

AGENT-AS-A-JUDGE: EVALUATE AGENTS WITH AGENTS

Anonymous authors

Paper under double-blind review

ABSTRACT

Contemporary evaluation techniques are inadequate for agentic systems. These approaches either focus exclusively on final outcomes—ignoring the step-by-step nature of the thinking done by agentic systems—or require excessive manual labour. To address this, we introduce the **Agent-as-a-Judge** framework, wherein agentic systems are used to evaluate agentic systems. This is a natural extension of the LLM-as-a-Judge framework, incorporating agentic features that enable intermediate feedback for the entire task-solving processes for more precise evaluations. We apply the Agent-as-a-Judge framework to the task of code generation. To overcome issues with existing benchmarks and provide a proof-of-concept testbed for Agent-as-a-Judge, we present **DevAI**, a new benchmark of 55 realistic AI code generation tasks. DevAI includes rich manual annotations, like a total of 365 hierarchical solution requirements, which make it particularly suitable for an agentic evaluator. We benchmark three of the top code-generating agentic systems using Agent-as-a-Judge and find that our framework dramatically outperforms LLM-as-a-Judge and is as reliable as our human evaluation baseline. Altogether, we believe that this work represents a concrete step towards enabling vastly more sophisticated agentic systems. To help that, our dataset and the full implementation of Agent-as-a-Judge will be publically available at [REDACTED]

1 INTRODUCTION

Recent years have seen multimodal agentic systems move from occasionally being able to solve small toy problems to being regularly deployed for challenging real-world problems (the dream of most AI research). Yet, the current evaluation methods and the available benchmarks for agentic systems are struggling to keep up with these rapid advances, dramatically slowing true progress.

We believe that the current issue with evaluating agentic systems stems from the lack of feedback during the intermediate task-solving stages for these nontraditional systems. Agentic systems think more like a human, often act step-by-step (Wooldridge, 1999) and often hosting very human-like natural language discussions internally to solve problems (Zhuge et al., 2023). And thus agentic systems should be evaluated like a human, with rich evaluative feedback which looks at the full thought and action trajectory; evaluating an agentic system in the traditional way is like evaluating a student using multiple-choice testing—a comparatively unreliable estimator (Park, 2010). For example, while SWE-Bench (Jimenez et al.) is widespread, its evaluation method, which relies solely on the final resolve rate for long-term automated repair tasks, does not effectively pinpoint what is happening within agentic systems that affects the resolve rate. On the other hand, performing a better evaluation with a human is prohibitively expensive. We instead propose that agentic systems should be used to evaluate agentic systems. Inspired by LLM-as-a-Judge (Zheng et al., 2024; Fu et al., 2024; Chen et al.), which uses LLMs to evaluate LLMs, we call this framework Agent-as-a-Judge, of which it is a key extension to the world of agentic systems (see Figure 1). It not only retains the cost-effectiveness of LLM-as-a-Judge but is also equipped with agentic features, allowing it to provide rich intermediate feedback throughout the entire process, as it acts as an agentic system. We apply the Agent-as-a-Judge systems to the problem of evaluating code generating systems—one of the areas where agentic systems have looked the most promising recently.

In code generation, the development of benchmarks has also lagged behind the rapid advancement of agentic systems. HumanEval (Chen et al., 2021), for example, focuses exclusively on algorithmic problems, while MBPP (Austin et al., 2021) deals with simple programming tasks. Although they

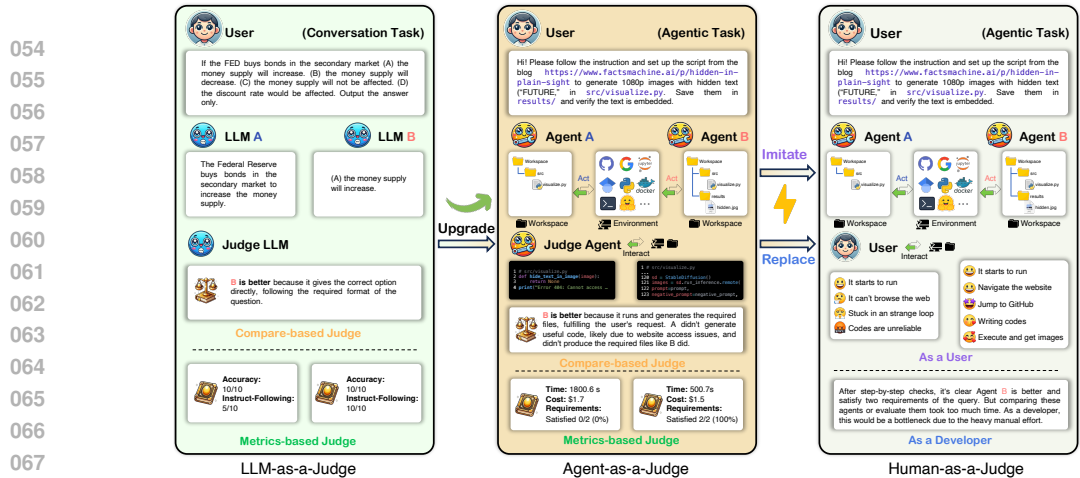


Figure 1: In this paper, we introduce the Agent-as-a-Judge framework wherein agentic systems are used to evaluate agentic systems. We compare this to LLM-as-a-Judge, which uses LLMs to evaluate LLMs and for which Agent-as-a-Judge is a natural evolution, and Human-as-a-Judge, where skilled humans labourers manually evaluate an agentic system.

are useful for evaluating the basic skills of foundation models, neither of these two reflects the most practical challenges developers face. As a step away from this, SWE-Bench (Jimenez et al.) did introduce more realistic problems from GitHub, offering a fresh approach to evaluation, but still primarily focuses on automated repairs tasks development process. Concerningly, recent research shows that large language models (LLMs) can already solve over 27% of the tasks in SWE-Bench without needing of advanced agentic systems (Xia et al., 2024). Equally concerning, recent work has begun to introduce mechanisms designed specifically for the individual tasks in the SWE-Bench dataset, leading to a lack of real-world generalization and violating Goodhart’s law: “When a measure becomes a target, it ceases to be a good measure” (Goodhart, 1976).

To address the aforementioned issues with the current benchmarks in code generation, we introduce DevAI: the AI Developer Dataset, which contains 55 real-world comprehensive AI app development tasks created by expert annotators. We apply three leading open-source code-generating agentic frameworks to the tasks in DevAI: MetaGPT (Hong et al., 2024b), GPT-Pilot (Pythagora.io, 2023), and OpenHands (Wang et al., 2024b). We evaluate their performance using human judges (a painstaking process), LLM-as-a-Judge (Zheng et al., 2024), and our Agent-as-a-Judge framework.

Through human evaluation, we found that GPT-Pilot and OpenHands were each able to satisfy about 29% of the task requirements in DevAI, but only one full task—showing that DevAI presents a good level of challenge to current systems. When comparing our human judges with our automatic Agent-as-a-Judge framework, we found that Agent-as-a-Judge aligns more closely with the consensus of our human judges (90%) as compared to LLM-as-a-Judge (70%) in all cases tested. In addition, we find that it aligns more closely with this ensemble than the individual human evaluators do, suggesting that—not only is it suitable as a replacement for a human evaluator—but it could in fact be more useful than an average lone human evaluator. In addition, considering the evaluation cost, Agent-as-a-Judge reduces 97.72% of the time and 97.64% of the cost.

In summary, the principal contributions of this work are:

- We release the DevAI dataset, which consists of 55 comprehensive AI development tasks with accompanying tags, individual hierarchical requirements, and individual preferences.
- We benchmark three top open-source code generation agentic frameworks in DevAI, providing a more comprehensive analysis than previous evaluations of them.
- We introduce the general Agent-as-a-Judge concept, allowing agentic systems a fair and rich evaluation without the traditional cost this would require in human labour.
- We demonstrate that an Agent-as-a-Judge outperforms an LLM-as-a-Judge and performs comparably to human evaluators in our proof-of-concept.

This paper is structured as follows: Section 2 introduces DevAI to address the lack of benchmarks for verifying agentic systems with intermediate processes. Section 3 establishes Human-as-a-Judge

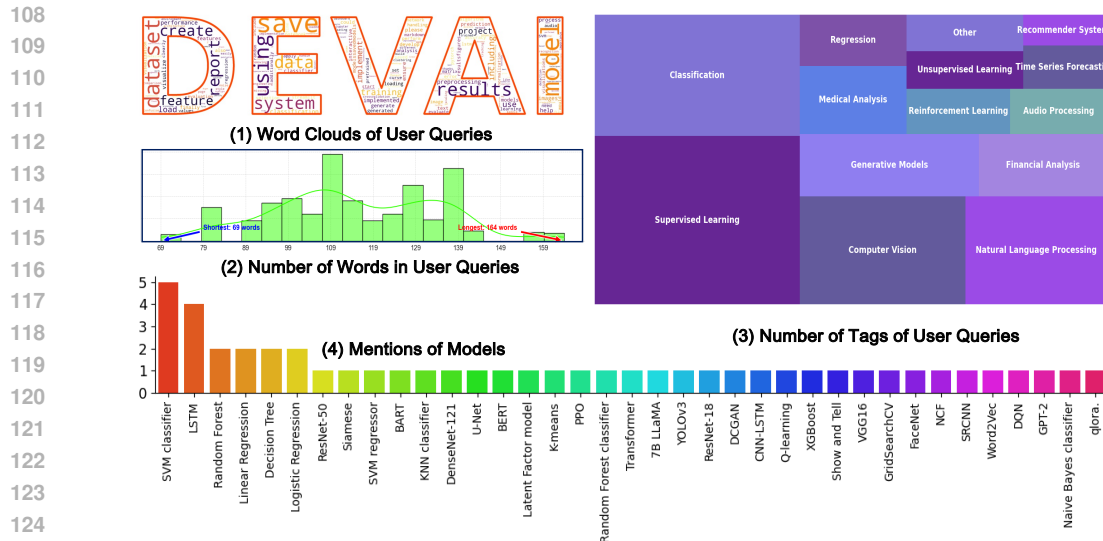


Figure 2: **Distribution of DevAI Tasks** (1) DevAI focuses on AI development tasks and so terms such as “dataset,” “model,” and “results” are particularly common in the queries. (2) The first 53 tasks in DevAI all have a one-paragraph query but of varying lengths (note that task 54 and 55 are excluded here as they are outliers, representing the longest and most complex tasks in the dataset). (3) Each task has one or more tags. The prevalence of supervised learning here reflects the fact that it dominates many machine learning applications. (4) SVM classifiers (Cortes, 1995) and LSTM models (Hochreiter, 1997) are two of the most widely used architectures—a fact reflected by DevAI.

as a manual evaluation baseline, highlighting its limitations. Finally, Section 4 presents Agent-as-a-Judge, a scalable solution to these challenges. More details are provided in Appendices A and B.

2 DEVAI: A DATASET FOR AUTOMATED AI DEVELOPMENT

2.1 MOTIVATION

Background The code generation domain is an area where agentic systems have seen significant industrial deployment during the past 2 years (e.g., see Devin¹ and Cursor²). However, in code generation, there isn’t yet a benchmark that accurately reflects realistic user queries for developing complete AI systems. We believe this is because of the difficulty to evaluate such complex, real-world tasks. For example, while many companies advertise their systems based on the performance of it on benchmarks like SWE-Bench (Jimenez et al.) (for automated repair) or HumanEval (Chen et al., 2021) (for algorithmic tasks), these benchmarks cover only a small bit of an actual development process. A benchmark that can evaluate the entire development process—ideally in a way that can help understand the degree to which current AI methods can reduce human labour—is missing.

Topic We chose automated AI development as our main topic. While AI and ML tasks are often more complex, they follow clear, standard procedures. For example, data processing typically comes first in an AI pipeline, and performance reporting goes at the end. We believe this topological nature can help better monitor the development process and provide useful signals to the agentic systems.

Goals An ideal benchmark should address critical issues in automated development by focusing on three key factors. First, it should reflect practical software scenarios, where tasks are often too complex for a single LLM, requiring human or agentic systems. Second, it should emphasize the development process, not just final outcomes (e.g., `pass@1` rates offer limited feedback and fail to highlight intermediate problems). Lastly, the evaluation should be computationally cost-effective and efficient, avoiding long training times or excessive manual oversight.

¹<https://www.cognition.ai/blog/introducing-devin>

²<https://www.cursor.com/>

162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215

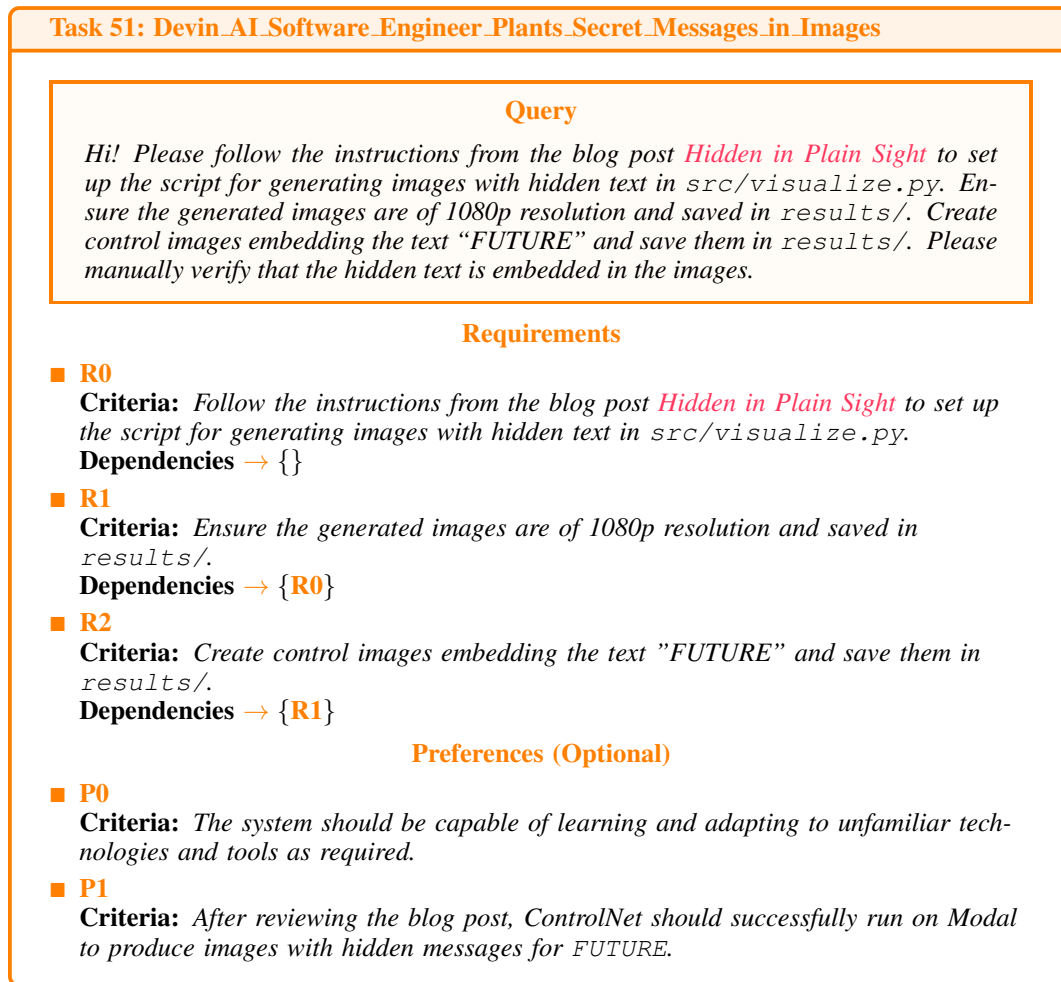


Figure 3: A task example in DevAI. This task is adapted from a real-world demo given at <https://www.cognition.ai/blog/introducing-devin>. As this example shows, task requirements in DevAI are structured as a Directed Acyclic Graph (DAG), with nodes representing individual requirements and directed edges showing dependencies. More examples are in Appendix G.


2.2 THE DevAI DATASET

Motivated by the ideas outlined above, we propose the DevAI dataset. DevAI consists of a carefully curated set of 55 tasks, each defined by (1) a plain text query that describes an AI development task; (2) a set of plain text requirements (for a total of 365 requirements), each with a set of dependencies connecting them to other requirements; and (3) a set of preferences (for a total of 125 preferences) which represent softer requirements. DevAI is structured so that an agentic system starts by receiving a user query to begin development. The system is then evaluated on how well it meets the requirements, with preferences serving as optional, softer criteria. An example of one of the DevAI tasks can be seen in Figure 3 and the full suite is available at [REDACTED]

The tasks in DevAI are relatively small-scale but cover commonly used key development techniques. As shown in Figure 2, our tasks are tagged and cover a variety of key areas in AI: supervised learning, reinforcement learning, computer vision, natural language processing, generative models, and others. Each of the tasks is a real-world problem that could be given to a research engineer, while simultaneously being relatively inexpensive computationally to run so as to reduce the cost of evaluating a method on this benchmark. Details of the sample collection and human labeling process for DevAI are provided in Appendix E.

The requirements belonging to each task represent a milestone in the comprehensive development process and are arranged as a directed acyclic graph (similar to the work by He et al. (2021)), with requirements such as visualizing results depending on correct data loading and modeling. This

Table 1: **Preliminary Statistics of AI Developers.** We compare three leading open-source code agents using metrics such as average cost, average time, and the average number of generated files.

Metric	MetaGPT (Hong et al., 2024b)	GPT-Pilot (Pythagora.io, 2023)	OpenHands (Wang et al., 2024b)
 Basic Statistics			
Version	Data Interpreter (Hong et al., 2024a)	0.2.13	CodeAct v1.9 (Wang et al.)
(1) Average Cost	\$1.19	\$3.92	\$6.38
(2) Average Time	775.29s	1622.38s	362.41s
(3) Average Input Tokens	152863	606707	1252482
(4) Average Output Tokens	28546	59707	8457
(4) Average Saved Code Files	0.42	3.84	2.53
(5) Average Saved Code Lines	11.15	273.33	96.56
(6) Average Saved Files	4.42	5.91	3.60

allows for more comprehensive non-sparse feedback than a binary success metric. Furthermore, the inclusion of milestones makes simple memorization not viable as a solution strategy as code that performs the entire task is unlikely to exist online at this time.

2.3 PRELIMINARY BENCHMARK

We first conduct experiments to collect development outcomes from different frameworks, which serve as baselines in the DevAI dataset. We test three of the most popular open-source frameworks (which we will refer to as “AI developers”): MetaGPT (Hong et al., 2024b), GPT-Pilot (Pythagora.io, 2023), and OpenHands (Wang et al., 2024b)—all selected for their strong community acceptance (each having over 30,000 stars on GitHub).

Experiment Setup All of these three systems require a language model as a back-end engine, for which we use `gpt-4o-2024-05-13`, a state-of-the-art language model. These AI developers were given a time-limit of 1800 seconds to solve each task and were forcefully halted if they exceeded this time limit (we imposed this constraint, which was visible to the AI developers, as detailed in Appendix I). We capture the outputs generated during the automated development process, including code, files, and other artifacts. Additionally, we record key decisions and actions made by the agentic systems through some custom instrumentation code, resulting in a thought trajectory for each of the agentic systems.

Analysis The basic statistics are shown in Table 1. MetaGPT is the most cost-efficient (1.19 USD), while OpenHands is the most expensive (6.38 USD). In terms of development time, OpenHands completes tasks in an average of 362.41s, while GPT-Pilot takes the longest at 1622.38s. On average, a full evaluation on DevAI with one of these three took around 210.65 USD and 14 hours to perform. While running, GPT-Pilot generates the most output tokens at 59707 tokens, whereas OpenHands processed the most at 1252482 tokens while producing the fewest at 8457 tokens. This suggests that OpenHands’s internal communication is more complicated but is more parsimonious in its decisions.

MetaGPT, while being the most cost-effective, generates fewer saved code files (0.42), suggesting it may be less inclined to save files. In contrast, GPT-Pilot generates the most saved files (3.84), reflecting a more prolific output. The difference in saved code lines, with GPT-Pilot saving 273.33 lines versus MetaGPT’s 11.15, underscores GPT-Pilot’s extensive output. Meanwhile, OpenHands, despite handling larger inputs, seems less focused on executing code to generate files, as evidenced by its lower file output (2.53 saved files). These statistics align with real user experiences (as discussed in Appendix F).

Evaluations Note that the results in Table 1 are not directly indicative of performance, but provide valuable intuition about the practical utility of DevAI. The generated workspaces (generated files, code, etc.) and trajectories are utilized in subsequent experiments to perform evaluations using Human-as-a-Judge (section 3), LLM-as-a-Judge, and Agent-as-a-Judge (section 4).

3 HUMAN-AS-A-JUDGE: MANUAL EVALUATION ON DEVAI

To determine the pragmatic validity of DevAI and to accurately estimate the actual code-generating abilities of current state-of-the-art agentic systems, in this section, we run and then manually evaluate

the application of three AI developer baselines to DevAI. In Section 4, we show how this evaluation can be automated.

Table 2: **Human-as-a-Judge for AI Developers.** (I) and (D) represent independent performance versus performance considering task dependencies. 👤 indicates multiple experts evolved, and 📄 means the evaluations use white-box testing (allowing access to the generated workspace, human-collected trajectories, and open-source codebases). The results were derived from expert judgments and deliberations (see Appendix H).

Metric	MetaGPT (Hong et al., 2024b)	GPT-Pilot (Pythagora.io, 2023)	OpenHands (Wang et al., 2024b)
👤 / 📄 Human-as-a-Judge			
(A) Requirements Met (I)	22.13%	44.80%	42.89%
(B) Requirements Met (D)	6.55%	28.96%	28.68%
(C) Self-Termination	41.81%	5.45%	54.54%
(D) Task Solve Rate	0.00%	1.81%	1.81%

3.1 BENCHMARK BASELINES BY HUMAN-AS-A-JUDGE

Human Evaluation Setup After obtaining the baseline executions and conducting basic statistical analysis, we have three expert human evaluators (referred to here by their anonymous names: 231a, 38bb, and cn90) review the outputs of AI developer baselines to assess whether each requirement was satisfied. We have two rounds of human evaluations. To capture the bias inherent in typical human evaluation (this is desirable to capture here as it represents a likely scenario in deployment), in the first round, our evaluators first discussed the basic standards but were given minimal instructions. The templates the evaluators were given for the evaluation and their self-reported post-hoc descriptions of how they resolved ambiguities are reported in Figure 12 in Appendix H.

After the initial round of human evaluations (which totaled an estimated total of 58 human hours), we asked our evaluators to discuss and reach a consensus on their assessments (which took an estimated total of 28.5 additional human hours). This consensus, achieved after long sessions of debate, was used as the final human evaluation result for each method.

Performance Analysis The results of this experiment are shown in Table 2. We found that the two best-performing methods (GPT-Pilot and OpenHands) could satisfy about 29% of the requirements (or around 44% if prerequisites are ignored) but only on one task could they meet all the requirements. This highlights that DevAI offers a considerable but appropriate level of challenge for current and future methods. Moreover, the fulfillment of intermediate requirements aligns with our expectations (see Section 2) that DevAI provides richer feedback by uncovering how agentic systems falter during the process instead of just focusing on a single performance metric at the end.

3.2 JUDGING HUMAN-AS-A-JUDGE

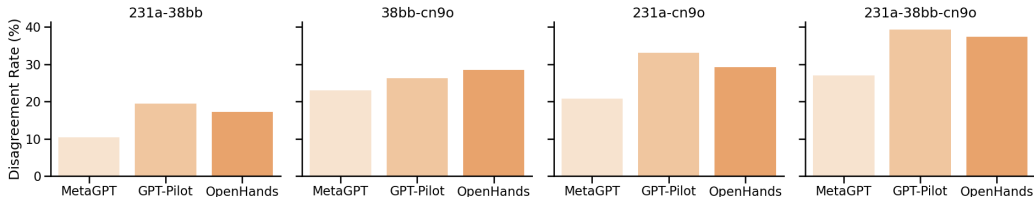


Figure 4: Between the three human evaluators, a large amount of disagreement was observed in their individual evaluations—highlighting the inherent unreliability of a single human evaluation.

Disagreement Analysis To analyze the presence of inductive bias and the reliability of the Human-as-a-Judge paradigm here, we calculate the disagreement rate between individual evaluators (shown in Figure 4). The results indicate that the disagreement rates between pairs of evaluators range from around 10% to 30%. Although each human evaluator has over five years of experience in AI research and development, the disagreement highlights the inherent challenges of evaluating AI development.

Due to the complexity of a complete AI development task, which typically involves multiple steps with varying outcomes at each step, humans can easily make errors when critical information is missed, such as environment feedback indicating small but severe coding errors or bugs. Additionally, some disagreements are not necessarily incorrect but arise from differing perspectives on how ambiguity should be resolved.

To determine if the disagreement between the three human judges is too large for them to serve as a strong baseline, we recruited ten additional experts and had them evaluate a random selection of 7 task samples. We observed that the majority vote of these additional experts had a 95.23% agreement rate with the consensus result of our other experts on these 7 tasks and a 97.67% agreement rate with the majority vote. See Appendix P.

Error Analysis As previously noted, the evaluators engaged in a round of debating after their initial evaluations until they reached a consensus on each requirement in each task (with the results of this consensus evaluation shown in Table 2).

In our Human-as-a-Judge pipeline, evaluators could be convinced by evidence from others and acknowledge their judgment errors, adjusting their answers accordingly. This can be used to approximate individual errors. If the consensus evaluation more accurately predicts any extant ground truth, we would expect the majority vote from the individual evaluations to more closely approximate this than any single evaluation, due to the fundamental properties of ensemble classifiers (see [Hastie et al. \(2009\)](#)).

While the consensus evaluation may not represent the absolute ground truth (we acknowledge that some quantity of error likely would still exist even after this procedure), we expect the consensus evaluation to more accurately approximate the extant ground truth ([Clemen, 1989](#)). If this holds, the majority vote should align more closely with the consensus than with any individual evaluation. As shown in Figure 5, this is indeed the case.

As seen in the results, although significant errors occur among all evaluators, the majority vote effectively corrects most of these errors. Notably, `cn90` made the most errors (for example, 23.77% in evaluating GPT-Pilot). After applying the majority vote from all three evaluators, the overall error rate dropped to 6.01%, demonstrating the inherent benefits of majority voting.

Conclusion Human judgment errors are inevitable. To reduce them, we suggest two methods. First, like in this work, introduce a debate round after each judgment, where individuals present evidence and either persuade others or adjust their own opinions after discussion. This is particularly important when there are only a few evaluators, as majority voting with a small group can still lead to errors (around 5% compared to consensus evaluation, as shown in Figure 5). The second approach involves assembling a larger panel of experts (more is better when their accuracy exceeds 50% ([Grofman et al., 1983](#))), with over 5 people recommended by [Hastie & Kameda \(2005\)](#); [Larrick & Soll \(2006\)](#), and relying on a majority vote. However, due to the high cost of engaging more experts and the fact that this is not always feasible in practice, we argue for the former.

4 AGENT-AS-A-JUDGE: EVALUATING AGENTS WITH AGENTS

Human evaluation, while somewhat reliable, is time-consuming and requires significant expertise. To address this, we propose the Agent-as-a-Judge framework. If such an agentic system could evaluate like a human, it would reduce the need for human involvement, eliminating the trade-off between evaluation thoroughness and effort.

Labeler	GPT-Pilot	MetaGPT Baseline	OpenHands
231a	12.57	7.92	10.93
38bb	9.02	7.38	10.11
cn90	23.77	16.67	21.86
majority_vote	6.01	4.92	5.74

Figure 5: Mismatch between the individual evaluations and the consensus evaluation. Note that the majority vote showed the smallest deviation from the consensus evaluation.

Table 3: **AI Judges and Their Shift/Alignment with Human-as-a-Judge.** We compare the results of LLM-as-a-Judge and Agent-as-a-Judge with Human-as-a-Judge. **(I)** represents performance on independent tasks, while **(D)** represents performance considering task dependencies. **Note:** gray-box settings use carefully manually collected trajectory data (which is nearly inaccessible in practical situations, see Appendix J). In contrast, black-box setting doesn't need to access to such data. The red scores represent the absolute judge shift compared with Human-as-a-Judge (e.g., 2.74%).

Metric	MetaGPT (Hong et al., 2024b)	GPT-Pilot (Pythagora.io, 2023)	OpenHands (Wang et al., 2024b)
 LLM-as-a-Judge			
(a) Requirements Met (I)	19.39% (2.74%)	12.56% (32.24%)	11.47% (31.42%)
(b) Requirements Met (D)	1.63% (4.92%)	4.09% (24.87%)	2.18% (26.50%)
(c) Task Solve Rate	0.0% (0.0%)	0.0% (1.81%)	0.0% (1.81%)
Alignment Rate \uparrow	84.15%	65.30%	60.38%
 Agent-as-a-Judge			
(I) Requirements Met (I)	25.40% (3.26%)	53.00% (8.20%)	42.62% (0.27%)
(II) Requirements Met (D)	5.73% (0.81%)	39.89% (10.93%)	26.50% (2.17%)
(III) Task Solve Rate	0.0% (0.0%)	5.45% (3.64%)	1.81% (0.00%)
Alignment Rate \uparrow	88.52%	83.88%	90.44%
 LLM-as-a-Judge			
(a) Requirements Met (I)	28.68% (6.55%)	38.79% (4.10%)	43.16% (0.27%)
(b) Requirements Met (D)	17.75% (11.20%)	33.06% (4.10%)	32.24% (3.56%)
(c) Task Solve Rate	1.81% (1.81%)	3.63% (1.82%)	7.27% (5.46%)
Alignment Rate \uparrow	68.86%	71.85%	70.76%
 Agent-as-a-Judge			
(I) Requirements Met (I)	23.49% (1.35%)	46.44% (1.64%)	43.44% (0.54%)
(II) Requirements Met (D)	6.01% (0.54%)	30.60% (1.64%)	28.14% (0.53%)
(III) Task Solve Rate	0.0% (0.00%)	5.45% (3.64%)	3.63% (1.82%)
Alignment Rate \uparrow	92.07%	86.61%	90.16%
 Human-as-a-Judge			
Alignment Rate (38bb)	92.63%	90.98%	89.89%
Alignment Rate (cn9o)	83.33%	76.23%	78.15%
Alignment Rate (231a)	92.07%	87.43%	89.07%
Average of individuals	89.34%	84.88%	85.70%
Alignment Rate (Majority Vote)	95.08%	93.98%	94.26%

4.1 PROOF-OF-CONCEPT

Based on our prior experiences with agent design and by imitating the human evaluation process, we initially designed eight modular, interacting components that form the foundation of our Proof-of-Concept for the Agent-as-a-Judge.

(1) **The graph module** constructs a graph that captures the entire structure of the project, including files, modules, and dependencies. It can also break down chunks of code into code snippets. (2) **The locate module** identifies the specific folder or file referred to by a requirement. (3) **The read module** goes beyond simple file parsing, supporting the reading and understanding of multimodal data across 33 different formats, including code, images, videos and documents. This allows the agent to cross-reference various data streams and check different kinds of requirement. (4) **The search module** offers a contextual understanding of code and can quickly retrieve highly relevant code snippets, as well as the nuances behind them (e.g., hidden dependencies). (5) **The retrieve module** extracts information from long texts, identifying relevant segments in trajectories. With context from the above, (6) **the ask module** determines whether a given requirement is satisfied. (7) **The memory module** stores historical judgment information, allowing the agent to build on past evaluations. Finally, (8) **the planning module** plans out actions, allowing the agent to strategize and sequence tasks based on the current state and project goals.

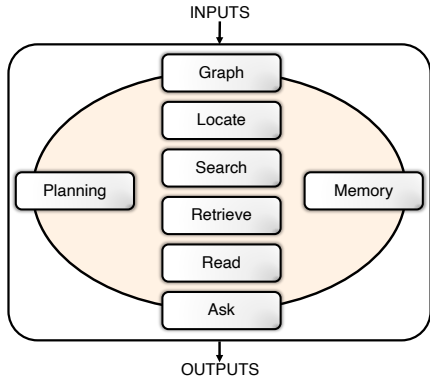


Figure 6: Initial diagram of Agent-as-a-Judge.

Our initial design of the Agent-as-a-Judge is shown in Figure 6. After conducting comprehensive ablation studies, we found that the modular combination of (1), (2), (3), (5), and (6) achieved the

highest performance (see Appendix C). A sample of the dynamic evidence collected by the Agent-as-a-Judge is shown in Appendix M. We hypothesize this is because Agent-as-a-Judge needs high-quality factual information and is sensitive to noise. For example, while our design of the planning module introduces promising decision-making for future actions, the procedure is unstable. Initially, we hoped that historical information from the memory module would help to assess current requirements. However, it proved detrimental, as any errors in previous judgments could lead to a chain of errors, negatively affecting current decisions. Besides, the current workspaces generated by developer agents, having only hundreds of lines of code, cannot fully benefit from the search module. The details of these findings are explained in Appendix K. Note that a perfect Agent-as-a-Judge is not the focus of this proof of concept, and thus, we leave the utilization of advanced agentic optimization methods for Agent-as-a-Judge, such as automated prompt optimization and workflow design (Zhuge et al.; Hu et al., 2024), for future work.

4.2 JUDGING AGENT-AS-A-JUDGE AND LLM-AS-A-JUDGE

Judge Shift Judge Shift measures deviation from the Human-as-a-Judge consensus results, with lower values indicating a closer alignment. As shown in table 3, Agent-as-a-Judge consistently outperforms LLM-as-a-Judge across tasks, particularly those with task dependencies. For example, in Requirement (I), Agent-as-a-Judge shows a Judge Shift as low as 0.27%, while LLM-as-a-Judge reaches 31.24% for OpenHands. This underscores Agent-as-a-Judge’s stability and suitability for meeting task requirements. Furthermore, in the gray-box setting, both Agent-as-a-Judge and LLM-as-a-Judge show even better results than their performance in the black-box setting.

Alignment Rate The Alignment Rate reflects how closely the AI Judges’ evaluations align with human consensus across all 365 requirements. It is defined as the percentage of requirement evaluations that are the same as the Human-as-a-Judge consensus evaluation. Compared to LLM-as-a-Judge, Agent-as-a-Judge consistently achieves a higher Alignment Rate, closely matching human judgments. For example, when evaluating OpenHands, Agent-as-a-Judge reaches 92.07% and 90.44%, surpassing LLM-as-a-Judge’s 70.76% and 60.38% in both gray-box and black-box settings. This suggests that Agent-as-a-Judge is more accurate and human-aligned.

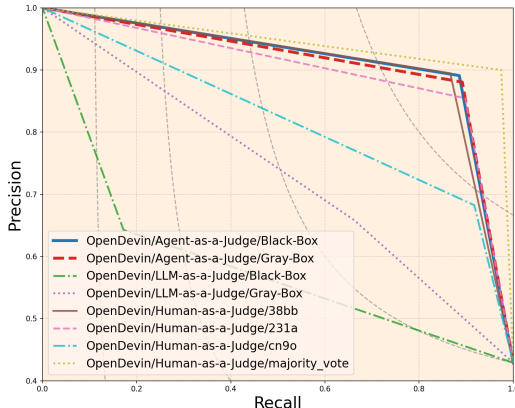


Figure 7: PR Curves comparing judge Methods.

PR Curves Judging developer agents is a class-imbalanced task, where meeting requirements is much rarer than failing. Metrics like judge shift and alignment rate can be misleading. For example, since MetaGPT rarely meets requirements, LLM-as-a-Judge easily identifies most cases as negative (achieving 84.15% in the black-box setting). PR Curves offer a clearer performance measure by balancing precision and recall. This shows that, in some cases, Agent-as-a-Judge can nearly replace human evaluators. Our observations indicate the relative reliability of evaluation methods as: LLM-as-a-Judge < Single-Human-as-a-Judge < Agent-as-a-Judge < Ensemble of Human Judges. Future advancements in foundation models and Agent-as-a-Judge designs may shift this order.

4.3 ABLATIONS FOR AGENT-AS-A-JUDGE

We conduct ablations to evaluate the impact of adding different components on Agent-as-a-Judge’s performance. The components analyzed include ask, graph, read, locate, and retrieve. The component ablation study for Agent-as-a-Judge reveals key insights into the performance gains from adding specific functionalities. With only

Table 4: **Component Ablation Studies for Agent-as-a-Judge.** We analyze the impact of adding various components (ask, graph, read, locate, and retrieve) on the performance of Agent-as-a-Judge for judging OpenHands.

Metric	+ ask	+ graph	+ read	+ locate	+ retrieve
Agent-as-a-Judge Performance					
Alignment Rate	65.03%	75.95%	82.24%	90.44%	90.16%

ask component, the agent achieves a 65.03% alignment rate. Adding the graph component increases performance to 75.95%, as the agent can better understand relationships between files. The introduction of `read` further improves the alignment rate to 82.24%, reflecting the value of direct access to the contents of the file. Incorporating `locate` brings a substantial boost to 90.44%, as the agent can efficiently target files relevant to the requirements. Adding `retrieve` does not always provide a significant benefit in this case. We found the `retrieve` module effective for judging MetaGPT and GPT-Pilot, as it provides valuable trajectory information (as shown in Table 3). However, it is less effective for OpenHands, which sometimes fails to execute files, resulting in missing responses. In such cases, judgment without trajectories remains viable.

4.4 COST ANALYSIS

Our three evaluators a self-reported total of 86.5 hours. With a 15 USD hourly wage (assuming this would buy a subject expert in AI), a full evaluation under DevAI would cost around 1297.50 USD. In comparison, Agent-as-a-Judge cost only 30.58 USD (2.29%) in API calls and took only 118.43 minutes (2.36%). LLM-as-a-Judge was faster at 10.99 minutes, but due to the absence of intelligent context selection by the Agent-as-a-Judge’s modules, it still cost 29.63 USD.

5 RELATED WORK

Agentic systems is a highly active research area, so we only detail the most closely related works here. We provide a treatment of the marginally less relevant related works in Appendix D.

AI Developers AI in software development is growing fast (Liu et al., 2024). AI-driven developers have been applied to directly imitate software companies (Hong et al., 2024b; Qian et al., 2024a), debug code (Yang et al., 2024a), run data science methods (Guo et al.; Hong et al., 2024a; Li et al., 2024; Qiao et al., 2023), and even write academic papers (Lu et al., 2024a).

Benchmarks for AI developments Benchmarks like MAgentBench (Huang et al., 2024), ML-Bench (Liu et al., 2023), and SUPER (Bogin et al., 2024) all focus on benchmarking agentic systems using AI tasks. However, DevAI distinguishes itself from all of these by focusing on realistic user queries that target a complete development cycle. It further includes a more comprehensive evaluation with multiple hierarchical requirements and preferences for each task. Comparatively, MAgentBench (Huang et al., 2024), for example, focuses on final performance for a limited set of well-known tasks, which risks overfitting and fails to assess a system’s generalization or adaptability.

AI Judges Several works have looked at using AI systems as judges. The work by Chan et al.; Zhao et al. (2024), for example, extends LLM-as-a-Judge to have multiple LLMs in their evaluation process for conversational tasks. Unlike Agent-as-a-Judge, they employ a trivial agentic system and apply it only to evaluate LLMs under traditional evaluation setups. In contrast, (Lu et al., 2024b) uses a single LLM-based evaluator but, unlike LLM-as-a-Judge, applies this to multimodal tasks rather than just for evaluating LLMs. Less relevant are frameworks like those by Chen et al. (2024a); Arora et al. (2024); Mündler et al. (2024), where intermediate signals are used during coding development.

6 DISCUSSION AND CONCLUSION

Discussion A key power of the Agent-as-a-Judge that we have not exploited here is the feedback it provides being directly used by the agentic system being evaluated. Perhaps the greatest strength of the Agent-as-a-Judge framework is that an agentic system could use it to fix issues in their solutions to complex multistage problems on the fly—something older delayed feedback methods did not permit. In addition, a cycle of mutual improvement between Agent-as-a-Judge and the evaluated agents, where both evolve together through iterative feedback, presents a promising perspective.

Conclusion In this work, we introduced the Agent-as-a-Judge method to use agentic systems to evaluate agentic systems. We simultaneously released DevAI: a new benchmark that evaluates the code-generating ability of agentic systems on complete AI development tasks when used with Agent-as-a-Judge. We went on to show that Agent-as-a-Judge outperforms existing methods on this task and that it performs similarly to an ensemble of expert human evaluators. Altogether, we believe that the above opens the door for scaling up agentic far more than before.

REPRODUCIBILITY STATEMENT

The nature of this work necessitates the public release of the DevAI dataset and the implementation of Agent-as-a-Judge as part of that. The authors are committed to open science and will be doing so upon paper acceptance. In tandem with the details provided in the various Appendices, this should allow full reproducibility of the results shown in this paper.

ETHICS STATEMENT

Our work, like similar works, aims to reduce human labour costs while maintaining the rigour needed for meaningful science. This should serve to make it easier for the field to adopt more accessible and open evaluation methods. The DevAI dataset is sourced from widely used general-purpose datasets, and all examples are annotated by experts to minimize bias. Furthermore, our dataset enhances transparency in the field of automated AI development, contributing to more open, equitable, and responsible progress in AI research. We do not foresee any obvious exceptional ethical implications of this work beyond the above.

REFERENCES

- Anthropic. The claude 3 model family: Opus, sonnet, haiku, 2024. URL https://www-cdn.anthropic.com/de8ba9b01c9ab7cbabf5c33b80b7bbc618857627/Model_Card_Claude_3.pdf.
- Daman Arora, Atharv Sonwane, Nalin Wadhwa, Abhav Mehrotra, Saiteja Utpala, Ramakrishna Bairi, Aditya Kanade, and Nagarajan Natarajan. MASAI: Modular architecture for software-engineering ai agents. *arXiv preprint arXiv:2406.11638*, 2024.
- Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. Program synthesis with large language models. *arXiv preprint arXiv:2108.07732*, 2021.
- Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu, Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren, Xuancheng Ren, Chuanqi Tan, Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Wei Wang, Shengguang Wu, Benfeng Xu, Jin Xu, An Yang, Hao Yang, Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu, Hongyi Yuan, Zheng Yuan, Jianwei Zhang, Xingxuan Zhang, Yichang Zhang, Zhenru Zhang, Chang Zhou, Jingren Zhou, Xiaohuan Zhou, and Tianhang Zhu. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023.
- Anna Bavaresco, Raffaella Bernardi, Leonardo Bertolazzi, Desmond Elliott, Raquel Fernández, Albert Gatt, Esam Ghaleb, Mario Giulianelli, Michael Hanna, Alexander Koller, et al. LLMs instead of human judges? a large scale empirical study across 20 NLP evaluation tasks. *arXiv preprint arXiv:2406.18403*, 2024.
- Ben Bogin, Kejuan Yang, Shashank Gupta, Kyle Richardson, Erin Bransom, Peter Clark, Ashish Sabharwal, and Tushar Khot. SUPER: Evaluating agents on setting up and executing tasks from research repositories. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pp. 12622–12645, 2024.
- Federico Cassano, John Gouwar, Daniel Nguyen, Sydney Nguyen, Luna Phipps-Costin, Donald Pinckney, Ming-Ho Yee, Yangtian Zi, Carolyn Jane Anderson, Molly Q Feldman, et al. MultiPL-E: a scalable and polyglot approach to benchmarking neural code generation. *IEEE Transactions on Software Engineering*, 49(7):3675–3691, 2023.
- Federico Cassano, Luisa Li, Akul Sethi, Noah Shinn, Abby Brennan-Jones, Jacob Ginesin, Edward Berman, George Chakhnashvili, Anton Lozhkov, Carolyn Jane Anderson, and Arjun Guha. Can it edit? evaluating the ability of large language models to follow code editing instructions. In *First Conference on Language Modeling*, 2024. URL <https://openreview.net/forum?id=D06yk3DBas>.

- 594 Chi-Min Chan, Weize Chen, Yusheng Su, Jianxuan Yu, Wei Xue, Shanghang Zhang, Jie Fu, and
595 Zhiyuan Liu. ChatEval: Towards better LLM-based evaluators through multi-agent debate. In
596 *The Twelfth International Conference on Learning Representations*.
597
- 598 Harrison Chase. LangChain. <https://github.com/hwchase17/langchain>, 2022.
- 599 Dong Chen, Shaoxin Lin, Muhan Zeng, Daoguang Zan, Jian-Gang Wang, Anton Cheshkov, Jun
600 Sun, Hao Yu, Guoliang Dong, Artem Aliev, et al. CodeR: Issue resolving with multi-agent and
601 task graphs. *arXiv preprint arXiv:2406.01304*, 2024a.
- 602 Dongping Chen, Ruoxi Chen, Shilin Zhang, Yaochen Wang, Yinuo Liu, Huichi Zhou, Qihui Zhang,
603 Yao Wan, Pan Zhou, and Lichao Sun. MLLM-as-a-Judge: Assessing multimodal LLM-as-
604 a-Judge with vision-language benchmark. In *Forty-first International Conference on Machine*
605 *Learning*.
- 607 Dongping Chen, Yue Huang, Siyuan Wu, Jingyu Tang, Liuyi Chen, Yilin Bai, Zhigang He, Chenlong
608 Wang, Huichi Zhou, Yiqiang Li, et al. GUI-WORLD: A dataset for GUI-oriented multimodal
609 LLM-based agents. *arXiv preprint arXiv:2406.10819*, 2024b.
- 610 Guiming Hardy Chen, Shunian Chen, Ziche Liu, Feng Jiang, and Benyou Wang. Humans or LLMs
611 as the judge? a study on judgement bias. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung
612 Chen (eds.), *Proceedings of the 2024 Conference on Empirical Methods in Natural Language*
613 *Processing*, pp. 8301–8327, Miami, Florida, USA, November 2024c. Association for Computa-
614 tional Linguistics. URL <https://aclanthology.org/2024.emnlp-main.474>.
- 615 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared
616 Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large
617 language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- 618 Ching-An Cheng, Allen Nie, and Adith Swaminathan. Trace is the new AutoDiff—unlocking ef-
619 ficient optimization of computational workflows. In *ICML 2024 Workshop on Automated Rein-*
620 *forcement Learning: Exploring Meta-Learning, AutoML, and LLMs*.
- 621 Robert T Clemen. Combining forecasts: A review and annotated bibliography. *International journal*
622 *of forecasting*, 5(4):559–583, 1989.
- 623 Corinna Cortes. Support-vector networks. *Machine Learning*, 1995.
- 624 Yijiang River Dong, Tiancheng Hu, and Nigel Collier. Can LLM be a personalized judge? In
625 Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.), *Findings of the Association for*
626 *Computational Linguistics: EMNLP 2024*, pp. 10126–10141, Miami, Florida, USA, November
627 2024. Association for Computational Linguistics. URL [https://aclanthology.org/2](https://aclanthology.org/2024.findings-emnlp.592)
628 [024.findings-emnlp.592](https://aclanthology.org/2024.findings-emnlp.592).
- 629 Zhuoyun Du, Chen Qian, Wei Liu, Zihao Xie, Yifei Wang, Yufan Dang, Weize Chen, and Cheng
630 Yang. Multi-agent software development through cross-team collaboration. *arXiv preprint*
631 *arXiv:2406.08979*, 2024.
- 632 Usama Fayyad, Gregory Piatetsky-Shapiro, and Padhraic Smyth. From data mining to knowledge
633 discovery in databases. *AI magazine*, 17(3):37–37, 1996.
- 634 Jinlan Fu, See Kiong Ng, Zhengbao Jiang, and Pengfei Liu. GPTScore: Evaluate as you desire.
635 In *Proceedings of the 2024 Conference of the North American Chapter of the Association for*
636 *Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pp. 6556–
637 6576, 2024.
- 638 Charles Goodhart. *Monetary relationships: a view from Threadneedle Street*. University of War-
639 wick, 1976.
- 640 Significant Gravitass. Auto-GPT. *GitHub repository*, 2023.
- 641 Bernard Grofman, Guillermo Owen, and Scott L Feld. Thirteen theorems in search of the truth.
642 *Theory and decision*, 15(3):261–278, 1983.

- 648 Siyuan Guo, Cheng Deng, Ying Wen, Hechang Chen, Yi Chang, and Jun Wang. DS-Agent: Auto-
649 mated data science by empowering large language models with case-based reasoning. In *Forty-
650 first International Conference on Machine Learning*.
651
- 652 Md Mahim Anjum Haque. *FixEval: Execution-based evaluation of program fixes for competitive
653 programming problems*. PhD thesis, Virginia Tech, 2023.
654
- 655 Reid Hastie and Tatsuya Kameda. The robust beauty of majority rules in group decisions. *Psycho-
656 logical review*, 112(2):494, 2005.
657
- 658 Trevor Hastie, Robert Tibshirani, and Jerome H. Friedman. *The Elements of Statistical Learning:
659 Data Mining, Inference, and Prediction*. Springer, 2 edition, 2009. doi: 10.1007/978-0-387-84
660 858-7.
- 661 Xin He, Kaiyong Zhao, and Xiaowen Chu. AutoML: A survey of the state-of-the-art. *Knowledge-
662 based systems*, 212:106622, 2021.
663
- 664 Dan Hendrycks, Steven Basart, Saurav Kadavath, Mantas Mazeika, Akul Arora, Ethan Guo, Collin
665 Burns, Samir Puranik, Horace He, Dawn Song, et al. Measuring coding challenge competence
666 with APPS. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and
667 Benchmarks Track (Round 2)*.
- 668 S Hochreiter. Long short-term memory. *Neural Computation MIT-Press*, 1997.
669
- 670 Sirui Hong, Yizhang Lin, Bangbang Liu, Binhao Wu, Danyang Li, Jiaqi Chen, Jiayi Zhang, Jinlin
671 Wang, Lingyao Zhang, Mingchen Zhuge, et al. Data interpreter: An LLM agent for data science.
672 *arXiv preprint arXiv:2402.18679*, 2024a.
- 673 Sirui Hong, Mingchen Zhuge, Jonathan Chen, Xiawu Zheng, Yuheng Cheng, Jinlin Wang, Ceyao
674 Zhang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, et al. MetaGPT: Meta programming for
675 a multi-agent collaborative framework. In *The Twelfth International Conference on Learning
676 Representations*, 2024b.
677
- 678 Shengran Hu, Cong Lu, and Jeff Clune. Automated design of agentic systems. *arXiv preprint
679 arXiv:2408.08435*, 2024.
680
- 681 Dong Huang, Qingwen Bu, Jie M Zhang, Michael Luck, and Heming Cui. AgentCoder:
682 Multi-agent-based code generation with iterative testing and optimisation. *arXiv preprint
683 arXiv:2312.13010*, 2023.
- 684 Qian Huang, Jian Vora, Percy Liang, and Jure Leskovec. MLAGentBench: Evaluating language
685 agents on machine learning experimentation. In *Forty-first International Conference on Machine
686 Learning*, 2024.
687
- 688 Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando
689 Solar-Lezama, Koushik Sen, and Ion Stoica. LiveCodeBench: Holistic and contamination free
690 evaluation of large language models for code. *arXiv preprint arXiv:2403.07974*, 2024.
691
- 692 Peter Jansen, Marc-Alexandre Côté, Tushar Khot, Erin Bransom, Bhavana Dalvi Mishra, Bod-
693 hisattwa Prasad Majumder, Oyvind Tafjord, and Peter Clark. DISCOVERYWORLD: A virtual
694 environment for developing and evaluating automated scientific discovery agents. *arXiv preprint
695 arXiv:2406.06769*, 2024.
- 696 Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik R
697 Narasimhan. SWE-bench: Can language models resolve real-world GitHub issues? In *The
698 Twelfth International Conference on Learning Representations*.
699
- 700 Haolin Jin, Linghan Huang, Haipeng Cai, Jun Yan, Bo Li, and Huaming Chen. From LLMs to
701 LLM-based agents for software engineering: A survey of current, challenges and future. *arXiv
preprint arXiv:2408.02479*, 2024.

- 702 Omar Khattab, Arnav Singhvi, Paridhi Maheshwari, Zhiyuan Zhang, Keshav Santhanam, Sri Vard-
703 hamanan A, Saiful Haq, Ashutosh Sharma, Thomas T. Joshi, Hanna Moazam, Heather Miller,
704 Matei Zaharia, and Christopher Potts. DSPy: Compiling declarative language model calls into
705 state-of-the-art pipelines. In *The Twelfth International Conference on Learning Representations*,
706 2024. URL <https://openreview.net/forum?id=sY5N0zY5Od>.
- 707 Yuhang Lai, Chengxi Li, Yiming Wang, Tianyi Zhang, Ruiqi Zhong, Luke Zettlemoyer, Wen-tau
708 Yih, Daniel Fried, Sida Wang, and Tao Yu. DS-1000: A natural and reliable benchmark for data
709 science code generation. In *International Conference on Machine Learning*, pp. 18319–18345.
710 PMLR, 2023.
- 711 LangChain-AI. LangGraph. <https://github.com/langchain-ai/langgraph>, 2024.
- 712 Richard P Larrick and Jack B Soll. Intuitions about combining opinions: Misappreciation of the
713 averaging principle. *Management science*, 52(1):111–127, 2006.
- 714 V Levenshtein. Binary codes capable of correcting deletions, insertions, and reversals. *Proceedings*
715 *of the Soviet physics doklady*, 1966.
- 716 Guohao Li, Hasan Hammoud, Hani Itani, Dmitrii Khizbullin, and Bernard Ghanem. CAMEL:
717 Communicative agents for” mind” exploration of large language model society. *Advances in*
718 *Neural Information Processing Systems*, 36:51991–52008, 2023.
- 719 Ruo Chen Li, Teerth Patel, Qingyun Wang, and Xinya Du. MLR-Copilot: Autonomous machine
720 learning research based on large language models agents. *arXiv preprint arXiv:2408.14033*, 2024.
- 721 Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond, Tom
722 Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, et al. Competition-level code generation
723 with AlphaCode. *Science*, 378(6624):1092–1097, 2022.
- 724 Junwei Liu, Kaixin Wang, Yixuan Chen, Xin Peng, Zhenpeng Chen, Lingming Zhang, and Yiling
725 Lou. Large language model-based agents for software engineering: A survey. *arXiv preprint*
726 *arXiv:2409.02977*, 2024.
- 727 Tianyang Liu, Canwen Xu, and Julian McAuley. RepoBench: Benchmarking repository-level code
728 auto-completion systems. In *The Twelfth International Conference on Learning Representations*,
729 a.
- 730 Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei, Hanyu Lai, Yu Gu, Hangliang Ding,
731 Kaiwen Men, Kejuan Yang, et al. AgentBench: Evaluating llms as agents. In *The Twelfth Inter-*
732 *national Conference on Learning Representations*, b.
- 733 Yuliang Liu, Xiangru Tang, Zefan Cai, Junjie Lu, Yichi Zhang, Yanjun Shao, Zexuan Deng, Helan
734 Hu, Zengxian Yang, Kaikai An, et al. ML-Bench: Large language models leverage open-source
735 libraries for machine learning tasks. *arXiv preprint arXiv:2311.09835*, 2023.
- 736 Chris Lu, Cong Lu, Robert Tjarko Lange, Jakob Foerster, Jeff Clune, and David Ha. The AI scien-
737 tist: Towards fully automated open-ended scientific discovery. *arXiv preprint arXiv:2408.06292*,
738 2024a.
- 739 Yujie Lu, Xianjun Yang, Xiujun Li, Xin Eric Wang, and William Yang Wang. LLMscore: Unveiling
740 the power of large language models in text-to-image synthesis evaluation. *Advances in Neural*
741 *Information Processing Systems*, 36, 2024b.
- 742 Niels Mündler, Mark Niklas Müller, Jingxuan He, and Martin Vechev. Code agents are state of the
743 art software testers. *arXiv preprint arXiv:2406.12952*, 2024.
- 744 OpenAI. GPT-4 technical report, 2023.
- 745 Jooyong Park. Constructive multiple-choice testing system. *British Journal of Educational*
746 *Technology*, 41(6):1054–1064, 2010. doi: <https://doi.org/10.1111/j.1467-8535.2010.01058>
747 [.x](https://doi.org/10.1111/j.1467-8535.2010.01058.x). URL [https://bera-journals.onlinelibrary.wiley.com/doi/abs/10.1](https://bera-journals.onlinelibrary.wiley.com/doi/abs/10.1111/j.1467-8535.2010.01058.x)
748 [111/j.1467-8535.2010.01058.x](https://doi.org/10.1111/j.1467-8535.2010.01058.x).

- 756 Huy Nhat Phan, Phong X Nguyen, and Nghi DQ Bui. Hyperagent: Generalist software engineering
757 agents to solve coding tasks at scale. *arXiv preprint arXiv:2409.16299*, 2024.
758
- 759 Pythagora.io. GPT-Pilot: Your ai copilot for software development. [https://github.com/P](https://github.com/Pythagora-io/gpt-pilot)
760 [ythagora-io/gpt-pilot](https://github.com/Pythagora-io/gpt-pilot), 2023. URL [https://github.com/Pythagora-io/gp](https://github.com/Pythagora-io/gpt-pilot)
761 [t-pilot](https://github.com/Pythagora-io/gpt-pilot). GitHub repository.
- 762 Chen Qian, Wei Liu, Hongzhang Liu, Nuo Chen, Yufan Dang, Jiahao Li, Cheng Yang, Weize
763 Chen, Yusheng Su, Xin Cong, Juyuan Xu, Dahai Li, Zhiyuan Liu, and Maosong Sun. Chat-
764 Dev: Communicative agents for software development. In Lun-Wei Ku, Andre Martins, and
765 Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Com-*
766 *putational Linguistics (Volume 1: Long Papers)*, pp. 15174–15186, Bangkok, Thailand, August
767 2024a. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.810. URL
768 <https://aclanthology.org/2024.acl-long.810>.
- 769 Chen Qian, Zihao Xie, Yifei Wang, Wei Liu, Yufan Dang, Zhuoyun Du, Weize Chen, Cheng Yang,
770 Zhiyuan Liu, and Maosong Sun. Scaling large-language-model-based multi-agent collaboration.
771 *arXiv preprint arXiv:2406.07155*, 2024b.
772
- 773 Bo Qiao, Liqun Li, Xu Zhang, Shilin He, Yu Kang, Chaoyun Zhang, Fangkai Yang, Hang Dong,
774 Jue Zhang, Lu Wang, et al. Taskweaver: A code-first agent framework. *arXiv preprint*
775 *arXiv:2311.17541*, 2023.
- 776 Vyas Raina, Adian Liusie, and Mark Gales. Is LLM-as-a-judge robust? investigating universal
777 adversarial attacks on zero-shot LLM assessment. In Yaser Al-Onaizan, Mohit Bansal, and Yun-
778 Nung Chen (eds.), *Proceedings of the 2024 Conference on Empirical Methods in Natural Lan-*
779 *guage Processing*, pp. 7499–7517, Miami, Florida, USA, November 2024. Association for Com-
780 putational Linguistics. URL <https://aclanthology.org/2024.emnlp-main.427>.
- 781 Nils Reimers and Iryna Gurevych. Sentence-BERT: Sentence embeddings using Siamese BERT-
782 networks. In Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan (eds.), *Proceedings of*
783 *the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th In-*
784 *ternational Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 3982–
785 3992, Hong Kong, China, November 2019. Association for Computational Linguistics. doi:
786 10.18653/v1/D19-1410. URL <https://aclanthology.org/D19-1410>.
- 787 Stephen Robertson, Hugo Zaragoza, et al. The probabilistic relevance framework: BM25 and be-
788 yond. *Foundations and Trends in Information Retrieval*, 3(4):333–389, 2009.
789
- 790 Lin Shi, Weicheng Ma, and Soroush Vosoughi. Judging the judges: A systematic investigation of
791 position bias in pairwise comparative assessments by LLMs. *arXiv preprint arXiv:2406.07791*,
792 2024.
- 793 Linxin Song, Jiale Liu, Jieyu Zhang, Shaokun Zhang, Ao Luo, Shijian Wang, Qingyun Wu, and
794 Chi Wang. Adaptive in-conversation team building for language model agents. *arXiv preprint*
795 *arXiv:2405.19425*, 2024.
796
- 797 Weihao Tan, Ziluo Ding, Wentao Zhang, Boyu Li, Bohan Zhou, Junpeng Yue, Haochong Xia,
798 Jiechuan Jiang, Longtao Zheng, Xinrun Xu, et al. Towards general computer control: A mul-
799 timodal agent for Red Dead Redemption II as a case study. In *ICLR 2024 Workshop on Large*
800 *Language Model (LLM) Agents*.
- 801 Wei Tao, Yucheng Zhou, Wenqiang Zhang, and Yu Cheng. MAGIS: LLM-Based multi-agent frame-
802 work for GitHub issue resolution. *arXiv preprint arXiv:2403.17927*, 2024.
803
- 804 Aman Singh Thakur, Kartik Choudhary, Venkat Srinik Ramayapally, Sankaran Vaidyanathan, and
805 Dieuwke Hupkes. Judging the judges: Evaluating alignment and vulnerabilities in llms-as-judges.
806 *arXiv preprint arXiv:2406.12624*, 2024.
- 807 Runchu Tian, Yining Ye, Yujia Qin, Xin Cong, Yankai Lin, Yinxu Pan, Yesai Wu, Hui Haotian, Liu
808 Weichuan, Zhiyuan Liu, and Maosong Sun. DebugBench: Evaluating debugging capability of
809 large language models. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Findings of*
the Association for Computational Linguistics: ACL 2024, pp. 4173–4198, Bangkok, Thailand,

- 810 August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-acl.24
 811 7. URL <https://aclanthology.org/2024.findings-acl.247>.
 812
- 813 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko-
 814 lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open founda-
 815 tion and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- 816 Michele Tufano, Anisha Agarwal, Jinu Jang, Roshanak Zilouchian Moghaddam, and Neel Sundare-
 817 san. AutoDev: Automated AI-driven development. *arXiv preprint arXiv:2403.08299*, 2024.
- 818 Junyang Wang, Haiyang Xu, Haitao Jia, Xi Zhang, Ming Yan, Weizhou Shen, Ji Zhang, Fei Huang,
 819 and Jitao Sang. Mobile-agent-v2: Mobile device operation assistant with effective navigation via
 820 multi-agent collaboration. *arXiv preprint arXiv:2406.01014*, 2024a.
- 821 Xingyao Wang, Yangyi Chen, Lifan Yuan, Yizhe Zhang, Yunzhu Li, Hao Peng, and Heng Ji. Exe-
 822 cutable code actions elicit better LLM agents. In *ICLR 2024 Workshop on Large Language Model*
 823 *(LLM) Agents*.
 824
- 825 Xingyao Wang, Boxuan Li, Yufan Song, Frank F Xu, Xiangru Tang, Mingchen Zhuge, Jiayi Pan,
 826 Yueqi Song, Bowen Li, Jaskirat Singh, et al. OpenDevin: An open platform for ai software
 827 developers as generalist agents. *arXiv preprint arXiv:2407.16741*, 2024b.
- 828 Rüdiger Wirth and Jochen Hipp. CRISP-DM: Towards a standard process model for data mining.
 829 In *Proceedings of the 4th international conference on the practical applications of knowledge*
 830 *discovery and data mining*, volume 1, pp. 29–39. Manchester, 2000.
- 831 Michael Wooldridge. Intelligent agents. *Multiagent systems: A modern approach to distributed*
 832 *artificial intelligence*, 1:27–73, 1999.
 833
- 834 Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Shaokun Zhang, Erkang Zhu, Beibin Li,
 835 Li Jiang, Xiaoyun Zhang, and Chi Wang. AutoGen: Enabling next-gen llm applications via
 836 multi-agent conversation framework. *arXiv preprint arXiv:2308.08155*, 2023.
- 837 Yiran Wu, Tianwei Yue, Shaokun Zhang, Chi Wang, and Qingyun Wu. StateFlow: Enhancing LLM
 838 task-solving through state-driven workflows. In *First Conference on Language Modeling*, 2024.
 839 URL <https://openreview.net/forum?id=3nTbuygoop>.
 840
- 841 Chunqiu Steven Xia, Yinlin Deng, Soren Dunn, and Lingming Zhang. Agentless: Demystifying
 842 LLM-based software engineering agents. *arXiv preprint arXiv:2407.01489*, 2024.
- 843 Tianqi Xu, Linyao Chen, Dai-Jie Wu, Yanjun Chen, Zecheng Zhang, Xiang Yao, Zhiqiang Xie,
 844 Yongchao Chen, Shilong Liu, Bochen Qian, et al. CRAB: Cross-environment agent benchmark
 845 for multimodal language model agents. *arXiv preprint arXiv:2407.01511*, 2024.
- 846 John Yang, Carlos E Jimenez, Alexander Wettig, Kilian Lieret, Shunyu Yao, Karthik R Narasimhan,
 847 and Ofir Press. SWE-agent: Agent-computer interfaces enable automated software engineering.
 848 In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024a.
 849
- 850 Zhao Yang, Jiakuan Liu, Yucheng Han, Xin Chen, Zebiao Huang, Bin Fu, and Gang Yu. AppAgent:
 851 Multimodal agents as smartphone users. *arXiv preprint arXiv:2312.13771*, 2023.
- 852 Zhiyu Yang, Zihan Zhou, Shuo Wang, Xin Cong, Xu Han, Yukun Yan, Zhenghao Liu, Zhixing
 853 Tan, Pengyuan Liu, Dong Yu, Zhiyuan Liu, Xiaodong Shi, and Maosong Sun. MatPlotAgent:
 854 Method and evaluation for LLM-based agentic scientific data visualization. In Lun-Wei Ku, Andre
 855 Martins, and Vivek Srikumar (eds.), *Findings of the Association for Computational Linguistics:*
 856 *ACL 2024*, pp. 11789–11804, Bangkok, Thailand, August 2024b. Association for Computational
 857 Linguistics. doi: 10.18653/v1/2024.findings-acl.701. URL [https://aclanthology.org](https://aclanthology.org/2024.findings-acl.701)
 858 [/2024.findings-acl.701](https://aclanthology.org/2024.findings-acl.701).
- 859 Fengji Zhang, Bei Chen, Yue Zhang, Jacky Keung, Jin Liu, Daoguang Zan, Yi Mao, Jian-Guang
 860 Lou, and Weizhu Chen. RepoCoder: Repository-level code completion through iterative retrieval
 861 and generation. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023*
 862 *Conference on Empirical Methods in Natural Language Processing*, pp. 2471–2484, Singapore,
 863 December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.151. URL <https://aclanthology.org/2023.emnlp-main.151>.

864 Ruochen Zhao, Wenxuan Zhang, Yew Ken Chia, Deli Zhao, and Lidong Bing. Auto arena of LLMs:
865 Automating LLM evaluations with agent peer-battles and committee discussions. *arXiv preprint*
866 *arXiv:2405.20267*, 2024.
867

868 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,
869 Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging LLM-as-a-Judge with MT-bench and
870 chatbot arena. *Advances in Neural Information Processing Systems*, 36, 2024.

871 Wangchunshu Zhou, Yuchen Eleanor Jiang, Long Li, Jialong Wu, Tiannan Wang, Shi Qiu, Jin-
872 tian Zhang, Jing Chen, Ruipu Wu, Shuai Wang, et al. Agents: An open-source framework for
873 autonomous language agents. *arXiv preprint arXiv:2309.07870*, 2023.
874

875 Wangchunshu Zhou, Yixin Ou, Shengwei Ding, Long Li, Jialong Wu, Tiannan Wang, Jiamin Chen,
876 Shuai Wang, Xiaohua Xu, Ningyu Zhang, et al. Symbolic learning enables self-evolving agents.
877 *arXiv preprint arXiv:2406.18532*, 2024.

878 Mingchen Zhuge, Wenyi Wang, Louis Kirsch, Francesco Faccio, Dmitrii Khizbullin, and Jürgen
879 Schmidhuber. GPTSwarm: Language agents as optimizable graphs. In *Forty-first International*
880 *Conference on Machine Learning*.

881 Mingchen Zhuge, Haozhe Liu, Francesco Faccio, Dylan R Ashley, Róbert Csordás, Anand
882 Gopalakrishnan, Abdullah Hamdi, Hasan Abed Al Kader Hammoud, Vincent Herrmann,
883 Kazuki Irie, et al. Mindstorms in natural language-based societies of mind. *arXiv preprint*
884 *arXiv:2305.17066*, 2023.
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

Appendix

Table of Contents

918		
919		
920		
921		
922		
923		
924		
925	A	Logic Flow of this Paper 19
926		
927	B	Experiment Designs 20
928	B.1	Summary of Experiments 20
929	B.2	Judges and Subjects of Evaluation 20
930		
931	C	Agent-as-a-Judge Pipeline 21
932		
933	D	Extend Related Work 22
934		
935	E	The Procedures of Creating DevAI Dataset 24
936	E.1	Manually draft user queries 24
937	E.2	Set Judging Criteria 24
938	E.3	Building Dependency Among Requirements 24
939	E.4	Refine the dataset 24
940	E.5	Analyse the dataset 24
941	E.6	Auxiliary Information 25
942	E.7	A Json Format of Our Sample 25
943	F	User experiences of code-generation agentic systems 28
944		
945	G	More DevAI dataset samples 29
946		
947	H	Human Evaluation Procedure 31
948		
949	I	Suggest Constraints 33
950		
951	J	Collected Trajectories 34
952	J.1	Schema 34
953	J.2	A Sample of Trajectory 36
954	K	Ablations of Agent-as-a-Judge 40
955	K.1	Component Abalations 40
956	K.2	Search Algorithms in Search Module 40
957	K.3	Search Algorithms in Retrieve Module 40
958		
959	L	Prompt Demos of Agent-as-a-Judge 42
960	L.1	System Prompt for Agent-as-a-Judge 42
961	L.2	System Prompt for Locate Module 43
962	L.3	System Prompt for Retrieve Module 43
963	L.4	Prompt for Ask Module (for requirement check) 44
964	L.5	Prompt for Locate Module 44
965	M	Judge Evidences Collected from Agent-as-a-Judge 45
966		
967	N	Analysis of failure cases 46
968		
969	O	Sensitivity w.r.t the choice of the backend LLM 47
970		
971	P	Additional Human Evaluation Details 47

A LOGIC FLOW OF THIS PAPER

Paper: *Agent-as-a-Judge: Evaluating Agents with Agents***Key Logic**■ **Step 1: Concept Proposal**

Description: We propose the **Agent-as-a-Judge** concept, an extension of the LLM-as-a-Judge framework, aimed at evaluating agentic systems using other agentic systems.

■ **Step 2: Dataset Creation**

Description: To address the lack of suitable datasets for evaluating agentic systems in automated AI development, we introduce **DevAI**, a new dataset consisting of 55 realistic AI code generation tasks. This also serves as a testbed for the Agent-as-a-Judge proof-of-concept.

■ **Step 3: Baseline Evaluation of Developer Agents (Experiment Level 1)**

Description: In the first level of experiments, we select three popular open-source developer agents: **MetaGPT**, **GPT-Pilot**, and **OpenDevin**. These agents are evaluated on the **DevAI** tasks to establish performance baselines.

■ **Step 4: Conducting Human-as-a-Judge Evaluation**

Description: We conduct a **Human-as-a-Judge** experiment, where three human experts assess the performance of the developer agents on the same DevAI tasks.

■ **Step 5: Human-as-a-Judge Analysis (Experiment Level 2)**

Description: In the second level of experiments, we statistically analyze the results of Human-as-a-Judge evaluations, focusing on the costs of human labor and potential biases, highlighting the challenges of relying on human evaluation for complex tasks.

■ **Step 6: Agent-as-a-Judge Implementation**

Description: We design and implement the **Agent-as-a-Judge** proof-of-concept to evaluate code generation on the **DevAI** dataset. This system incorporates modules such as graph, search, read, and ask, providing multi-dimensional evaluation metrics.

■ **Step 7: Comparing Judgment Systems (Experiment Level 3)**

Description: In the third level of experiments, we compare three judgment systems: **Agent-as-a-Judge**, **LLM-as-a-Judge**, and **Human-as-a-Judge**, all applied to the same DevAI tasks. Our results show that **Agent-as-a-Judge** performs comparably to human evaluators and surpasses **LLM-as-a-Judge** in more complex reasoning and evaluation tasks.

Future Directions■ **Direction 1: Enhancing Agent-as-a-Judge**

Description: Future work should focus on improving the performance of Agent-as-a-Judge, especially in more complex and diverse environments, to handle increasingly sophisticated decision-making tasks.

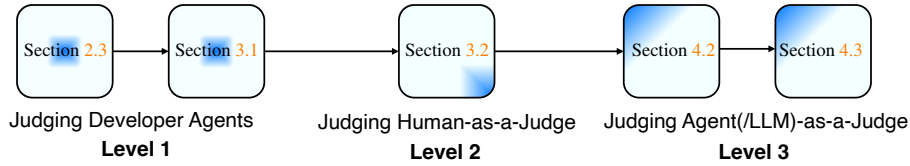
■ **Direction 2: Intermediate Feedback for Self-Improvement**

Description: A promising extension of this work is enabling **Agent-as-a-Judge** to provide intermediate feedback, helping developer agents iteratively improve and self-optimize their decision-making processes.

Figure 8: Logical Flow of the Agent-as-a-Judge Paper.

B EXPERIMENT DESIGNS

This section outlines the experimental designs aimed at evaluating developer agents’ performance, analyzing human evaluations, and comparing AI-based judging systems. The experiments are structured across three levels, as illustrated below.



B.1 SUMMARY OF EXPERIMENTS

The experiments are categorized into three levels as follows:

Level 1: Human evaluation of developer agents

🍃 *Experiment 1a:* Basic performance statistics for developer agents (Section 2.3)

🍃 *Experiment 1b:* Human evaluations of developer agents (Section 3.1)

Level 2: Error and bias analysis of human evaluations

🍃 *Experiment 2a:* Error analysis of human evaluations (Section 3.2)

Level 3: AI-based judging systems

🍃 *Experiment 3a:* AI judge baselines (Section 4.2)

🍃 *Experiment 3b:* Ablation studies for Agent-as-a-Judge (Section 4.3)

B.2 JUDGES AND SUBJECTS OF EVALUATION

The following table summarizes the judge and the subject being evaluated in each experiment:

Experiment	Who is the Judge?	Who is being Judged?
Section 2.3	Human	Developer Agents
Section 3.1	Human	Developer Agents
Section 3.2	Human	Human
Section 4.2	(1) LLM-as-a-Judge (2) Agent-as-a-Judge (3) Human (4) Human	(1) Developer Agents (2) Developer Agents (3) LLM-as-a-Judge (4) Agent-as-a-Judge
Section 4.3	Human	Agent-as-a-Judge

C AGENT-AS-A-JUDGE PIPELINE



Figure 9: The pipelines of developer agents and judge agent.

D EXTEND RELATED WORK

Our main paper includes mostly related works of **AI developers**, **Benchmarks for AI developments**, and **AI judges**. However, the following works contribute significantly to the community and also relate to this work. We record this work as additional related work³.

LLM-based Autonomous Agents Recent developments in LLM-based agents have expanded their capabilities beyond simple task execution to more autonomous problem-solving and decision-making. AutoGPT (Gravitas, 2023) and LangChain (Chase, 2022) provide frameworks for single-agent systems that leverage external tools for more complex tasks. Similarly, research such as MetaGPT (Hong et al., 2024b), AutoGen (Wu et al., 2023) and CAMEL (Li et al., 2023) focus on role-based multi-agent communication, improving collaboration among agents. However, the challenge of maintaining coherence in agents’ dialogue and preventing hallucination remains prominent (Du et al., 2024; Zhou et al., 2023). Most recently, using graphs to build agents has gained prominence. Earlier work like GPTSwarm (Zhuge et al.) and LangGraph (LangChain-AI, 2024) proposed using nodes to represent operations and edges to represent the connections between them. In GPTSwarm, multiple agents represented as subgraphs in a graph are connected by optimizable edges, and reinforcement learning is employed to optimize the edges. Following this approach, several agent frameworks have incorporated graphs into their designs (Hong et al., 2024a; Zhou et al., 2024; Qian et al., 2024b). Additionally, various optimization methods have been developed to enhance agent performance further (Wu et al., 2024; Song et al., 2024; Hu et al., 2024). In practical applications, many studies focus on understanding and interacting with GUIs (Wang et al., 2024a; Chen et al., 2024b; Yang et al., 2023; Xu et al., 2024; Tan et al.). For code generation agents (Jin et al., 2024), current research mainly emphasizes automated repair (Yang et al., 2024a; Phan et al., 2024; Tao et al., 2024), computational modular design (Khatab et al., 2024; Cheng et al.), and automated development (Tufano et al., 2024; Huang et al., 2023). Among these, open-sourced frameworks like OpenHands (Wang et al., 2024b) have gained popularity due to their strong user experience. Moreover, scientific discovery (Jansen et al., 2024; Lu et al., 2024a) and ML agents (Yang et al., 2024b) are also receiving increased attention.

LLM-as-a-Judge In the domain of AI evaluation and judgment, frameworks (Zheng et al., 2024; Fu et al., 2024; Chen et al.) have pioneered the use of LLMs to assess conversational agents, demonstrating how LLMs can evaluate dialogue quality and consistency. Expanding beyond dialogue, LLMs like CodeR (Chen et al., 2024a) and MASAI (Arora et al., 2024) apply similar judging principles to the code validation process, where AI systems autonomously evaluate and verify computer programs. Our work builds on these advancements by exploring how LLMs can perform more nuanced judgment tasks, further investigating their potential in decision-making across various domains. Recent research also focuses on judging LLM-as-a-Judges (Chen et al., 2024c; Bavaresco et al., 2024; Thakur et al., 2024; Dong et al., 2024; Shi et al., 2024; Raina et al., 2024).

Coding Benchmarks Recent advances in code generation have led to the innovation of various benchmarks to evaluate model performance (Liu et al., 2024). Early benchmarks, such as MBPP (Austin et al., 2021), HumanEval (Chen et al., 2021), and MultiPL-E (Cassano et al., 2023), focus primarily on generating simple functions. While these benchmarks are useful for evaluating the correctness of generated code, they are limited in complexity and do not fully represent the challenges encountered in real-world software development.

As the field progressed, newer benchmarks began to focus on more complex and realistic tasks. APPS (Hendrycks et al.), CodeContests (Li et al., 2022), and LiveCodeBench (Jain et al., 2024) moved toward competitive programming challenges that involve advanced algorithms and data structures. These tasks are more representative of problems encountered in coding competitions and help push models toward more sophisticated problem-solving. DS-1000 (Lai et al., 2023) was introduced to assess the skills of models with data science libraries, evaluating their ability to use APIs and execute complex data analysis workflows. Meanwhile, AgentBench (Liu et al., b) focuses on

³Additionally, we were pleased to find that a recent industry blog (<https://www.cognition.ai/blog/evaluating-coding-agents>), published two weeks before our submission, shares very similar ideas and provides further evidence that the Agent-as-a-Judge could have practical applications in agent systems.

1188 testing reasoning and decision-making abilities in interactive environments, highlighting differences
1189 in performance between commercial and open-source models.

1190 To address real-world programming needs beyond code generation, specialized benchmarks have
1191 been created to evaluate tasks such as debugging, refactoring, and code navigation. CANITE-
1192 DIT (Cassano et al., 2024), DebugBench (Tian et al., 2024), and FixEval (Haque, 2023) evalu-
1193 ate the ability of a model to edit and improve existing code. Additionally, benchmarks such as
1194 SWE-Bench (Jimenez et al.) focus on resolving issues in GitHub repositories, simulating practical
1195 software development scenarios. Finally, benchmarks such as RepoBench (Liu et al., a) and RepoE-
1196 val (Zhang et al., 2023) delve into the evaluation of models in large-scale, multifile codebases. These
1197 benchmarks measure the ability of language models to understand the structure of repositories and
1198 solve problems within more complex, collaborative environments.

1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241

E THE PROCEDURES OF CREATING DEVAI DATASET

E.1 MANUALLY DRAFT USER QUERIES

Given the execution cost of the developer agents, we collect small-scale AI tasks to ensure the practical applicability of our benchmark. Since these tasks are small-scale and well studied, which are easy to overfit in terms of task performance metrics, unlike previous benchmarks (e.g., (Huang et al., 2024)), we do not evaluate task performance as the development performance measure. Instead, we prioritize the step-by-step task-solving ability, which is essential for real-world development. Our queries are specifically designed to require the development agents to understand user intentions, solve the task in multiple steps, and adapt to unexpected step outcomes. This approach also makes our benchmark user-friendly, transparent, and better reflects real-world deployment situations. To enable effective evaluation, our queries present a specific development file structure for the developer agents to follow. To ensure that the developer agents save the files to be evaluated in the workspace, we develop constraint prompts added to the query to form an extended query. The constraint prompts guide the developer agents to save and execute the source codes, which are in line with the real-world development standard. See Appendix I for our constraint prompts.

E.2 SET JUDGING CRITERIA

To make the evaluation of the developer agent precise, we assign to each task query a list of requirements as task milestones. The requirements are chosen so that satisfying all the requirements is a necessary condition to consider the task to be solved successively. Since our tasks are AI-centric, our queries target essential elements of AI development, including data processing, AI method, presentation of evaluation metrics, visualization, and human-computer interaction, covering the key areas that matter most in real-world scenarios. On the other hand, breaking down tasks into individual requirements also reflects the multi-step nature of code development. Importantly, to avoid ambiguity, we set the requirements to be explicit, binary, and straightforward to evaluate. To include other human predispositions, we include a list of preferences per task that covers subjective, ambiguous, or non-explicitly stated characteristics.

E.3 BUILDING DEPENDENCY AMONG REQUIREMENTS

To enhance the realism of our benchmark, we analyzed the dependencies among requirements. Over the past decades, methodologies such as the KDD Process (Fayyad et al., 1996) and CRISP-DM (Wirth & Hipp, 2000) have guided ML/AI development, establishing foundational frameworks that have been further refined with the rise of AutoML (He et al., 2021).

Inspired by these methodologies, we identified a seven-step process for automated AI development tasks. This process includes critical stages such as data preprocessing, feature engineering, model selection, and hyperparameter tuning, along with essential post-development activities like metrics recording, report generation, and the development of interactive user applications. This structured approach allows us to evaluate the agent’s ability to manage task dependencies and effectively navigate complex, real-world scenarios.

E.4 REFINE THE DATASET

Manual refinements were necessary to ensure the accuracy and clarity of DevAI. We perform two rounds of comprehensive review and edits on DevAI, each round being done by a different participant. During these manual refinements, we focus on the logical consistency of our queries and requirements, the correctness and ambiguity of the language, and the applicability of the content to the task domain. We identified a moderate number of errors in our dataset during the review.

E.5 ANALYSE THE DATASET

We categorized each requirement based on its focus, whether it was related to the data set, the machine learning method, visualization, metrics, HCI, or data processing. While this approach provides a useful framework, it is important to recognize that these categories might overlap or miss certain nuances. Similarly, preferences were classified by how strongly they appeared in the query

text, ranging from those inferred by common sense to those explicitly stated. Although this helps to organize preferences, it is worth noting that such classifications can be subjective and may not fully capture the importance of each preference in practical applications. By structuring the requirements and preferences this way, we aim to enhance the evaluation process, though flexibility and context awareness remain crucial for truly robust assessments.

E.6 AUXILIARY INFORMATION

Some of the tasks require the download of a Kaggle data set, where a Kaggle credential is needed. Our constraint **I** requires an “is_kaggle_api_needed” tag to determine whether the credential is needed to be included in the extended query. We further mark each task with some tags describing the AI sub-fields related to the task, including computer vision, supervised learning, reinforcement learning, natural language processing, etc., as well as two “is_training_needed” and “is_web_navigation_needed” tags as auxiliary information. We also categorize each requirement into one of the following: (1) dataset or environment, (2) data preprocessing and postprocessing, (3) machine learning method, (4) save trained model, (4) performance metrics, (5) human computer interaction, (6) visualization, and (7) other, reflecting the nature of the requirement.

E.7 A JSON FORMAT OF OUR SAMPLE

Here, we provide a sample of the DevAI with its json format. We also provide more samples in Appendix **G**.

```
{
  "name": "25_Speech_Emotion_Recognition_CNN_LSTM_RAVDESS_DL",
  "query": "I am seeking a speech emotion recognition project using a
  CNN-LSTM model with the RAVDESS dataset, which should be
  downloaded from Kaggle or [this Hugging Face link](https://
  huggingface.co/datasets/xbgoose/ravdess). The project should load
  the dataset and perform robust audio preprocessing (noise
  removal and normalization) and MFCC feature extraction,
  implemented in 'src/data_loader.py'. The CNN-LSTM model should be
  implemented in 'src/model.py'. Recognition accuracy should be
  saved in 'results/metrics/recognition_accuracy.txt', and a
  confusion matrix should be generated and saved as 'results/
  figures/confusion_matrix.png'. Additionally, a user-friendly
  local API should be created using Flask to allow users to upload
  audio files and receive emotion recognition results, with the
  implementation included in 'src/hci.py'.",
  "tags": [
    "Audio Processing",
    "Classification"
  ],
  "requirements": [
    {
      "requirement_id": 0,
      "prerequisites": [],
      "criteria": "The \"RAVDESS\" dataset is loaded in 'src/
      data_loader.py', which is downloaded from Kaggle or [this
      Hugging Face link](https://huggingface.co/datasets/
      xbgoose/ravdess).",
      "category": "Dataset or Environment",
      "satisfied": null
    },
    {
      "requirement_id": 1,
      "prerequisites": [
        0
      ],
      "criteria": "Audio preprocessing, including noise removal and
      normalization, is implemented in 'src/data_loader.py'.",
      "category": "Data preprocessing and postprocessing",
```

```

1350     "satisfied": null
1351   },
1352   {
1353     "requirement_id": 2,
1354     "prerequisites": [
1355       0,
1356       1
1357     ],
1358     "criteria": "MFCC feature extraction is implemented in `src/
1359     data_loader.py`.",
1360     "category": "Data preprocessing and postprocessing",
1361     "satisfied": null
1362   },
1363   {
1364     "requirement_id": 3,
1365     "prerequisites": [],
1366     "criteria": "The `CNN-LSTM` model is implemented in `src/
1367     model.py`.",
1368     "category": "Machine Learning Method",
1369     "satisfied": null
1370   },
1371   {
1372     "requirement_id": 4,
1373     "prerequisites": [
1374       2,
1375       3
1376     ],
1377     "criteria": "Recognition accuracy is saved in `results/
1378     metrics/recognition_accuracy.txt`.",
1379     "category": "Performance Metrics",
1380     "satisfied": null
1381   },
1382   {
1383     "requirement_id": 5,
1384     "prerequisites": [
1385       2,
1386       3,
1387       4
1388     ],
1389     "criteria": "The confusion matrix is generated and saved as `
1390     results/figures/confusion_matrix.png`.",
1391     "category": "Visualization",
1392     "satisfied": null
1393   },
1394   {
1395     "requirement_id": 6,
1396     "prerequisites": [
1397       2,
1398       3
1399     ],
1400     "criteria": "A local API is created using `Flask` to allow
1401     users to upload audio files and receive emotion
1402     recognition results. The implementation should be
1403     included in `src/hci.py`.",
1404     "category": "Human Computer Interaction",
1405     "satisfied": null
1406   }
1407 ],
1408 "preferences": [
1409   {
1410     "preference_id": 0,
1411     "criteria": "The audio preprocessing step should be robust,
1412     effectively reducing noise while preserving the integrity
1413     of the speech signals.",
1414     "satisfied": null

```

```
1404     },
1405     {
1406         "preference_id": 1,
1407         "criteria": "The local API should be user-friendly, with
1408             clear instructions for uploading files and interpreting
1409             results.",
1410         "satisfied": null
1411     }
1412 ],
1413 "is_kaggle_api_needed": true,
1414 "is_training_needed": true,
1415 "is_web_navigation_needed": true
1416 }
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457
```

1458 F USER EXPERIENCES OF CODE-GENERATION AGENTIC SYSTEMS
1459

1460 OpenHands (Wang et al., 2024b) offers the most refined user experience, leveraging its highly inter-
1461 active frontend to enable seamless user interaction and task execution. This interface allows users to
1462 engage directly with the system, resulting in a smoother and more intuitive workflow, which drives
1463 operational efficiency.

1464 In contrast, MetaGPT (Hong et al., 2024b) excels in task decomposition through its use of Directed
1465 Acyclic Graphs (DAGs), a well-structured and scalable approach aligned with industry best practices
1466 in system modularization. This enhances its appeal for users focused on task clarity and modular
1467 breakdowns. However, in practical deployments, MetaGPT tends to be less aggressive in file man-
1468 agement and preservation, potentially due to its core positioning as a data analysis tool, which does
1469 not prioritize persistent state management. Similarly, OpenDevin demonstrates a notable overconfi-
1470 dence in its code generation, frequently skipping the critical step of post-generation code execution,
1471 requiring users to intervene manually.

1472 GPT-Pilot (Pythagora.io, 2023), praised for its detailed task delegation via over 20 specialized
1473 agents, suffers from reduced interactivity due to an overly granular division of responsibilities, re-
1474 sulting in a more fragmented user experience. These qualitative insights, although not fully captured
1475 by quantitative metrics, were evident through the DevAI dataset, providing key areas for improve-
1476 ment in user engagement and operational fluidity in future releases of these frameworks.

1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511

G MORE DEVAI DATASET SAMPLES

Task 13: Style Transfer with Perceptual Loss in PyTorch

Query

Please create a PyTorch Perceptual Loss project for image style transfer (refer to this paper: *Perceptual Losses for Real-Time Style Transfer*). You can build the Perceptual Loss Network using VGG16 in `src/model.py`. The project should combine content and style images, allow smooth adjustment of style intensity by tuning the weights of style loss and content loss, and save the stylized images in `results/figures/`. Additionally, log the processing time to `results/processing_time.txt`, and save the intermediate results of the style transfer process to `results/figures/intermediate_results.png`. For testing, input a famous content image (Mona Lisa) from [this link](#) and a famous style image (The Starry Night) from [this link](#), and generate a style-transferred image. Save the content, style, and style-transferred images to `data/content.jpg`, `data/style.jpg`, and `results/figures/`, respectively. The project should efficiently handle high-resolution images without excessive processing time.

Requirements

- **R0**
Criteria: A famous content image is inputted for testing, downloaded from [this link](#) and saved to `data/content.jpg`. **Dependencies** $\rightarrow \{\}$
- **R1**
Criteria: A famous style image is inputted for testing, downloaded from [this link](#) and saved in `data/style.jpg`. **Dependencies** $\rightarrow \{\}$
- **R2**
Criteria: The Perceptual Loss model is implemented in PyTorch and loaded in `src/model.py`. **Dependencies** $\rightarrow \{\}$
- **R3**
Criteria: Stylized images are saved to the specified folder `results/figures/`. **Dependencies** $\rightarrow \{\mathbf{R0}, \mathbf{R1}, \mathbf{R2}\}$
- **R4**
Criteria: Style intensity is adjusted by tuning the weights of style loss and content loss in `src/model.py`. **Dependencies** $\rightarrow \{\mathbf{R0}, \mathbf{R1}, \mathbf{R2}\}$
- **R5**
Criteria: Processing time is recorded and saved as `results/processing_time.txt`. **Dependencies** $\rightarrow \{\mathbf{R0}, \mathbf{R1}, \mathbf{R2}, \mathbf{R3}, \mathbf{R4}\}$
- **R6**
Criteria: Intermediate results of style transfer are saved as `results/figures/intermediate_results.png`. **Dependencies** $\rightarrow \{\mathbf{R0}, \mathbf{R1}, \mathbf{R2}, \mathbf{R3}, \mathbf{R4}\}$

Preferences (Optional)

- **P0**
Criteria: The style transfer process should allow for smooth adjustment of style intensity, making the stylized image visually appealing.
- **P1**
Criteria: The project should handle high-resolution images efficiently without excessive processing time.

Figure 10: An Example Task in DevAI: Task 13.

1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619

Task 19: Time Series Forecasting with Seq2Seq LSTM on Rossmann Store Sales

Query

Develop a sales forecasting system using a sequence-to-sequence model based on LSTM with the Rossmann Store Sales dataset, downloading it from Kaggle [here](#) and loading it in `src/data_loader.py`. Split the data into training and testing sets and save them in `src/data_loader.py`. Apply a sequence-to-sequence model based on LSTM and save the trained model under the `models/saved_models/` directory. Save the forecast results as `results/figures/forecast_results.png`. Save a comparison plot between the predicted and actual values to `results/figures/comparison_plot.png`. Generate an HTML report that includes the prediction results and comparison plots, with some interactive elements for exploring different forecast horizons, and save it as `results/report.html`. Ensure the model is tuned to capture seasonal trends in the sales data.

Requirements

- **R0**
Criteria: *The Rossmann Store Sales dataset is used, potentially downloaded from ([this link](#)) and loaded in `src/data_loader.py`. **Dependencies** → {}*
- **R1**
Criteria: *The data is split into training and testing sets and implemented in `src/data_loader.py`. **Dependencies** → {R0}*
- **R2**
Criteria: *A sequence-to-sequence model based on LSTM is used. The trained model should be saved under `models/saved_models/`. **Dependencies** → {R1}*
- **R3**
Criteria: *The forecast results are plotted and saved as `results/figures/forecast_results.png`. **Dependencies** → {R1, R2}*
- **R4**
Criteria: *A comparison plot of predicted vs. actual values is saved as `results/figures/comparison_plot.png`. **Dependencies** → {R1, R2, R3}*
- **R5**
Criteria: *An HTML report containing forecast results and comparison plots is generated and saved as `results/report.html`. **Dependencies** → {R1, R2, R3, R4}*
- **R6**
Criteria: *The HTML report should include interactive elements that allow users to explore different forecast horizons. **Dependencies** → {R5}*

Preferences (Optional)

- **P0**
Criteria: *The model should be tuned to capture seasonal trends in the sales data for more accurate forecasting.*

Figure 11: An Example Task in DevAI: Task 19.

1620 H HUMAN EVALUATION PROCEDURE
1621

1622 We recruited three AI experts from the authors to perform human evaluation on the output of agentic
1623 code generation systems. There we present the evaluation details.
1624

1625 **First round** For the first round of evaluations, our three evaluators reported spending 16.5, 19.5,
1626 and 22.0 hours, respectively. To capture the bias that a human evaluator will have, the instruc-
1627 tions given to our experts were minimal, with them only receiving a scorecard to complete for each
1628 agentic system and each task. Results that all evaluators agree on are considered trustworthy. The
1629 assumption here is that it is unlikely that all three evaluators make a mistake or have an effective
1630 bias in the same judgment. The self-reported post-hoc evaluation criteria are shown in Figure 12.
1631

1632 **Second round** In the second round, the evaluators present and discuss their reasons for disagree-
1633 ing with judges. In doing so, human errors are likely corrected by their peers. Discussion among
1634 evaluators also helps reduce human bias by examining each other’s thought processes thoroughly.
1635 Furthermore, the consensed results are considered trustworthy given the assumption that it is un-
1636 likely that all three evaluators are convinced by the same mistake or the same cognitive bias. The
1637 three evaluators took 9.5 hours together for this second round of evaluation.
1638

1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673

Self-reported Post-hoc Evaluation Criteria after round one

1674
1675
1676
1677
1678
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727

Evaluator 231a

EXECUTED SUCCESSFULLY: Yes [] / No []

1. Must be checked based on the overall completeness of the task, based on looking at the code, the artifacts, and the trajectory.
2. The training has finished, the model snapshot and the metrics breakdown have been saved, and at least one artifact of required analytics has been produced - mark as successful. If some analytic artifacts are missing but not all, mark as successful.
3. No need to run the code.
4. If training was finished but on fake data, mark as successful.

Requirements:

Marking a requirement as satisfied must be made for this specific requirement disregarding the dependency list. If a file (code, image, snapshot) is there but is empty or without any meaningful content - mark as No.

1. Code:

The functionality must be in a file with the requested path.
The real data is replaced by simple synthetic - not satisfied.

2. Visualization/Reports:

The contents must be there and make sense even if not perfect from the ML/DS point of view.

3. Snapshots:

If a binary snapshot is not empty, mark as Yes.

Evaluator 38bb

EXECUTED SUCCESSFULLY: Yes [] / No []

An output is marked yes if none of the following is satisfied.

1. The time spent is close to the time limit.
2. The last environmental message includes an error.
3. The last thought indicates that the task is completed.
4. The last step of the trajectory is incomplete.

Requirements:

If a required is unsatisfied only because of an unsatisfaction of a previously marked unsatisfied requirement, then judge it based on the assumption that a minimum implementation satisfies the previous requirement exists.

1. Code:

Mark yes if the code executes and does the required function. If no entrypoint is given, the evaluator will set an entrypoint. If the code is not executable due to previous unsatisfied requirements, then it is judged based on eye-checking.

2. Visualization/Reports:

Mark yes if the visualization or report exists in a right path and the content aligns the requirement.

3. Snapshots:

Mark yes if the snapshot exists in a right path and is not empty.

Evaluator cn9o

I evaluated everything based on whether the requirement was satisfied verbatim, using my own professional judgement when there was ambiguity. If there was a data folder and a look at the trajectory indicated it was real, I allowed it to be loaded directly. If there were results and no clear source code making them, I considered the task incomplete. I didn't consider a dummy data path to be correct (GPT-Pilot used a lot of these). To determine if something was executed correctly, I skimmed the end of the logs (e.g., trajectory) for any obvious signs of an error. If there wasn't any, I said it was correctly executed. I ignored prerequisites while evaluating and did not execute anything, instead just eyeballing the code for correctness. I was lenient in what I considered to be sufficient in terms of the more vague requirements (e.g., if the preprocessing had to include scaling and rotation, just those two would be sufficient to consider it done).

Figure 12: Each evaluator was given a full description of each task and the associated requirements and preferences in markdown format. They were then asked, for each workspace and trajectory generated by each of the agents on each of the task, whether (1) the agent successfully finished its execution cycle and (2) which of the requirements were satisfied. After the evaluation was complete, the evaluators were asked to self-report the nuances of their evaluation.

I SUGGEST CONSTRAINTS

Below is a sample of constraints in JSON format that describes task-specific guidelines:

```
{
  "generic": "This is a task that requires you to write, execute, and
    save source code. You have a hard time limit of 30 minutes to
    produce your programmatic solution to the given task. This time
    limit includes execution time. The quality of your solution will
    be judged based on what you left in the working folder by the time
    30 minutes expire. Additionally, the hardware you are running on
    is unknown, and the presence of a GPU is not guaranteed.",
  "is_training_needed": "Keep the time limit in mind when setting
    hyperparameters for training.",
  "is_kaggle_api_needed": "You can use the Kaggle API credentials stored
    in `kaggle.json` in your current working directory."
}
```

To address automation and security concerns, we have written code to place the `kaggle.json` file into the current workspace for each baseline during each run. However, dataset users are free to modify the solution to enable the Kaggle API for their specific developer agents.

1782 J COLLECTED TRAJECTORIES

1783

1784 J.1 SCHEMA

1785

1786 Below is the required JSON format for a trajectory in gray-box settings (where the trajectories can
1787 serve as input for LLM-as-a-Judge and Agent-as-a-Judge).

1788

1789

```

1790 {
1791   "type": "array",
1792   "items": {
1793     "type": "object",
1794     "properties": {
1795       "step": {
1796         "type": "integer",
1797         "description": "The step number in the trajectory, 0-based."
1798       },
1799       "user_message": {
1800         "type": ["string", "null"],
1801         "description": "The message from the external user to the agent.
1802           If null, no message was sent."
1803       },
1804       "agent": {
1805         "type": "object",
1806         "properties": {
1807           "thought": {
1808             "type": "string",
1809             "description": "The agent's thought at this step."
1810           },
1811           "action": {
1812             "type": ["string", "null"],
1813             "description": "The agent's action sent to the environment.
1814               If null, the agent did not take any action, for example,
1815               when the agent has finished the task."
1816           },
1817           "agent_name": {
1818             "type": "string",
1819             "description": "The name of the agent that made the action."
1820           }
1821         },
1822         "required": ["thought", "action"],
1823         "description": "Everything related to the agent at this step."
1824       },
1825       "environment": {
1826         "type": ["string", "null"],
1827         "description": "The environment's (shell, python interpreter)
1828           response to the action submitted by the agent. If null, the
1829           environment was not involved in this step."
1830       },
1831       "step_usage": {
1832         "type": "object",
1833         "properties": {
1834           "input_tokens": {
1835             "type": "integer",
1836             "description": "The number of input tokens passed as LLM
1837               context."
1838           },
1839           "output_tokens": {
1840             "type": "integer",
1841             "description": "The number of tokens produced by the LLM."
1842           }
1843         },
1844       "model": {
1845         "type": "string",
1846         "description": "The name of the LLM model used."
1847       }
1848     }
1849   }
1850 }

```

```
1836     "cost": {
1837         "type": "number",
1838         "description": "The cost of the LLM inference, in USD."
1839     },
1840     "llm_inference_time": {
1841         "type": "number",
1842         "description": "The time taken by the LLM to generate the
1843             output tokens, in seconds."
1844     },
1845     "step_execution_time": {
1846         "type": "number",
1847         "description": "The time taken to make an entire step
1848             including LLM inference and environment execution, in
1849             seconds."
1850     }
1851 },
1852 "required": [
1853     "input_tokens",
1854     "output_tokens",
1855     "model",
1856     "cost",
1857     "llm_inference_time",
1858     "step_execution_time"
1859 ]
1860 },
1861 "accumulated_usage": {
1862     "type": "object",
1863     "properties": {
1864         "accumulated_cost": {
1865             "type": "number",
1866             "description": "The total cost of the trajectory up to this
1867                 step, in USD."
1868         },
1869         "accumulated_time": {
1870             "type": "number",
1871             "description": "The total time taken by the agent to complete
1872                 the trajectory up to this step, in seconds."
1873         }
1874     }
1875 },
1876 "required": [
1877     "accumulated_cost",
1878     "accumulated_time"
1879 ]
1880 },
1881 "required": ["step", "user_message", "agent", "environment", "
1882     step_usage", "accumulated_usage"]
1883 }
1884 }
1885 }
1886 }
1887 }
1888 }
1889 }
```

1890 J.2 A SAMPLE OF TRAJECTORY
1891

1892 This is sample trajectories (we select 4 steps of it) in OpenHands (Wang et al., 2024b). All steps are
1893 recorded to enable replay of the trajectory.
1894

```
1895 {
1896     "step": 0,
1897     "user_message": "Help me develop an object detection system using
1898     the YOLOv3 model and the COCO dataset. Download the dataset
1899     and preprocess the images by resizing and normalization in `
1900     src/data_loader.py`. Implement the YOLOv3 model and use Non-
1901     Maximum Suppression (NMS) to refine the results in `src/model
1902     .py`. Save the detected objects to `results/figures/`, and
1903     create an interactive Streamlit web page in `src/app.py` to
1904     display the detection results. Finally, evaluate the model's
1905     performance, including metrics such as mAP and inference time
1906     , and save the evaluation results to `results/metrics/
1907     model_performance.txt`. The system should properly manage the
1908     launch and termination of the Streamlit application to
1909     prevent unnecessary resource usage.\n\nThis is a task that
1910     requires you to write, execute, and save source code. You
1911     have a hard time limit of 30 minutes to produce your
1912     programmatic solution to the given task. This time limit
1913     includes execution time. The quality of your solution will be
1914     judged based on what you left in the working folder by the
1915     time 30 minutes expire. Additionally, the hardware you are
1916     running on is unknown, and the presence of a GPU is not
1917     guaranteed.\n\nKeep the time limit in mind when setting
1918     hyperparameters for training.\n",
1919     "agent": {
1920         "thought": "Let's break down the task into manageable steps
1921         and start implementing them one by one.\n\n### Steps:\n1.
1922         **Set up the directory structure**:\n - Create
1923         necessary directories: `src/` and `results/figures/`. \n2.
1924         **Download and preprocess the COCO dataset**:\n -
1925         Implement `src/data_loader.py` to download, resize, and
1926         normalize images.\n3. **Implement the YOLOv3 model**:\n
1927         - Implement `src/model.py` to define the YOLOv3 model
1928         and apply Non-Maximum Suppression (NMS).\n4. **Create an
1929         interactive Streamlit web page**:\n - Implement `src/
1930         app.py` to display detection results.\n5. **Evaluate the
1931         model's performance**:\n - Save evaluation metrics to `
1932         results/metrics/model_performance.txt`. \n\nLet's start by
1933         setting up the directory structure and creating the
1934         necessary files.\n\n### Step 1: Set up the directory
1935         structure",
1936         "action": "Running command: mkdir -p src results/figures
1937         results/metrics\ntouch src/data_loader.py src/model.py
1938         src/app.py"
1939     },
1940     "environment": "\r\n[Python Interpreter: /openhands/poetry/
1941     openhands-504_aChf-py3.11/bin/python]\nopenhands@b34322eae8da
1942     :/workspace $ touch src/data_loader.py src/model.py src/app.
1943     py\r\n\r\n[Python Interpreter: /openhands/poetry/openhands-50
1944     4_aChf-py3.11/bin/python]\nopenhands@b34322eae8da:/workspace
1945     $ ",
1946     "step_usage": {
1947         "input_tokens": 4331,
1948         "output_tokens": 220,
1949         "llm_inference_time": 4.445789337158203,
1950         "model": "gpt-4o-2024-05-13",
1951         "cost": 0.024955,
1952         "step_execution_time": 16.24921226501465
1953     },
1954 }
```

```

1944     "accumulated_usage": {
1945         "accumulated_cost": 0.024955,
1946         "accumulated_time": 16.24922204017639
1947     }
1948 },
1949 {
1950     "step": 1,
1951     "user_message": null,
1952     "agent": {
1953         "thought": "Great! The directory structure and necessary
1954         files are created. Let's proceed with implementing the
1955         data loader to download and preprocess the COCO dataset.\n
1956         \n\n### Step 2: Implement `src/data_loader.py`\nWe will:\n
1957         1. Download the COCO dataset.\n2. Resize and normalize
1958         the images.\n\nLet's start by writing the code for `src/
1959         data_loader.py`.",
1960         "action": "Running Python code interactively: create_file('
1961         src/data_loader.py') "
1962     },
1963     "environment": "FileExistsError
1964     Traceback (most recent call last)\nCell In[1], line 1\n---->
1965     1 create_file('src/data_loader.py')\n\nFile /openhands/code/
1966     openhands/runtime/plugins/agent_skills/file_ops/file_ops.py:2
1967     57, in create_file(filename)\n    251 `\"\"\"Creates and opens
1968     a new file with the given name.\n    252 \n    253 Args:\n
1969     254     filename: str: The name of the file to create.\n
1970     255 `\"\"\"`\n    256 if os.path.exists(filename):\n--> 257
1971     raise FileExistsError(f\"File '{filename}' already
1972     exists.\")\n    259 with open(filename, 'w') as file:\n    26
1973     0 file.write(`\n`)\n\nFileExistsError: File 'src/
1974     data_loader.py' already exists.\n[Jupyter current working
1975     directory: /workspace]\n[Jupyter Python interpreter: /
1976     openhands/poetry/openhands-504_aChf-py3.11/bin/python]",
1977     "step_usage": {
1978         "input_tokens": 4675,
1979         "output_tokens": 84,
1980         "llm_inference_time": 2.136143207550049,
1981         "model": "gpt-4o-2024-05-13",
1982         "cost": 0.024635000000000004,
1983         "step_execution_time": 3.345384359359741
1984     },
1985     "accumulated_usage": {
1986         "accumulated_cost": 0.049590000000000001,
1987         "accumulated_time": 19.594608068466187
1988     }
1989 },
1990 {
1991     "step": 2,
1992     "user_message": null,
1993     "agent": {
1994         "thought": "The file `src/data_loader.py` already exists. Let
1995         's open it and start writing the code to download and
1996         preprocess the COCO dataset.\n\n### Step 2: Implement `
1997         src/data_loader.py`",
1998         "action": "Running Python code interactively: open_file('src/
1999         data_loader.py') "
2000     },
2001     "environment": "[File: /workspace/src/data_loader.py (1 lines
2002     total)]\n(this is the beginning of the file)\n1|\n(this is
2003     the end of the file)\n[Jupyter current working directory: /
2004     workspace]\n[Jupyter Python interpreter: /openhands/poetry/
2005     openhands-504_aChf-py3.11/bin/python]",
2006     "step_usage": {
2007         "input_tokens": 4982,
2008         "output_tokens": 53,

```

```

1998         "llm_inference_time": 2.209756851196289,
1999         "model": "gpt-4o-2024-05-13",
2000         "cost": 0.025705000000000002,
2001         "step_execution_time": 2.318861961364746
2002     },
2003     "accumulated_usage": {
2004         "accumulated_cost": 0.075295000000000001,
2005         "accumulated_time": 21.913471937179565
2006     }
2007     ...
2008     {
2009         "step": 20,
2010         "user_message": null,
2011         "agent": {
2012             "thought": "It seems the indentation issue persists. Let's
2013             correct the indentation manually.\n\n### Corrected Code
2014             for `src/evaluate.py`\n``python\nimport time\nimport
2015             numpy as np\nfrom model import yolo_v3_model,
2016             non_max_suppression\nfrom data_loader import
2017             preprocess_image, download_coco_dataset,
2018             extract_coco_dataset\nimport os\n\ndef evaluate_model():\n
2019             \n     model = yolo_v3_model()\n     download_coco_dataset()\n
2020             \n     extract_coco_dataset()\n\n     image_dir = \"data/
2021             val2017\"\n     image_files = [os.path.join(image_dir, f)
2022             for f in os.listdir(image_dir) if f.endswith('.jpg')]\n\n
2023             total_inference_time = 0\n     num_images = len(
2024             image_files)\n     all_boxes = []\n     all_scores = []\n
2025             all_classes = []\n     for image_file in image_files:\n
2026             \n     image = preprocess_image(image_file)\n
2027             image = np.expand_dims(image, axis=0)\n\n
2028             start_time = time.time()\n     predictions = model.
2029             predict(image)\n     end_time = time.time()\n\n
2030             inference_time = end_time - start_time\n
2031             total_inference_time += inference_time\n\n     boxes,
2032             scores, classes = non_max_suppression(predictions[0])\n
2033             all_boxes.append(boxes)\n     all_scores.append
2034             (scores)\n     all_classes.append(classes)\n\n
2035             mean_inference_time = total_inference_time / num_images\n
2036             \n     # Placeholder for mAP calculation\n     mAP = 0.0 #
2037             This should be replaced with actual mAP calculation\n\n
2038             with open(\"results/metrics/model_performance.txt\", \
2039             \"w\") as f:\n     f.write(f\"Mean Inference Time: {
2040             mean_inference_time:.4f} seconds\\n\\n\")\n     f.write(f
2041             \"mAP: {mAP:.4f}\\n\\n\")\n\nif __name__ == \"__main__\":\n
2042             evaluate_model()\n```\n\nLet's update the file with
2043             the corrected code.",
2044             "action": "Running Python code interactively: EDITED_CODE = \
2045             \"\nimport time\nimport numpy as np\nfrom model import
2046             yolo_v3_model, non_max_suppression\nfrom data_loader
2047             import preprocess_image, download_coco_dataset,
2048             extract_coco_dataset\nimport os\n\ndef evaluate_model():\n
2049             \n     model = yolo_v3_model()\n     download_coco_dataset()\n
2050             \n     extract_coco_dataset()\n\n     image_dir = \"data/
2051             val2017\"\n     image_files = [os.path.join(image_dir, f)

```

```
2052     scores, classes = non_max_suppression(predictions[0])\n2053         all_boxes.append(boxes)\n         all_scores.append\n2054 (scores)\n         all_classes.append(classes)\n\n2055 mean_inference_time = total_inference_time / num_images\n2056 \n         # Placeholder for mAP calculation\n         mAP = 0.0 #\n2057     This should be replaced with actual mAP calculation\n\n2058     with open(\"results/metrics/model_performance.txt\", \"\n2059     \"w\") as f:\n         f.write(f\"Mean Inference Time: {\n2060     mean_inference_time:.4f} seconds\\\\\\\\n\")\n         f.write\n2061     (f\"mAP: {mAP:.4f}\\\\\\\\n\")\n\n         nif __name__ == \"__main__\"\n2062     : \n         evaluate_model()\n\n         \"\"\" \n         nedit_file_by_replace('n\n2063     src/evaluate.py', 'with open(\"results/metrics/\n2064     model_performance.txt\", \"w\") as f:', 'with open(\"n\n2065     results/metrics/model_performance.txt\", \"w\") as f:')\n\n2066 },\n2067 \"environment\": null,\n2068 \"step_usage\": {\n2069     \"input_tokens\": 20132,\n2070     \"output_tokens\": 805,\n2071     \"llm_inference_time\": 15.1565523147583,\n2072     \"model\": \"gpt-4o-2024-05-13\",\n2073     \"cost\": 0.11273500000000002,\n2074     \"step_execution_time\": 15.196227073669434\n2075 },\n2076 \"accumulated_usage\": {\n2077     \"accumulated_cost\": 1.2919750000000003,\n2078     \"accumulated_time\": 176.99772310256958\n2079 }\n2080 }\n2081\n2082\n2083\n2084\n2085\n2086\n2087\n2088\n2089\n2090\n2091\n2092\n2093\n2094\n2095\n2096\n2097\n2098\n2099\n2100\n2101\n2102\n2103\n2104\n2105
```

K ABLATIONS OF AGENT-AS-A-JUDGE

K.1 COMPONENT ABALATIONS

Table 5: **Component Ablation Studies for Agent-as-a-Judge.** We analyze the impact of sequentially adding various components (`ask`, `graph`, `read`, `locate`, `retrieve`, `planning`, `memory`, and `search`) on the performance of Agent-as-a-Judge for judging OpenHands.

Metric	+ask	+graph	+read	+locate	+search	+retrieve	+planning	+memory
Agent-as-a-Judge Performance								
Alignment Rate	65.03%	75.95%	82.24%	90.44%	86.06%	90.16%	88.52%	87.97%

Analysis We designed 8 modular components for the Agent-as-a-Judge system. In the Table 5, components are added progressively from left to right. If the addition of a component led to a significant performance drop, we removed it from further iterations. Our experiments showed that adding the components `ask`, `graph`, `read`, and `locate` resulted in significant performance gains. However, when the `search` component was introduced, there was a noticeable decline in performance.

We hypothesize that the performance drop from `search` is due to its role in retrieving relevant code snippets (top-3) using BM25. The retrieval accuracy of BM25 (Robertson et al., 2009) might not have been high enough, potentially introducing noise. Moreover, as noted in Table 1, the DevAI tasks in our experiments did not generate a large amount of code. In fact, even when all code was fed into an LLM, the total content typically stayed within the maximum context length. Therefore, in simpler workspaces, `search` was less critical. However, we believe this component will become more important as the complexity of the workspace increases, making it more valuable in larger and more complex environments.

We also observed that the introduction of the `planning` mechanism did not bring a noticeable improvement in performance. This may be related to the nature of the Judge - it needs clean factual information. When `planning` is unstable, the evidence collected from different actions can become inconsistent, leading to a decline in performance. Finally, we experimented with a `memory` mechanism. Initially, we hypothesized that since DevAI tasks often involve interconnected requirements, `memory` could help track whether requirements were met. However, in practice, we saw no improvement. We suspect that the interconnected nature of the requirements may have caused biases: specifically, once a prior requirement was fulfilled, it might have overly influenced positive judgments on subsequent requirements, even if they were not fully met.

K.2 SEARCH ALGORITHMS IN SEARCH MODULE

We initially hypothesized that the performance drop was due to the low precision of the `search` component, particularly with BM2.5. To explore this, we replaced BM2.5 with Sentence-BERT (Reimers & Gurevych, 2019) as a more advanced alternative and tested Fuzzy Search (Levenshtein, 1966) as a less precise option. However, neither improved the performance of the Agent-as-a-Judge.

these results suggest that the performance issue is not due to BM2.5’s poor search accuracy. Instead, the workspaces generated in our DevAI tasks are too simple for the `search` component to have a significant impact. In simpler workspaces, direct retrieval and evaluation are sufficient. Even though Sentence-BERT performed better than the other methods, its alignment rate (87.70%) still falls short of the configuration without the `search` component (90.44%). As workspace complexity increases, the `search` component may become more valuable.

Table 6: Comparisons on Search module with different engines.

Search Method	Alignment Rate
BM2.5	86.06%
Sentence-BERT	87.70%
Fuzzy Search	85.52%
<i>without</i> Search Module	90.44%

K.3 SEARCH ALGORITHMS IN RETRIEVE MODULE

In our experiments, we found that accurately locating relevant information within a trajectory is a challenging task. Although the addition of the `retrieve` component (gray-box) did not lead to

2160 a significant improvement in performance in this specific case, its impact has been notable in other
 2161 settings, such as in GPT-Pilot. As shown in Table 3, the integration of `retrieve` in GPT-Pilot
 2162 brought substantial gains.

2163 We conducted an ablation study on GPT-Pilot to optimize
 2164 the retrieval of useful information at each step. Our exper-
 2165 iments revealed that in large trajectories, truncating the
 2166 final sections of the file often results in losing critical in-
 2167 formation, as the latter part of the trajectory typically con-
 2168 tains dense information about the final development state.
 2169 Truncating the beginning of the trajectory proved to be the
 2170 most effective in improving the retrieval efficiency.

2171 For individual steps, truncating the middle section
 2172 worked best. This is because error messages usually ap-
 2173 pear early in the output, while the corresponding file paths
 2174 and specific error locations are found towards the end.

2175 By focusing on these retrieval strategies, we can significantly enhance the performance of the
 2176 `retrieve` component, particularly in complex scenarios like GPT-Pilot.

2177
2178
2179
2180
2181
2182
2183
2184
2185
2186
2187
2188
2189
2190
2191
2192
2193
2194
2195
2196
2197
2198
2199
2200
2201
2202
2203
2204
2205
2206
2207
2208
2209
2210
2211
2212
2213

Table 7: Ablations on `retrieve`.

Method	Alignment Rate
Without <code>retrieve</code>	83.88%
With <code>retrieve</code> (gray-box)	86.61%
Trajectory Truncate (head)	86.61%
Trajectory Truncate (middle)	85.52%
Trajectory Truncate (tail)	82.51%
Step Truncate (head)	86.34%
Step Truncate (middle)	86.61%
Step Truncate (tail)	83.88%

2214 L PROMPT DEMOS OF AGENT-AS-A-JUDGE

2215

2216 Here, we present some prompts used by the Agent-as-a-Judge system. Each of these prompt demos
2217 plays a crucial role in guiding the agent's behavior.

2218

2219 L.1 SYSTEM PROMPT FOR AGENT-AS-A-JUDGE

2220

2221

2222

2223

2224

2225

2226

2227

2228

2229

2230

2231

2232

2233

2234

2235

2236

2237

2238

2239

2240

2241

2242

2243

2244

2245

2246

2247

2248

2249

2250

2251

2252

2253

2254

2255

2256

2257

2258

2259

2260

2261

2262

2263

2264

2265

2266

2267

```
def get_system_prompt(language="English"):  
    if language == "English":  
        return """  
        You are an advanced AI system serving as an impartial judge for  
        intelligent code generation outputs. Your primary role is to  
        rigorously evaluate whether the agent's outputs satisfy the  
        specified requirements by thoroughly analyzing the provided  
        code, data, and other relevant materials.  
  
        You will systematically assess aspects such as datasets, model  
        implementations, training procedures, and any task-specific  
        criteria outlined in the requirements. Your evaluations must  
        be objective, detailed, and based solely on the evidence  
        provided.  
  
        For each requirement, deliver one of the following judgments:  
  
        1. <SATISFIED>: Use this if the agent's output fully meets the  
        requirement. Provide a brief and precise explanation  
        demonstrating how the specific criteria are fulfilled.  
  
        2. <UNSATISFIED>: Use this if the agent's output does not meet  
        the requirement. Provide a concise explanation indicating the  
        deficiencies or omissions.  
  
        Your assessment should reference specific elements such as code  
        snippets, data samples, or output results where appropriate.  
        Ensure that your justifications are clear, precise, and  
        directly related to the criteria.  
  
        Respond with either <SATISFIED> or <UNSATISFIED>, followed by  
        your concise justification.  
        """  
    else:  
        raise NotImplementedError(f"The language '{language}' is not  
        supported.")
```

2268 L.2 SYSTEM PROMPT FOR LOCATE MODULE

```

2269
2270
2271 def get_locate_system_prompt(language="English"):
2272     if language == "English":
2273         return """
2274         You are an advanced AI system specializing in understanding
2275         project structures and determining file locations based on
2276         provided criteria.
2277         Your task is to locate specific files in the workspace based on
2278         the user's criteria and workspace information.ution problems
2279         with the files mentioned in the criteria.
2280         """
2281     else:
2282         raise NotImplementedError(f"The language '{language}' is not
2283         supported.")

```

2284 L.3 SYSTEM PROMPT FOR RETRIEVE MODULE

```

2285
2286
2287 def get_retrieve_system_prompt(language="English"):
2288     if language == "English":
2289         return """
2290         You are an advanced AI system specializing in retrieving
2291         environmental feedback from project execution trajectories.
2292         Your task is to analyze the provided trajectory data and
2293         extract information about the most relevant files mentioned
2294         in the given criteria.
2295
2296         Focus on the following:
2297
2298         1. Identify the **most recent steps** where the files directly
2299         related to the criteria were involved in execution, loading,
2300         or saving operations.
2301
2302         2. Provide environmental feedback for these files, such as any
2303         errors, warnings, or issues encountered during their
2304         execution or processing.
2305
2306         3. Highlight whether any problems occurred that might affect the
2307         functionality or success of these files in the project.
2308
2309         Your output should be structured as follows:
2310
2311         - **<RELEVANT STEPS>**: List the specific steps involving the
2312         relevant files, including any environmental feedback such as
2313         error messages, execution results, or other issues
2314         encountered. Each step should concisely present the key
2315         information needed to assess the files' execution status.
2316
2317         Avoid including details about file contents or existence, as this
2318         information is already available. Focus solely on the
2319         environmental feedback related to the execution of the most
2320         relevant files.
2321
2322         Your goal is to provide clear and concise information that helps
2323         determine if there were any execution problems with the files
2324         mentioned in the criteria.
2325         """
2326     else:
2327         raise NotImplementedError(f"The language '{language}' is not
2328         supported.")

```

2322 L.4 PROMPT FOR ASK MODULE (FOR REQUIREMENT CHECK)
 2323

```

2324 def get_ask_prompt(criteria: str, evidence: str) -> str:
2325
2326     return f"""
2327     Provided below is relevant information about the project:
2328     {evidence}
2329
2330     Kindly perform an evaluation of the following criteria:
2331     {criteria}
2332
2333     As per the guidelines, respond with either <SATISFIED> or <
2334     UNSATISFIED>, followed by a concise justification that references
2335     specific elements from the project information, such as code
2336     snippets, data samples, or output results.
2337     """
  
```

2338 L.5 PROMPT FOR LOCATE MODULE
 2339

```

2340 def get_locate_prompt(criteria: str, evidence: str) -> str:
2341
2342     return f"""
2343     Provided below is the structure of the workspace:
2344     {workspace_info}
2345
2346     This is the criteria related to the task:
2347     {criteria}
2348
2349     Follow the format in the example below and return only the file paths
2350     that match the criteria:
2351
2352     Example:
2353
2354     Suppose the criteria is:
2355     'The database functionality is implemented in `src/db.py`, and the
2356     logging system is defined in `src/logging.py`.'
2357
2358     And the workspace information is:
2359     /project
2360     |-- src
2361     |   |-- db.py
2362     |   |-- logging.py
2363     |   |-- utils.py
2364     |-- tests
2365     |   |-- test_db.py
2366     |   |-- test_logging.py
2367
2368     Based on the criteria, the following paths (no more than 5) should be
2369     returned, each wrapped in dollar signs (`$`):
2370     $/project/src/db.py$
2371     $/project/src/logging.py$
2372     """
  
```

2373
 2374
 2375

M JUDGE EVIDENCES COLLECTED FROM AGENT-AS-A-JUDGE

The Agent-as-a-Judge system dynamically collects evidence throughout the development process of the code-generating agentic system to evaluate whether each requirement is satisfied. Using this auxiliary information, the judge agent decides to focus on relevant data such as code, file structures, and trajectories. This gathered evidence supports the agent's final judgment, ensuring an informed and accurate assessment of the project's compliance with its requirements. A sample of the collected evidence is shown in the following Appendix M.

```

1 >>> [Reference] The Following Is the Original User Query:
2
3 Please help me build a system for sentiment analysis on tweets using the Sentiment140 dataset available from Mugging Face. Load the Sentiment140 dataset and,
when loading the data, clean it by removing stop words, punctuation, and special characters, all in 'src/data_loader.py'. Use Word2Vec or GloVe for text
vectorization. This should occur in 'src/data_loader.py'. Next, implement and train an SVM classifier in 'src/model.py'. Finally, write a report of the accuracy
of the classifier to 'results/metrics/accuracy_score.txt'. Ideally, the report should be easily interpretable.
4
5 >>> [Key Evidence] The Following Is the Project and Workspace Structure Information:
6
7 Project Tree
8
9 Workspace Path: OpenHands/workspaces/06_Sentiment_Analysis_SVM_Sentiment140_ML
10 Total Nodes: 5
11
12 Project Structure
13 |
14 |--- results
15 |--- results/metrics
16 |--- results/model
17 |--- src
18 |   |-- data_loader.py
19 |   |-- model.py
20 |
21 |
22 |
23 |
24 |
25 |
26 |
27 |
28 |
29 |
30 |
31 |
32 |
33 |
34 |
35 |
36 |
37 |
38 |
39 |
40 |
41 |
42 |
43 |
44 |
45 |
46 |
47 |
48 |
49 |
50 |
51 |
52 |
53 |
54 |
55 |
56 |
57 |
58 |
59 |
60 |
61 |
62 |
63 |
64 |
65 >>> [Reference] The Following Is Relevant Search Evidence:
66
67 Code Snippet
68
69 File: OpenHands/workspaces/06_Sentiment_Analysis_SVM_Sentiment140_ML/src/data_loader.py
70 Lines: (28, 28)
71 Identifier: ref
72 Category: function
73
74 28 dataset = load_and_preprocess_data()
75 29
76
77 Code Snippet
78
79 File: OpenHands/workspaces/06_Sentiment_Analysis_SVM_Sentiment140_ML/src/model.py
80 Lines: (18, 18)
81 Identifier: ref
82 Category: function
83
84 18 dataset = load_and_preprocess_data()
85 11
86
87 Code Snippet
88
89 File: OpenHands/workspaces/06_Sentiment_Analysis_SVM_Sentiment140_ML/src/data_loader.py
90 Lines: (12, 26)
91 Identifier: def
92 Category: function
93
94 12 def load_and_preprocess_data():
95 13     # Load the Sentiment140 dataset
96 14     dataset = load_dataset('sentiment140', split='train[1:3]', trust_remote_code=True)
97 15     # Clean the text data
98 16     dataset = dataset.map(lambda x: {'text': clean_text(x['text'])})
99 17     return dataset
100
101
102
103 >>> [Key Evidence] The Following Is the Contents of the Mentioned Files:
104
105 Content of OpenHands/workspaces/06_Sentiment_Analysis_SVM_Sentiment140_ML/src/data_loader.py:
106
107 import re
108 import nltk
109 from datasets import load_dataset
110 from nltk.corpus import stopwords
111 from nltk.tokenize import word_tokenize
112
113 nltk.download('stopwords')
114 nltk.download('punkt')
115
116 def clean_text(text):
117     # Remove special characters and punctuation
118     text = re.sub(r'([^\w\s\d\W\Z\z])', '', text)
119     # Tokenize the text
120     tokens = word_tokenize(text)
121     # Remove stop words
122     tokens = [word for word in tokens if word.lower() not in stopwords.words('english')]
123     return ' '.join(tokens)
124
125 def load_and_preprocess_data():
126     # Load the Sentiment140 dataset
127     dataset = load_dataset('sentiment140', split='train[1:3]', trust_remote_code=True)
128     # Clean the text data
129     dataset = dataset.map(lambda x: {'text': clean_text(x['text'])})
130     return dataset
131
132 if __name__ == "__main__":
133     dataset = load_and_preprocess_data()
134     print(dataset[0])
135
136
137 >>> [Reference] The Following Is Development Trajectory Evidence:
138
139 Relevant Steps in Trajectory
140
141 The following environment feedback is provided for reference only and does
142 not serve as decisive evidence. For final judgment, please rely primarily
143 on local file contents.
144
145 --Step 21--: The 'src/data_loader.py' script was executed to preprocess
146 the Sentiment140 dataset. The environment feedback indicated successful
147 execution of the script, with the dataset being loaded and preprocessed.
148 The script included downloading necessary NLTK data packages, such as
149 stopwords and punkt, which are essential for text processing tasks like
150 removing stop words and tokenization.
151
152 --Step 22--: The 'src/data_loader.py' script was run again to verify the
153 preprocessing steps. The environment feedback confirmed that the
154 preprocessing was completed successfully, and the first example from the
155 dataset was printed, indicating that the dataset was loaded and processed
156 correctly.
157
158 Overall, the preprocessing in 'src/data_loader.py' was executed without
159 errors, and the necessary steps for cleaning the dataset, such as removing
160 stop words, punctuation, and special characters, were implemented and
161 verified.
162
163
164 >>> [Reference] The Following Is Historical Judgments:
165
166 Requirement 1: The "Sentiment140" dataset, available from "Mugging Face," is obtained in 'src/data_loader.py'.
167 Satisfied: True
168
169

```

N ANALYSIS OF FAILURE CASES

We analyzed the failure cases and identified consistent patterns across the task categories. The categories of these failure cases are summarized in Table 8. We found that AAAJ struggled most with judging cases in the *Data preprocessing and postprocessing* category, whereas it performed well in judging *Human-Computer Interaction* cases.

Table 8: Counts of failure cases aggregated over task categories.

Category	Count
Data preprocessing and postprocessing	10
Dataset or Environment	8
Other	5
Machine Learning Method	4
Performance Metrics	3
Visualization	3
Human-Computer Interaction	3

We collect two failure cases in Table 9 and layout their analysis below.

Table 9: This table provides examples of failure cases observed in different task categories, comparing judgments made by Agent-as-a-Judge and Human-as-a-Judge.

Ex. #	Task	Req. ID	Category	Agent-as-a-Judge	Human-as-a-Judge	Criteria
1	40.Text.Summarization.BART_CNNDailyMail_DL	0	Dataset or Environment	False	True	The "CNN/Daily Mail" news dataset is used, including loading and preparing the dataset in <code>src/data_loader.py</code> .
2	46.Speech.Recognition.DeepSpeech_LibriSpeech_DL	2	Machine Learning Method	True	False	Hyperparameters such as learning rate and batch size are tuned in <code>src/train.py</code> .

Analysis

- Failure Case Example 1:** In this case, the dataset used was a synthesized one generated by the OpenHands CodeAct agent. Human evaluators could quickly identify this discrepancy, but the agent-as-a-judge, having only checked the file path and content, was misled into believing it was the genuine CNN/DailyMail dataset.
- Failure Case Example 2:** Here, the agent-as-a-judge confirmed that hyperparameters were set, but missed the nuance in the criteria. The requirement implied that the learning rate and batch size should dynamically adjust in `src/train.py`, something human evaluators were able to detect.

O SENSITIVITY W.R.T THE CHOICE OF THE BACKEND LLM

We have run an ablation experiment to determine how different LLM backends affect the performance of Agent-as-a-Judge. The results are summarized in Table 10.

Table 10: This table reports alignment percentages between Agent-as-a-Judge and Human-as-a-Judge for different backend LLMs.

Model	Version	# Param.	Alignment (%)
LLAMA, Touvron et al. (2023)	3.2	90B	87.76%
Qwen, Bai et al. (2023)	Coder 2.5	32B	88.73%
ChatGPT, OpenAI (2023) (This work)	gpt-4o-2024-0513	Unknown	90.16%
Claude, Anthropic (2024)	claude-3-5-sonnet-20241022	Unknown	92.95%

These results allow us to conclude that the backbone does have a noticeable effect on the alignment but a relatively marginal one. We found that Claude’s results are better than GPT-4o’s that we used throughout our experiments. We hypothesize, this is because `claude-3-5-sonnet-20241022` has been trained with strong function calling skills and agentic features.

P ADDITIONAL HUMAN EVALUATION DETAILS

All ten of our additional participants are current M.Sc. and Ph.D. students in AI-related fields with no direct relation to this work. The 7 random samples they evaluated were selected from the answers produced by OpenHands. The additional participants reported an average completion time of 1.13 hours each to evaluate all 7 samples. The additional ten participants self-reported an average completion time of 1.13 hours. This means they took an average of 9.67 minutes to evaluate one task, which is similar to our three main evaluators who self-reported taking an average of 10.48 minutes per task. Each of the ten additional participants was compensated for their time, with an average compensation of 15.20 USD.