganeh Kordi, Amirreza Mirzaei, Anjana Arunkumar, Arjun Ashok, Arut Selvan Dhanasekaran, Atharva Naik, David Stap, et al. 2022b. Super-naturalinstructions: Generalization via declarative instructions on 1600+ nlp tasks. *arXiv preprint arXiv:2204.07705*.

Jason Wei, Maarten Bosma, Vincent Y Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M Dai, and Quoc V Le. 2021. Finetuned language models are zero-shot learners. *arXiv preprint arXiv:2109.01652*.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*, 35:24824–24837.

Luca Weihs, Matt Deitke, Aniruddha Kembhavi, and Roozbeh Mottaghi. 2021. Visual room rearrangement. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.

Binfeng Xu, Zhiyuan Peng, Bowen Lei, Subhabrata Mukherjee, Yuchen Liu, and Dongkuan Xu. 2023a. Rewoo: Decoupling reasoning from observations for efficient augmented language models. *arXiv preprint arXiv:2305.18323*.

Yiheng Xu, Hongjin Su, Chen Xing, Boyu Mi, Qian Liu, Weijia Shi, Binyuan Hui, Fan Zhou, Yitao Liu, Tianbao Xie, et al. 2023b. Lemur: Harmonizing natural language and code for language agents. *arXiv preprint arXiv:2310.06830*.

John Yang, Akshara Prabhakar, Karthik Narasimhan, and Shunyu Yao. 2023. Intercode: Standardizing and benchmarking interactive coding with execution feedback. *arXiv preprint arXiv:2306.14898*.

Ke Yang, Jiateng Liu, John Wu, Chaoqi Yang, Yi R Fung, Sha Li, Zixuan Huang, Xu Cao, Xingyao Wang, Yiquan Wang, et al. 2024. If llm is the wizard, then code is the wand: A survey on how code empowers large language models to serve as intelligent agents. *arXiv preprint arXiv:2401.00812*.

Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W Cohen, Ruslan Salakhutdinov, and Christopher D Manning. 2018. Hotpotqa: A dataset for diverse, explainable multi-hop question answering. *arXiv preprint arXiv:1809.09600*.

Shunyu Yao, Howard Chen, John Yang, and Karthik Narasimhan. 2022a. Webshop: Towards scalable real-world web interaction with grounded language agents. *Advances in Neural Information Processing Systems*, 35:20744–20757.

Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. 2022b. React: Synergizing reasoning and acting in language models. *arXiv preprint arXiv:2210.03629*.

Da Yin, Faeze Brahman, Abhilasha Ravichander, Khyathi Chandu, Kai-Wei Chang, Yejin Choi, and Bill Yuchen Lin. 2023. Lumos: Learning agents with unified data, modular design, and open-source llms. *arXiv preprint arXiv:2311.05657*.

Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, Yu Zhang, James T Kwok, Zhenguo Li, Adrian Weller, and Weiyang Liu. 2023. Metamath: Bootstrap your own mathematical questions for large language models. *arXiv preprint arXiv:2309.12284*.

Xiang Yue, Xingwei Qu, Ge Zhang, Yao Fu, Wenhao Huang, Huan Sun, Yu Su, and Wenhu Chen. 2023. Mammoth: Building math generalist models through hybrid instruction tuning. *arXiv preprint arXiv:2309.05653*.

Aohan Zeng, Mingdao Liu, Rui Lu, Bowen Wang, Xiao Liu, Yuxiao Dong, and Jie Tang. 2023. Agenttuning: Enabling generalized agent abilities for llms. *arXiv preprint arXiv:2310.12823*.

Shengyu Zhang, Linfeng Dong, Xiaoya Li, Sen Zhang, Xiaofei Sun, Shuhe Wang, Jiwei Li, Runyi Hu, Tianwei Zhang, Fei Wu, et al. 2023. Instruction tuning for large language models: A survey. *arXiv preprint arXiv:2308.10792*.

Zilong Zheng, Mengmeng Wang, Zixia Jia, and Baichen Tong. 2023. Langsuite: Controlling, planning, and interacting with large language models in embodied text environments.

Shuyan Zhou, Frank F Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng, Yonatan Bisk, Daniel Fried, Uri Alon, et al. 2023. Webarena: A realistic web environment for building autonomous agents. *arXiv preprint arXiv:2307.13854*.

11

# A Details of Tasks in SUPERAGENT

**Reasoning Tasks** HotpotQA (Yang et al., 2018) is a question answering dataset featuring multi-hop reasoning. StrategyQA (Geva et al., 2021) is another question answering task where the required reasoning steps are implicit in the question and should be inferred using a strategy. TriviaQA (Joshi et al., 2017) is a dataset consisting of complex compositional questions that require multi-evidence reasoning. In our work, we repurpose these three datasets to interaction environments by incorporating a search engine tool. We employ the GPT-exploration pipeline and filter out failed cases to build the golden trajectories.

For our held-out evaluation, we use Bamboogle (Press et al., 2022), which is made up of questions that need compositional reasoning and are unable to be directly answered by Google.

**Math Tasks** GSM8K (Cobbe et al., 2021) is a dataset of diverse grade school math problems created by humans. Each problem in GSM8K comes with an official solution path. In our work, we leverage the power of GPT-3.5-Turbo to transform these solution paths into interaction trajectories.

MathQA (Amini et al., 2019) is a large-scale multiple-choice math problem dataset covering multiple math domains. MATH (Press et al., 2022) contains challenging mathematics problems from high school math competitions. To adapt these two datasets into interaction environments, we employ a Python interpreter and employ the GPT-exploration pipeline to construct the trajectories.

For the held-out task, we use SVAMP (Patel et al., 2021), which is a challenge set for elementary-level Math Word Problems.

**Programming Tasks** InterCode (Yang et al., 2023) is a benchmark for evaluating language models on interactive programming tasks. In this task, agents are required to respond to natural language requests by interacting with a software system, such as a database or terminal. Our work focuses on evaluating the programming ability of agents using two environments: IC-Bash and IC-SQL. IC-Bash is specifically used for the held-out evaluation of agents.

APPS (Hendrycks et al., 2021) is a benchmark focused on Python code generation, encompassing a range of difficulty levels from introductory to competition level. We utilize GPT-3.5-Turbo to reformat the instances in this dataset and construct the trajectories.

HumanEval (Chen et al., 2021) is a dataset designed to measure functional correctness for synthesizing programs from docstrings. MBPP (Austin et al., 2021) consists of around 1,000 crowd-sourced Python programming problems. For both of these datasets, we employ the GPT-exploration pipeline to annotate the interaction trajectories. Subsequently, we employ the answer forcing method to re-annotate the cases where GPT failed.

**Web Tasks** Mind2Web (Deng et al., 2023) is a dataset for developing and evaluating generalist agents for the web that can follow language instructions to complete complex tasks on any website. WebArena (Zhou et al., 2023) builds realistic web environments for agents to execute tasks. Even GPT-4 struggles with these tasks, so we utilize a teacher forcing and break down the complete interaction trajectory into multiple single steps. Then GPT-3.5-Turbo is employed to annotate the rationales.

WebShop (Yao et al., 2022a) is a simulated e-commerce website environment with real-world products and crowd-sourced text instructions. For 1571 official human annotated trajectories, we employ GPT-3.5-Turbo to reformat them and annotate rationales. Additionally, we incorporate trajectories generated through GPT-exploration, which have final rewards exceeding 0.3.

For our held-out task, we utilize Mini-WoB++ (Kim et al., 2023), a diverse collection of over 100 web interaction environments, to formulate our benchmark.

**Embodied AI Tasks** ALFWorld (Shridhar et al., 2020b) contains interactive TextWorld environments that parallel embodied worlds in the ALFRED dataset (Shridhar et al., 2020a). This dataset provides human-annotated golden trajectories for imitation learning. RoomR (Weihs et al., 2021) is an embodied AI dataset which requires agents to restore the initial configurations of all objects within a room. IQA (Gordon et al., 2018) is a question answering task that requires an agent to interact with a dynamic visual environment. In our work, we utilize the text versions of RoomR and IQA developed by Zheng et al. (2023). We employ a depth-first-search algorithm to build the golden action sequences for RoomR and IQA. We then leverage GPT-3.5-Turbo to annotate the corresponding rationales.

| Dataset | Model | $R_{\text{train}}$ | $R_{\text{pseudo}}$ | $R_{\text{test}}$ | $\Delta_1$ | $\Delta_2$ |
|---------|-------|--------|---------|--------|------------|------------|
| AgentInstruct (Zeng et al., 2023) | Llama-2-7B-Chat | 17.8 | 17.5 | 15.8 | -0.3 | -2.0 |
| | $+\mathcal{D}_{\text{train}}$ | 72.5 | 72.6 | 62.4 | +0.1 | -10.1 |
| SUPERAGENT (Ours) | Llama-2-7B-Chat | 16.2 | 16.5 | 16.0 | +0.3 | -0.2 |
| | $+\mathcal{D}_{\text{train}}$ | 73.3 | 62.3 | 62.8 | -11.0 | -10.5 |

Table 6: The average reward of WebShop on different instruction sets. We compare the reward $R_{\text{train}}$, $R_{\text{dev}}$, $R_{\text{test}}$ on the training set $\mathcal{D}_{\text{train}}$, a pseudo test set held-out from the original training set $\mathcal{D}_{\text{pseudo}}$, and original test set $\mathcal{D}_{\text{test}}$ respectively. We also reports two key metrics: $\Delta_1 = R_{\text{pseudo}} - R_{\text{train}}$ and $\Delta_2 = R_{\text{test}} - R_{\text{train}}$, as the indicators of the difficulty differences between datasets.

For the held-out evaluation, we utilize Science-World (Wang et al., 2022a), a text-based virtual environment which encompasses various elementary science experiment tasks, including thermodynamics and electrical circuits.

## B Difficulty Bias in Trajectory Collection

In this section, we conduct a experiment to verify the existence of difficulty bias introduced by the trajectory annotation pipeline widely used in recent studies (Chen et al., 2023; Zeng et al., 2023). Specifically, we choose WebShop trajectories in SUPERAGENT and AgentInstruct (Zeng et al., 2023) to conduct the experiment. For AgentInstruct and SUPERAGENT, we select 300 instances as the training set $\mathcal{D}_{\text{train}}$, 50 instances as the pseudo test set $\mathcal{D}_{\text{pseudo}}$. We also include the original WebShop test set $\mathcal{D}_{\text{test}}$.

For a dataset conforming to the *i.i.d.* assumption, the instances in $\mathcal{D}_{\text{train}}$, $\mathcal{D}_{\text{pseudo}}$, $\mathcal{D}_{\text{test}}$ are sampled from the same distribution. Therefore, the expected behavior is that the evaluation results on $\mathcal{D}_{\text{pseudo}}$ and $\mathcal{D}_{\text{test}}$ should be consistent. Furthermore, an agent trained on $\mathcal{D}_{\text{train}}$ should ideally perform better on $\mathcal{D}_{\text{train}}$ compared to $\mathcal{D}_{\text{pseudo}}$ and $\mathcal{D}_{\text{test}}$.

Table 6 illustrates the performance of untrained Llama-2-7B-Chat and the trained agent on different sets. For AgentInstruct, both models exhibit worse performance on $\mathcal{D}_{\text{test}}$ compared to $\mathcal{D}_{\text{pseudo}}$, indicating that instances in AgentInstruct are considerably easier than those in the original test set. Conversely, for SUPERAGENT, the agents have close performance on $\mathcal{D}_{\text{pseudo}}$ and $\mathcal{D}_{\text{test}}$, aligning with our expectations. The agent trained on our dataset also outperforms the agent trained on AgentInstruct when evaluated on $\mathcal{D}_{\text{test}}$. These experiments highlight that the GPT-exploration trajectory annotation pipeline can introduce difficulty bias in the training set, potentially compromising the generalizability of trained agents.

| Dataset | Win | Lose | Tie | Total |
|---------|-----|------|-----|-------|
| IC-SQL | 11 | 16 | 73 | 100 |
| WebShop | 12 | 10 | 58 | 80 |

Table 7: Human evaluation of the data quality for SUPERAGENT. For IC-SQL, we compare trajectories generated through answer forcing with those generated through exploration. For WebShop, we compare our constructed trajectories with the trajectories constructed by Zeng et al. (2023).

## C Quality Control of SUPERAGENT

In Section 3.4, we incorporate heuristic and GPT-based methods to construct SUPERAGENT, which can mitigate the difficulty bias problem in the previous annotation pipeline. In this section, we propose to perform a human evaluation to assess the quality of SUPERAGENT. To achieve this, we employ 5 human annotators who are instructed to choose the better trajectory from two anonymous candidate options. Here, we select two representative tasks: IC-SQL to assess the quality of answer forcing annotation, and WebShop to evaluate the quality of trajectory reformatting. For IC-SQL, we compare 100 trajectories generated by answer forcing with those generated through GPT exploration. For WebShop, we select 80 trajectories from SUPERAGENT and Zeng et al. (2023) which correspond to the same task instance.

As shown in Table 7, for most cases, trajectories generated by answer forcing or reformatting have the same quality as GPT exploration. Therefore, we can conclude that our trajectory annotation process can achieve comparable quality with previous methods (Chen et al., 2023; Zeng et al., 2023) while mitigating the difficulty bias.

## D Data Contamination

When training SAMOYED, we construct a data mixture consisting of trajectory data (SUPERAGENT), generalist instruction data (ShareGPT), and code data (Evol-CodeAlpaca). However, it is important to address the concern of potential data contamination, which could result in an overestimation of performance. Therefore, we perform a contamination analysis by comparing our evaluation set with SUPERAGENT, ShareGPT, and Evol-CodeAlpaca. Following Liang et al. (2022), we heuristically match 9-grams and 13-grams from the instances in the test set with the training set data. Table 8 displays the proportion of instances which exhibit an overlap with the training data.

First, we observe a high contamination rate for held-in tasks with SUPERAGENT. After manually examining these instances, we have some findings. In the case of StrategyQA, we discovered that all instances followed a question format that could be answered with a simple "yes" or "no," potentially resulting in a high n-gram overlap. For WebShop and ALFWorld, we found that the contamination may be attributed to the template-based data construction process. For instance, in WebShop, instructions consistently followed specific formats like "I would like <product> that is <size> and is the color <color>, and price lower than <price> dollars". Additionally, we observed that MBPP suffers from data contamination issues across all three training sets. After manual inspection, we determined that most of the overlap occurs in importing Python packages and commonly used code snippets, such as loops.

In summary, it can be concluded that the data contamination has a minimal impact on the experimental results. While some overlap exists between the held-in tasks and the training set, this is primarily a result of their data construction process. Moreover, by adhering to the original train-test split of the datasets, the extent of performance overestimation is reduced. Most importantly, the held-out tasks, which are used to assess the agents' generalized capabilities, do not suffer from the issue of data contamination. This ensures the trustworthiness and robustness of our evaluation.

## E Complete Experimental Results

Table 9 shows the complete results on held-in tasks.

## F Evaluation on AgentBench

AgentBench (Liu et al., 2023) is another evaluation benchmark for LLM agents, encompassing 8 agent tasks. However, it is worth noting that some tasks in AgentBench are already covered by SUPERAGENT, and some tasks may pose a risk of data contamination with our dataset. Nevertheless, to provide a comprehensive perspective, we have included the results of SAMOYED on AgentBench as a point of reference in Table 10.

14

| Dataset | #Inst | SUPERAGENT | | ShareGPT | | Evol-CodeAlpaca | |
|---|---|---|---|---|---|---|---|
| | | 9-Gram Rate | 13-Gram Rate | 9-Gram Rate | 13-Gram Rate | 9-Gram Rate | 13-Gram Rate |
| *Held-in Tasks* | | | | | | | |
| HotpotQA | 100 | 1% | 0% | 0% | 0% | 0% | 0% |
| StrategyQA | 100 | 20% | 12% | 0% | 0% | 0% | 0% |
| GSM8K | 100 | 3% | 0% | 0% | 0% | 0% | 0% |
| MATH | 100 | 15% | 4% | 0% | 0% | 2% | 0% |
| IC-SQL | 100 | 7% | 0% | 0% | 0% | 1% | 0% |
| MBPP | 100 | 12% | 1% | 7% | 3% | 18% | 4% |
| Mind2Web | 1173 | 8% | 3% | 0% | 0% | 0% | 0% |
| WebShop | 200 | 41% | 14% | 0% | 0% | 0% | 0% |
| ALFWorld | 134 | 14% | 8% | 0% | 0% | 0% | 0% |
| *Held-out Tasks* | | | | | | | |
| Bamboogle | 126 | 0% | 0% | 0% | 0% | 0% | 0% |
| SVAMP | 100 | 0% | 0% | 0% | 0% | 0% | 0% |
| IC-Bash | 200 | 0% | 0% | 0% | 0% | 0% | 0% |
| MiniWoB++ | 460 | 0% | 0% | 0% | 0% | 2% | 0% |
| SciWorld | 270 | 0% | 0% | 0% | 0% | 0% | 0% |

Table 8: Data contamination analysis.

| Model | Held-in Tasks | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | HotpotQA | StrategyQA | GSM8K | MATH | IC-SQL | MBPP | Mind2Web | WebShop | ALFWorld | Avg. |
| *Closed-Source Model* | | | | | | | | | | |
| GPT-4 | 52.1 | 71.0 | 87.0 | 59.0 | 37.8 | 72.0 | 22.6 | 58.6 | 77.8 | 59.8 |
| GPT-3.5-Turbo | 24.0 | 58.0 | 65.0 | 18.0 | 38.5 | 64.0 | 21.7 | 62.4 | 10.5 | 40.2 |
| *7B Open-Source Model* | | | | | | | | | | |
| Llama-2-7B-Chat | 3.0 | 5.0 | 15.0 | 0.0 | 4.0 | 1.0 | 11.9 | 15.8 | 0.0 | 6.2 |
| Vicuna-7B | 11.0 | 47.0 | 1.0 | 3.0 | 17.3 | 21.0 | 14.8 | 33.5 | 6.0 | 17.2 |
| CodeLlama-7B | 2.0 | 5.0 | 7.0 | 0.0 | 3.0 | 0.0 | 17.0 | 32.5 | 0.0 | 7.4 |
| AgentLM-7B | 10.0 | 49.0 | 14.0 | 6.0 | 13.9 | 10.0 | 10.6 | 63.7 | 63.4 | 26.7 |
| SAMOYED-7B | 30.0 | 66.0 | 43.0 | 18.0 | 59.2 | 24.0 | 12.2 | 60.5 | 61.2 | 41.6 |
| *13B Open-Source Model* | | | | | | | | | | |
| Llama-2-13B-Chat | 6.0 | 19.0 | 18.0 | 3.0 | 3.0 | 13.4 | 17.2 | 5.3 | 0.0 | 9.4 |
| Vicuna-13B | 15.0 | 36.0 | 9.0 | 4.0 | 37.0 | 23.7 | 15.2 | 53.3 | 2.2 | 21.7 |
| CodeLlama-13B | 7.0 | 20.0 | 29.0 | 8.1 | 3.0 | 7.2 | 7.6 | 23.0 | 0.0 | 11.7 |
| AgentLM-13B | 24.0 | 52.0 | 21.0 | 6.1 | 25.7 | 20.0 | 11.1 | 65.0 | 52.2 | 30.8 |
| SAMOYED-13B | 41.0 | 68.0 | 53.0 | 24.0 | 67.7 | 43.0 | 18.6 | 63.1 | 72.4 | 50.1 |

Table 9: Performance of SAMOYED and baseline LLMs on held-in tasks.

| Model | Code-grounded | | | Game-grounded | | | Web-grounded | | Overall |
|---|---|---|---|---|---|---|---|---|---|
| | OS[†] | DB[†] | KG[†] | DCG | LTP | HH[‡] | WS[‡] | WB[‡] | |
| GPT-4 | 42.4 | 32.0 | 58.8 | 74.5 | 16.6 | 78.0 | 61.1 | 29.0 | 4.01 |
| GPT-3.5-Turbo | 32.6 | 36.7 | 25.9 | 33.7 | 10.5 | 16.0 | 64.1 | 20.0 | 2.32 |
| Llama-2-7B-Chat | 4.2 | 8.0 | 2.1 | 6.9 | 0.0 | 0.0 | 11.6 | 7.0 | 0.34 |
| Vicuna-7B | 9.7 | 8.7 | 2.5 | 0.3 | 6.4 | 0.0 | 2.2 | 9.0 | 0.56 |
| CodeLlama-7B | 4.9 | 12.7 | 8.2 | 0.0 | 0.0 | 2.0 | 25.2 | 12.0 | 0.50 |
| SAMOYED-7B | 11.8 | 9.7 | 2.7 | 1.9 | 8.2 | 68.0 | 60.5 | 12.2 | 1.60 |

Table 10: Performance of SAMOYED and baseline LLMs on AgentBench (Liu et al., 2023). † means the test set may suffer data contamination with SUPERAGENT. ‡ means the task is already covered by SUPERAGENT.