

ClarEval: A Benchmark for Evaluating Clarification Skills of Code Agents under Ambiguous Instructions

Anonymous ACL submission

Abstract

To integrate seamlessly into real-world software engineering, Code Agents must evolve from passive instruction followers into proactive collaborative partners. However, current evaluation paradigms predominantly reward "guessing" user intent under ideal conditions, neglecting the agent's ability to align with users through dialogue—a critical trait for collaborative intelligence. In this work, we propose a paradigm shift in evaluation to drive this transition. We introduce ClarEval, a framework designed to assess an agent's "Collaborative Quotient" by simulating the inherent ambiguity of human communication. By systematically injecting three types of realistic ambiguity (missing goals, premises, and ambiguous terminology) into standard tasks, we force agents to step out of their "generator" role and engage in requirement elicitation. To quantify this capability, we propose a metric suite led by Average Turns to Clarify (ATC) and Key Question Coverage (KQC), which measure not just the correctness of the generated code, but the efficiency and precision of the collaboration. Our experiments on eleven state-of-the-art agents reveal a stark reality: while models like GPT-5-Coder excel at coding, they often lack the strategic communication skills required for efficient partnership. ClarEval thus serves as a crucial roadmap for bridging the gap between strong coders and capable collaborators.

1 Introduction

Large Language Models (LLMs)(Kasneci et al., 2023; Chang et al., 2024; Thirunavukarasu et al., 2023) have fundamentally reshaped software engineering, powering agents capable of automating coding, debugging, and maintenance. However, the prevailing evaluation paradigm(Guan et al., 2025; Chen et al., 2024) treats these agents as passive generators: given a static, self-contained instruction, the agent is expected to produce functionally

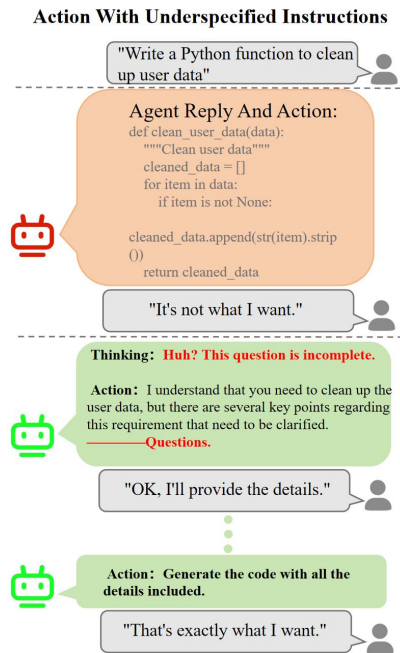


Figure 1: Contrasting agent behaviors when given an underspecified instruction. This example compares a passive generation approach, which fails due to ambiguity, with a proactive clarification approach, which successfully resolves uncertainty through dialogue before generating the correct code, highlighting the critical role of clarification skills.

correct code in a single turn. While this "Pass@1" centric approach has driven progress in syntactic correctness, it fails to capture a fundamental reality of software development: real-world requirements are rarely unambiguous or complete. In practice, instructions like "optimize system performance" are starting points for dialogue, not final specifications.

Current benchmarks(Liu et al., 2023b; Zhuo et al., 2024; Pan et al., 2024) inadvertently encourage agents to "guess" user intent to maximize pass rates, leading to misaligned outputs and eroded user trust. We argue that the next generation of Code Agents(Huang et al., 2023, 2024) must transcend passive generation to become Active Collaborative

057	Partners—agents capable of detecting ambiguity,	with ambiguous terminology and often resort to	106
058	inferring latent goals, and guiding users through	inefficient questioning strategies.	107
059	structured clarification. However, a critical gap	In summary, ClarEval establishes a standardized	108
060	remains: we lack a standardized diagnostic frame-	framework to advance Code Agents from passive	109
061	work to measure this “Collaborative Quotient”. Ex-	code generators to proactive collaborative partners.	110
062	isting evaluation paradigms either ignore interac-	By shifting the evaluation focus from correctness-	111
063	tion entirely or conflate communicative logic with	only to collaboration-aware, we provide the neces-	112
064	the noise of complex repository contexts, making	sary roadmap for developing trustworthy, human-	113
065	it difficult to isolate “why” an agent fails to clarify	aligned programming assistants.	114
066	effectively.		
067	To bridge this gap, we introduce ClarEval , a	2 Related Work	115
068	benchmark designed to systematically diagnose	2.1 Passive Code Generation vs. Real-World	116
069	the clarification logic of Code Agents. Unlike	Ambiguity	117
070	broad-spectrum evaluations, ClarEval functions	Traditional benchmarks like HumanEval (Chen	118
071	as a “unit test” for communicative intelligence.	et al., 2021), MBPP (Austin et al., 2021), and their	119
072	We deliberately construct our dataset by inject-	variants (Liu et al., 2023a ; Zhang et al., 2024a) evalu-	120
073	ing controlled ambiguities—Missing Goals (A_G),	uate functional correctness under the assumption of	121
074	Missing Premises (A_P), and Ambiguous Terminol-	ideal, self-contained specifications. This “Pass@1”	122
075	ogy (A_T)—into the canonical HumanEval dataset.	paradigm fundamentally diverges from real-world	123
076	While HumanEval focuses on algorithmic logic,	engineering, where requirements are rarely unam-	124
077	its isolated nature allows us to assess the agent’s	biguous. In these static settings, models are incen-	125
078	clarification strategy in a controlled environment,	tivized to “guess” latent intent to maximize pass	126
079	free from the confounding variables of massive	rates. Current benchmarks(Chen et al., 2025b ; Hu	127
080	codebases. If an agent cannot effectively clarify a	et al., 2025) lack mechanisms to penalize this guess-	128
081	sorting order in a list, it cannot be trusted to clarify	ing behavior or evaluate the agent’s initiative to	129
082	architectural decisions in a repository.	resolve underspecification through dialogue.	130
083	A key contribution of ClarEval is its focus on	2.2 Agentic Frameworks and Repository	131
084	interaction efficiency. Effective collaboration re-	Context	132
085	spects the user’s cognitive load; an agent that asks	The field has evolved from isolated generation to	133
086	ten vague questions is less useful than one that	autonomous agents (Schick et al., 2023 ; Yang et al.,	134
087	asks a single, precise question. We propose three	2024 ; Zhang et al., 2024b) operating within com-	135
088	targeted metrics to capture this nuance:	plex environments. Concurrently, benchmarks have	136
089		scaled to the repository level (Wu et al., 2024 ;	137
090	• Key Question Coverage (KQC): Measures	Puvvadi et al., 2025 ; Chen et al., 2025a), assessing	138
091	Intent Inference—the ability to identify the	code generation within broader project contexts.	139
092	high-level goal hidden behind a vague prompt.	However, these evaluations remain predominantly	140
093	• Missing Premises Recall (MPR): Assesses	execution-centric. They measure whether an agent	141
094	Requirement Identification—ensuring the	can solve a task given a fixed context, but over-	142
095	agent can surface necessary technical con-	look the communicative intelligence required to	143
096	straints before coding.	elicit and define the task itself, failing to capture	144
097	• Average Turns to Clarify (ATC): A novel	the collaborative dynamics essential for human-AI	145
098	metric quantifying Clarification Strategy. It	pairing.	146
099	penalizes inefficient, “lazy,” or redundant	2.3 Benchmarking Clarification and	147
100	questioning, favoring agents that reach clarity	Efficiency	148
101	with minimal user friction.	The evaluation of collaborative skills in code agents	149
102	We evaluate eleven state-of-the-art agents using	is an emerging frontier. Prior benchmarks for agen-	150
103	this framework. Our experiments reveal a stark real-	tic software engineering largely focus on complex,	151
104	ity: while models like GPT-5-Coder possess strong	multi-step implementation and real-world task suc-	152
105	reasoning capabilities, most agents lack system-	cess: SWE-bench-Live (Zhang et al., 2025) as-	153
105	atic clarification skills. They struggle significantly		

BenchMark	Task Type	Ambiguity	Single-round	Multi-round	Efficiency indicators
LoCoBench-Agent	software engineering	×	✓	✓	✓
SR-Eval	Code Generation	×	✓	×	×
CodeFlowBench	Code Generation	×	✓	✓	×
SWE-bench-Live	Issue Resolution	×	✓	✓	×
HumanEval	Code Generation	×	✓	×	×
HumanEvalComm	Code Generation	✓	✓	✓	×
ClarEval	Code Generation	✓	✓	✓	✓

Table 1: Benchmark comparison

154 assesses bug-fixing robustness, while SR-Eval (Zhan
155 et al., 2025), CodeFlowBench (Wang et al., 2025),
156 and LoCoBench-Agent (Qiu et al., 2025) evaluate
157 multi-turn, long-context code generation under iterative or long-context requirements. However, these
158 works do not quantify the quality or efficiency of
159 the communication itself.
160

161 Most notably, HumanEvalComm (Wu and Fard,
162 2025; McNall et al., 2024) pioneered the injection
163 of ambiguity to test clarification competence. While this work successfully highlights the importance of dialogue, it focuses primarily on the binary
164 capability of asking questions (e.g., Good-Question
165 metric) and the ultimate correctness of the code.
166
167

168 ClarEval distinguishes itself by shifting the focus from "Task Success" to "Communication Efficiency and Diagnosis." (See Table 1 for a detailed
169 comparison). We argue that an agent asking excessive or fragmented questions imposes a high cognitive tax on the user. Unlike prior works that may
170 reward inefficient questioning, ClarEval introduces
171 metrics like Average Turns to Clarify (ATC) and the aggregate Efficiency-Adjusted Recall (EAR) to explicitly penalize user friction. By employing a
172 controlled ambiguity taxonomy and a deterministic simulator, ClarEval provides a precise, reproducible unit test for collaborative intelligence, prioritizing
173 agents that achieve clarity with minimal interaction cost.
174
175
176
177
178
179
180
181
182

183 3 Method

184 In this section, we elaborate on the construction
185 method, evaluation scenarios, and evaluation metrics of the ClarEval benchmark. Our goal is to provide a systematic and reproducible evaluation
186 framework for assessing the proactive clarification capability of code agents. Figure 2 overviews the construction of the ClarEval.
187
188
189
190

3.1 Benchmark Construction

191 **Source Data:** To construct a robust and
192 contamination-aware benchmark, we integrate
193 the HumanEval dataset (Chen et al., 2021) with
194 a challenging subset from the LiveCodeBench
195 dataset (Jain et al., 2024). The resulting collection
196 comprises $N = 750$ curated programming
197 problems, consisting of 150 tasks from HumanEval
198 and 600 tasks from LiveCodeBench. We formalize
199 this as a collection of clear source tasks $\mathcal{T} =$
200 $\{T_1, T_2, \dots, T_N\}$, where each task $T_i = (D_i, S_i)$
201 consists of a clear natural language description D_i
202 and its corresponding canonical solution S_i . By
203 applying three distinct ambiguity injection strategies
204 to each of the 750 source tasks, we constructed
205 a comprehensive benchmark containing a total of
206 2,250 unique evaluation instances. **Ambiguity Injection:**
207 The core of our approach is the systematic
208 transformation of a clear task T_i into an ambiguous
209 task \tilde{T}_i . This is achieved through an ambiguity
210 injection function \mathcal{F} , formally defined as:
211

$$212 \tilde{T}_i = \mathcal{F}(T_i, A) = (\tilde{D}_i, S_i) \quad (1)$$

213 where $A \in \{A_G, A_P, A_T\}$ denotes the specific
214 type of ambiguity applied. We define three fundamental
215 ambiguity types:

- 216 • **Missing Goal (A_G):** The high-level objective
217 is omitted. Formally, if D_i can be decomposed
218 into context C and goal G , i.e., $D_i = C \oplus G$ (where
219 \oplus denotes string concatenation), then the ambiguous
220 description is defined as $\tilde{D}_i = C$.
221
- 222 • **Missing Premises (A_P):** Critical technical
223 constraints are removed. Let $\mathcal{P} = \{p_1, p_2, \dots, p_k\}$
224 be the set of missing premises. The ambiguous description
225 is generated by $\tilde{D}_i = \text{REMOVE}(D_i, \mathcal{P})$, where
226 REMOVE eliminates all textual references to
227 premises in \mathcal{P} .
228

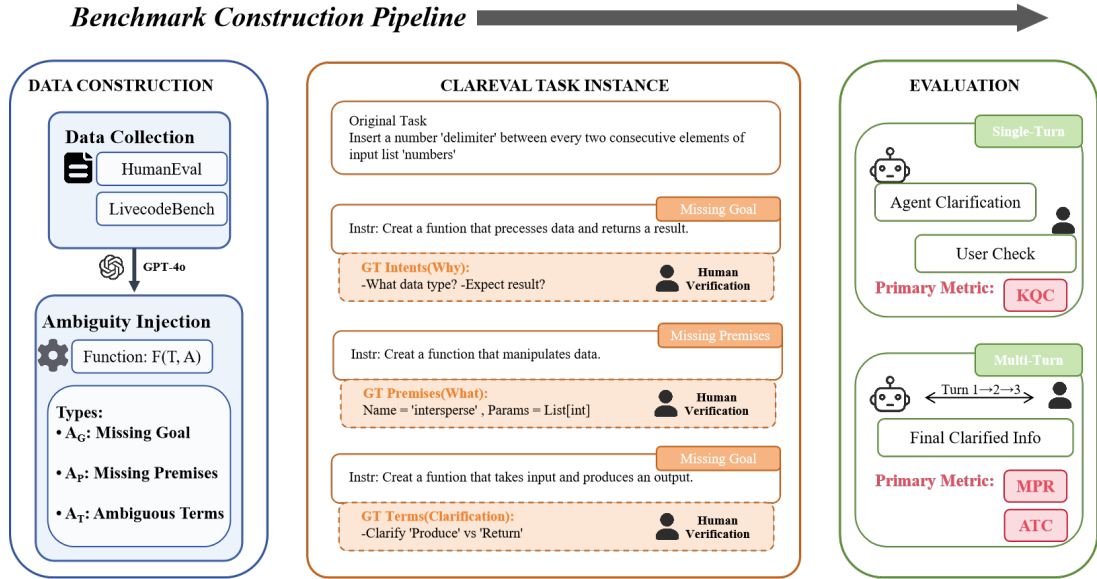


Figure 2: An overview for the ClarEval construction. First, original problems are obtained from the HumanEval dataset. Then, based on the three defined types of ambiguity—Missing Goal, Missing Premises, and Ambiguous Terminology—the original problems are transformed into ambiguous task descriptions using GPT-4o. Finally, the benchmark is refined and validated through human verification and annotation.

- **Ambiguous Terminology (A_T):** A precise term t is replaced with a vague term \tilde{t} . The transformation is defined as $\tilde{D}_i = \text{REPLACE}(D_i, t, \tilde{t})$.

For each original task T_i , we construct one or more ambiguous variants \tilde{T}_i by applying \mathcal{F} , creating our benchmark suite. Each ambiguous task is annotated with a golden standard that defines the information required for complete disambiguation.

3.2 Ambiguity Injection via Prompt Engineering

The core of ClarEval lies in the systematic transformation of clear tasks into ambiguous variants. To ensure reproducibility and diversity, we employed GPT-4o as a transformation engine, guided by specific Ambiguity Injection Prompts. Instead of random token deletion, we designed three distinct prompt strategies to simulate realistic patterns of human underspecification. Table 4 provides a comprehensive overview of these strategies and their corresponding transformation examples.

3.3 Evaluation Settings

We design two evaluation settings to model distinct interaction paradigms between users and code agents.

3.3.1 Single-Turn Setting

This setting evaluates an agent’s ability to identify ambiguity and pose critical clarification questions immediately upon receiving an initial ambiguous instruction. The agent is provided only with the ambiguous task description \tilde{D}_i and generates a single response. Performance is evaluated based on the clarification questions contained in this initial reply. This setting primarily assesses the agent’s initial intent inference capability (Key Question Coverage, KQC) and its zero-shot ability to infer missing premises (Single-Turn MPR), serving as a measure of initial planning quality.

3.3.2 Multi-Turn Setting

This setting assesses an agent’s capacity for sustained dialogue, including context management, iterative guidance, and probing follow-up questions. We pre-define structured conversation scripts for each task. Starting from the ambiguous instruction \tilde{D}_i , each script specifies potential user responses (which may provide partial information or introduce new ambiguities), enabling evaluation of the agent’s strategic questioning throughout extended interactions.

Crucially, the evaluation of the multi-turn setting is performed using a deterministic, rule-based User Simulator script (\mathcal{S}_i) derived from human-annotated ground truth intents (see Appendix C and D). This ensures the reproducibility and con-

sistency of agent evaluations across different models, isolating communicative logic from dynamic simulation noise. This setting focuses on the agent’s complete requirement elicitation (Dialogue-Informed Missing Premises Recall, MPR) and efficiency (Average Turns to Clarify, ATC) over the conversation .

3.4 Human Verification Protocol

To ensure validity, we implemented a stringent two-stage review process involving three senior software engineers.

Criteria Check. Reviewers assessed each entry based on three criteria: (1) **Ambiguity Validity:** Is the instruction genuinely unsolvable without assumptions? (2) **Plausibility:** Does the request resemble natural user prompts? (3) **Script Alignment:** Does the dialogue script accurately reflect the ground truth?

Ecological Validity Analysis. To validate the realism of synthetic ambiguity, we compared 50 ClarEval instances against real StackOverflow queries. Two experts blindly rated “Naturalness” on a Likert scale (1-5). Results show ClarEval’s score (4.12 ± 0.6) is statistically indistinguishable from real-world data (4.25 ± 0.8 , $p > 0.05$). This confirms that our ambiguity injection pipeline successfully mimics the “colloquial vagueness” found in human inquiries, ensuring the benchmark tests realistic communicative intelligence rather than synthetic pattern matching.

Consensus and Refinement. We employed a majority-vote mechanism for retaining samples, achieving a Fleiss’ Kappa of $\kappa = 0.82$ (substantial agreement). Ultimately, approximately 12% of the raw generated samples were discarded or rewritten to correct hallucinations, ensuring ClarEval serves as a reliable standard.

3.5 Evaluation Metrics

We propose a suite of metrics to quantify clarification capability across two distinct phases: initial intent inference (single-turn) and sustained dialogue (multi-turn).

3.5.1 Single-Turn Metrics

These metrics evaluate the agent’s immediate response to an ambiguous instruction, assessing its ability to identify the core problem without the benefit of dialogue.

Key Question Coverage (KQC) measures the agent’s success in inferring the high-level user intent obscured by ambiguity. Let $\mathcal{K} = \{k_1, k_2, \dots, k_m\}$ be the set of key intents required for disambiguation, as determined by expert annotation. Let $\mathcal{C} \subseteq \mathcal{K}$ be the set of intents covered by the agent’s clarification questions. KQC is defined as:

$$\text{KQC} = \frac{|\mathcal{C}|}{|\mathcal{K}|} \quad (2)$$

A higher KQC indicates superior initial intent inference capability, while discouraging agents from generating excessive, low-quality questions.

Premise Identification Rate (PIR) assesses the agent’s zero-shot ability to identify missing technical premises from the initial ambiguous instruction. Let $\mathcal{P} = \{p_1, p_2, \dots, p_n\}$ be the golden set of missing premises. Let $\mathcal{I} \subseteq \mathcal{P}$ be the set of premises detected in the agent’s first-turn questions. PIR is defined as:

$$\text{PIR} = \frac{|\mathcal{I}|}{|\mathcal{P}|} \quad (3)$$

PIR measures planning quality, indicating whether the agent recognizes what technical details are missing, even if it cannot yet resolve them through dialogue.

3.5.2 Multi-Turn Metrics

These metrics evaluate the agent’s performance throughout an extended dialogue, assessing its capacity for context management, iterative questioning, and efficient collaboration.

Key Question Coverage (KQC) in the multi-turn setting measures the agent’s success in inferring high-level user intent over the entire dialogue. Similar to the single-turn setting, let \mathcal{K} be the set of key intents and $\mathcal{C}_{\text{dialogue}} \subseteq \mathcal{K}$ be the set of intents covered by the agent’s questions across all turns. Multi-turn KQC is defined as:

$$\text{KQC} = \frac{|\mathcal{C}_{\text{dialogue}}|}{|\mathcal{K}|} \quad (4)$$

Missing Premises Recall (MPR) evaluates an agent’s completeness in resolving missing technical details over the entire dialogue. Let $\mathcal{P} = \{p_1, p_2, \dots, p_n\}$ be the golden set of missing premises. Let $\mathcal{R} \subseteq \mathcal{P}$ be the set of premises successfully clarified by the agent by the dialogue’s end. MPR is defined as:

$$\text{MPR} = \frac{|\mathcal{R}|}{|\mathcal{P}|} \quad (5)$$

While PIR focuses on "what is missing" in a zero-shot setting, MPR focuses on "what was successfully recovered" through sustained dialogue.

Average Turns to Clarify (ATC) quantifies interaction efficiency by measuring how quickly an agent secures necessary clarifications. For each clarified premise $p_i \in \mathcal{R}$, let $\text{Turn}(p_i)$ denote the dialogue turn index in which it was first successfully resolved. The ATC is calculated as:

$$\text{ATC} = \frac{1}{|\mathcal{R}|} \sum_{p_i \in \mathcal{R}} \text{Turn}(p_i) \quad (6)$$

A lower ATC indicates better efficiency, signifying that the agent reaches clarity more quickly. Note that in cases where no premises are clarified (i.e., $|\mathcal{R}| = 0$), ATC is considered undefined and treated as NaN (not included in averages).

Efficiency-Adjusted Recall (EAR) provides a unified measure that balances clarification completeness (MPR) with interaction efficiency (ATC), penalizing agents that achieve high recall through excessive dialogue turns. The formal definition is provided in Appendix B.

4 Experiment

4.1 Experimental Setup

4.1.1 Models

We evaluate eleven state-of-the-art code agents: Aider-GPT5 (Gauthier, 2025) (the Aider code editing framework powered by GPT-5), CodeX-GPT5 (OpenAI, 2025), GPT5-Coder (OpenAI, 2025a; Achiam et al., 2023), CodeLlama-70B (Roziere et al., 2023), Claude-Opus 4.1 (Anthropic, 2025a), Claude-Sonnet 4.5 (Anthropic, 2025b), O3, O4 (OpenAI, 2025b), DeepSeek-Coder (Guo et al., 2025), Qwen2.5-Coder (Hui et al., 2024), and Qwen3-Coder (Yang et al., 2025). These models represent the current spectrum of commercial and open-source code agents, spanning diverse architectures (e.g., decoder-only, tool-augmented) and training paradigms (e.g., supervised fine-tuning, reinforcement learning).

4.1.2 Evaluation Protocol

Each model is evaluated on the ClarEval benchmark under both single-turn and multi-turn settings. In the single-turn setting, models receive only the ambiguous instruction and generate a single response containing clarification questions. In the multi-turn setting, we employ pre-defined dialogue

scripts that simulate realistic user interactions, including partial answers and occasional introduction of new ambiguities. All tasks are derived from HumanEval via controlled injection of three ambiguity types: Missing Goal, Missing Premises, and Ambiguous Terminology. All experiments were conducted with temperature = 0.1 on NVIDIA A100 GPUs. Detailed hyperparameters and prompt templates are provided in Appendix A.

4.2 Overall Results

Table 2 presents the overall performance across all ambiguity types, revealing significant variations in clarification capability among state-of-the-art code agents.

4.2.1 Single-Turn Performance

In the single-turn setting, GPT5-Coder achieves the highest KQC (0.867), demonstrating superior capability in inferring high-level user intents from ambiguous instructions. Qwen2.5-Coder (0.855) and Claude-Sonnet 4.5 (0.852) also perform strongly, though they lag behind GPT5-Coder by 8–9 percentage points.

For MPR, GPT5-Coder again leads (0.661), substantially outperforming the next-best models—Claude-Sonnet 4.5 (0.631) and Claude-Opus 4.1 (0.629)—by approximately 20 percentage points. This indicates that while several models can identify missing technical premises in isolation, only top-tier agents do so consistently across diverse ambiguity patterns.

4.2.2 Multi-Turn Performance

In the multi-turn setting, performance dynamics shift significantly. Claude-Opus 4.1 achieves the highest KQC (0.754), suggesting strong contextual understanding and effective follow-up questioning. For MPR, GPT5-Coder maintains its lead (0.951), closely followed by Claude-Opus 4.1 (0.930).

The ATC metric reveals a critical efficiency gap: GPT5-Coder achieves the lowest ATC (1.493), indicating it secures necessary clarifications in fewer dialogue turns. In stark contrast, Qwen2.5-Coder exhibits the highest ATC (2.396), requiring nearly 60% more turns on average to clarify the same set of premises.

4.3 Key Findings and Analysis

The Single-Turn vs. Multi-Turn Dichotomy:

We observe a weak correlation between single-turn and multi-turn performance (Pearson $r = 0.32$,

Model	Single-Turn		Multi-Turn		
	KQC	PIR	KQC	MPR	ATC
Aider-GPT5	0.833	0.588	0.376	0.396	1.576
CodeX-GPT5	0.835	0.557	0.626	0.748	1.681
GPT5-Coder	<u>0.867</u>	<u>0.661</u>	0.529	<u>0.951</u>	<u>1.493</u>
CodeLlama-70B	0.817	0.535	0.427	0.385	1.602
Claude-opus4.1	0.847	0.629	<u>0.754</u>	0.930	1.900
Claude-sonnet-4.5	0.852	0.631	0.633	0.823	1.719
O3	0.794	0.572	0.573	0.791	1.771
O4	0.774	0.544	0.538	0.712	1.576
DeepSeek-Coder	0.792	0.588	0.675	0.796	1.772
Qwen2.5-Coder	0.855	0.622	0.544	0.643	2.396
Qwen3-Coder	0.816	0.592	0.668	0.640	1.645

Table 2: Performance comparison of code agents on ClarEval across Single-Turn and Multi-Turn settings. Metrics include KQC and MPR, where higher values indicate better clarification effectiveness, along with ATC, where lower values indicate superior efficiency. Underline indicates the best-performing model for each metric.

$p < 0.05$). For instance, Aider-GPT5 ranks 2nd in single-turn KQC (0.833) but drops to 10th in multi-turn KQC (0.376). Similarly, CodeLlama-70B shows a dramatic decline from single-turn (0.817) to multi-turn settings (0.427). This suggests that intent inference in isolation and dialogue-aware clarification are distinct skills, underscoring the necessity of evaluating both settings in a comprehensive benchmark.

The Clarification Capability Spectrum:

Our results reveal a clear spectrum of clarification competence. At the high end, GPT5-Coder and Claude-Opus 4.1 demonstrate balanced strength in intent understanding (KQC), technical detail recovery (MPR), and interaction efficiency (ATC). In the middle tier, models like DeepSeek-Coder and Claude-Sonnet 4.5 show competent but inconsistent performance. At the lower end, Aider-GPT5 and CodeLlama-70B struggle significantly in multi-turn settings, despite decent single-turn scores, revealing fundamental limitations in their interactive capabilities.

Ambiguous Terminology as the Primary Challenge: As detailed in the fine-grained analysis, Ambiguous Terminology emerges as the most challenging ambiguity type across models and metrics. This suggests that while agents can readily detect missing information, they struggle to recognize when existing terms are vague or underspecified (e.g., "fast", "clean up", "user-friendly"). Given that such terminology is pervasive in real developer requests, this represents a key bottleneck for

■ Baseline (Ambiguous) ■ Upper Bound (Clarified)

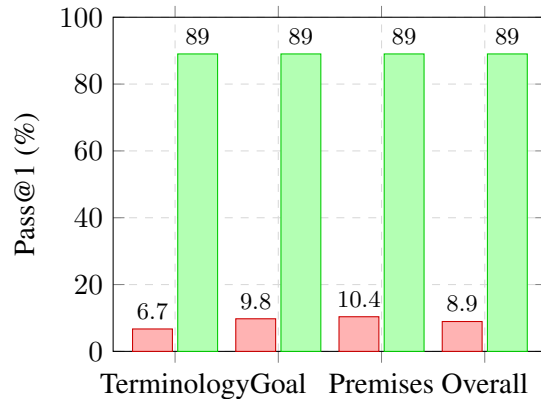


Figure 3: Performance Gap Analysis. Comparison of Pass@1 accuracy between the ambiguous baseline and the clarified upper bound across different ambiguity types. *Ambiguous Terminology* proves to be the most challenging scenario for passive code generation.

practical deployment. Detailed breakdowns in Appendix Tables 10-14 confirm this trend, showing lowest Pass@1 and highest ATC for terminology ambiguities across most models.

5 Discussion

5.1 The Impact of Clarification on Code Correctness

To validate the practical utility of ClarEval, we conducted a controlled ablation study to quantify the impact of ambiguity on code generation correctness (Pass@1). We compared the performance

512 of GPT-4o under two distinct conditions: an Am- 561
513 biguous Baseline (generating code directly from 562
514 underspecified instructions) and a Clarified Upper 563
515 Bound (generating code after resolving ambigu- 564
516 ities via oracle guidance). Note that for the baseline, 565
517 we applied an oracle patch to fix function name 566
518 mismatches, ensuring the scores reflect logic under- 567
519 standing rather than trivial interface misalignment.

520 **The Catastrophic Ambiguity Gap.** Figure 3 568
521 illustrates a stark reality: despite the reasoning 569
522 prowess of frontier models, ambiguity remains a 570
523 catastrophic failure mode. While GPT-4o achieves 571
524 a state-of-the-art Pass@1 of 89.02% on clarified 572
525 tasks, its performance plummets to 8.94% under 573
526 ambiguity. This substantial $\sim 80\%$ performance 574
527 gap empirically demonstrates that without proac- 575
528 tive clarification, even the most capable LLM- 576
529 based agents fail to generate functionally correct 577
530 code. Clarification is not merely an enhancement; 578
531 it is a prerequisite for correctness. 579

532 **Semantic Vagueness as the Primary Bottleneck.** 580
533 Decomposing the results by ambiguity type re- 581
534 veals that Ambiguous Terminology (A_T) poses the 582
535 most severe challenge, yielding the lowest Pass@1 583
536 of 6.71% (compared to 9.76% for Missing Goal 584
537 and 10.37% for Missing Premises). This indicates 585
538 that semantic vagueness—where user intent is ob- 586
539 scured by colloquialisms (e.g., “clean up”, “orga- 587
540 nize”)—introduces significantly higher variance in 588
541 implementation logic than missing technical con- 589
542 straints. 590

543 5.2 Beyond Detection: The Efficiency Frontier 591

544 The high KQC scores achieved by top-tier mod- 592
545 els (e.g., GPT-5-Coder at 0.867) suggest that the 593
546 community has largely solved the “Ambiguity De- 594
547 tectio” problem for standard tasks. However, this 595
548 does not render the benchmark obsolete; rather, 596
549 it shifts the evaluation focus from competency to 597
550 optimization. 598

551 The significant gaps exposed by ClarEval in the 599
552 multi-turn setting and ATC metric demonstrate that 600
553 current training paradigms prioritize correctness 601
554 over collaboration costs. Effective code agents 602
555 must evolve from “Passive Generators” that guess 603
556 intent, to “Active Partners” that navigate ambiguity 604
557 with minimal friction. 605

558 5.3 Implications for Agent Design 606

559 The performance variations across ClarEval reveal 607
560 critical directions for future agent development: 608

- 561 • **Training for User-Centric Planning:** The 562
563 poor efficiency of serial questioners points to 564
565 a lack of user simulation in the training loop. 566
567 Future training objectives should penalize ex-
568 cessive turns, encouraging models to antici-
569 pate multiple dependencies before yielding
570 control to the user. 571
- 572 • **Semantic Grounding for Terminology:** The 573
574 consistent failure on Ambiguous Terminol-
575 ogy (A_T) suggests a disconnect between code
576 reasoning and general semantic grounding. 577
578 Agents require improved grounding in both
579 technical domains and general world knowl-
580 edge to bridge the gap between vague human
581 intent and precise executable logic. 582
583
- 584 • **Adaptive Clarification Policies:** The ob- 585
586 served ATC variations indicate that a one-size-
587 fits-all questioning strategy is suboptimal. Fu-
588 ture agents should develop adaptive policies
589 that balance comprehensiveness (asking ev-
590 erything) and efficiency (asking only what is
591 critical) based on the complexity of the con-
592 text. 593

594 6 Conclusion 595

596 This work argues for a fundamental realignment 597
598 in how we evaluate Code Agents: moving beyond 599
600 the static “passive generation” paradigm toward 601
602 active, efficiency-aware collaboration. We intro-
603 duced ClarEval, not merely as a dataset, but as a
604 diagnostic framework to measure this transition. 605
606 By isolating communicative logic in a controlled
607 environment and introducing efficiency-penalized
608 metrics like ATC and EAR, we exposed a criti-
609 cal “Collaborative Gap” in current state-of-the-art
610 models. 611

612 Our findings reveal that while modern agents 613
614 have mastered the syntax of coding, they struggle 615
616 with the semantics of partnership. High-performing
617 coders often devolve into inefficient interrogators
618 under ambiguity, burdening users with redundant
619 dialogue. ClarEval provides the necessary granular-
620 ity to detect and correct these interaction behaviors.
621 Ultimately, we hope this benchmark serves as a
622 catalyst for the community to prioritize Collabora-
623 tive Quotient alongside code correctness, paving
624 the way for agents that are not just powerful tools,
625 but intuitive partners in the software engineering
626 lifecycle. 627

609 **Limitations**

610 While ClarEval provides comprehensive cover-
611 age of common ambiguity types, real-world sce-
612 narios may involve more complex or composite
613 forms of ambiguity. Future work could expand
614 the benchmark to include cross-cutting ambiguities
615 and domain-specific underspecification. Addition-
616 ally, exploring how clarification capability corre-
617 lates with ultimate task success in real development
618 environments would further validate the practical
619 significance of our metrics.

620

References

621
622
623
624
625

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, and 1 others. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.

626
627
628

Anthropic. 2025a. Claude opus 4.1. <https://www.anthropic.com/news/claude-opus-4-1>.

629
630

Anthropic. 2025b. Claude sonnet 4. Available online: <https://www.anthropic.com/claude/opus>.

631
632
633
634
635

Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, and 1 others. 2021. Program synthesis with large language models. *arXiv preprint arXiv:2108.07732*.

636
637
638
639
640
641

Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu, Linyi Yang, Kaijie Zhu, Hao Chen, Xiaoyuan Yi, Cunxiang Wang, Yidong Wang, and 1 others. 2024. A survey on evaluation of large language models. *ACM transactions on intelligent systems and technology*, 15(3):1–45.

642
643
644
645
646

Chen Chen, Xinlong Hao, Weiwen Liu, Xu Huang, Xingshan Zeng, Shuai Yu, Dexun Li, Shuai Wang, Weinan Gan, Yuefeng Huang, and 1 others. 2025a. Acebench: Who wins the match point in tool usage? *arXiv preprint arXiv:2501.12851*.

647
648
649
650
651
652

Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, and 1 others. 2021. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*.

653
654
655
656
657

Simin Chen, Xiaoning Feng, Xiaohong Han, Cong Liu, and Wei Yang. 2024. Ppm: Automated generation of diverse programming problems for benchmarking code generation models. *Proceedings of the ACM on Software Engineering*, 1(FSE):1194–1215.

658
659
660
661

Simin Chen, Pranav Pusarla, and Baishakhi Ray. 2025b. Dynamic benchmarking of reasoning capabilities in code large language models under data contamination. *arXiv preprint arXiv:2503.04149*.

662
663
664

Paul Gauthier. 2025. Aider: Ai-assisted coding in your terminal with gpt. <https://aider.chat/>, 2023. Accessed: 2025-05-15.

665
666
667
668

Batu Guan, Xiao Wu, Yuanyuan Yuan, and Shaohua Li. 2025. Is your benchmark (still) useful? dynamic benchmarking for code language models. *arXiv preprint arXiv:2503.06643*.

669
670
671
672
673
674

Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.

Wenhao Hu, Jinhao Duan, Chunchen Wei, Li Zhang, Yue Zhang, and Kaidi Xu. 2025. Dynacode: A dynamic complexity-aware code benchmark for evaluating large language models in code generation. *arXiv preprint arXiv:2503.10452*. 675
676
677
678
679

Dong Huang, Jie M Zhang, Michael Luck, Qingwen Bu, Yuhao Qing, and Heming Cui. 2023. Agent-coder: Multi-agent-based code generation with iterative testing and optimisation. *arXiv preprint arXiv:2312.13010*. 680
681
682
683
684

Yiming Huang, Jianwen Luo, Yan Yu, Yitong Zhang, Fangyu Lei, Yifan Wei, Shizhu He, Lifu Huang, Xiao Liu, Jun Zhao, and 1 others. 2024. Da-code: Agent data science code generation benchmark for large language models. *arXiv preprint arXiv:2410.07331*. 685
686
687
688
689

Binyuan Hui, Jian Yang, Zeyu Cui, Jiayi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang, Bowen Yu, Keming Lu, and 1 others. 2024. Qwen2. 5-coder technical report. *arXiv preprint arXiv:2409.12186*. 690
691
692
693
694

Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando Solar-Lezama, Koushik Sen, and Ion Stoica. 2024. Live-codebench: Holistic and contamination free evaluation of large language models for code. *arXiv preprint arXiv:2403.07974*. 695
696
697
698
699
700

Enkelejd Kasneci, Kathrin Seßler, Stefan Küchemann, Maria Bannert, Daryna Dementieva, Frank Fischer, Urs Gasser, Georg Groh, Stephan Günnemann, Eyke Hüllermeier, and 1 others. 2023. Chatgpt for good? on opportunities and challenges of large language models for education. *Learning and individual differences*, 103:102274. 701
702
703
704
705
706
707

Jiawei Liu, Chunqiu Steven Xia, Yuyao Wang, and Lingming Zhang. 2023a. Is your code generated by chatgpt really correct? rigorous evaluation of large language models for code generation. *Advances in Neural Information Processing Systems*, 36:21558–21572. 708
709
710
711
712
713

Tianyang Liu, Canwen Xu, and Julian McAuley. 2023b. Repobench: Benchmarking repository-level code auto-completion systems. *arXiv preprint arXiv:2306.03091*. 714
715
716
717

Kyle McNall, Yiyun Qiao, and 1 others. 2024. HumanEvalComm: Benchmarking the communication competence of code generation. *arXiv preprint arXiv:2406.00215*. 718
719
720
721

OpenAI. 2025a. Introducing gpt-5. <https://openai.com/zh-Hans-CN/index/introducing-gpt-5/>. 722
723
724

OpenAI. 2025b. Introducing openai o3 and o4-mini. <https://openai.com/zh-Hans-CN/index/introducing-o3-and-o4-mini/>. 725
726
727

OpenAI. 2025. Openai codex cli: Lightweight coding agent that runs in your terminal. <https://github.com/openai/codex>, 2025. 728
729
730

731	Zhenyu Pan, Rongyu Cao, Yongchang Cao, Yingwei Ma, Binhua Li, Fei Huang, Han Liu, and Yongbin Li. 2024. Codev-bench: How do llms understand developer-centric code completion? <i>arXiv preprint arXiv:2410.01353</i> .	787
732		788
733		789
734		790
735		
736	Meghana Puvvadi, Sai Kumar Arava, Adarsh Santoria, Sesha Sai Prasanna Chennupati, and Harsha Vardhan Puvvadi. 2025. Coding agents: A comprehensive survey of automated bug fixing systems and benchmarks. In <i>2025 IEEE 14th International Conference on Communication Systems and Network Technologies (CSNT)</i> , pages 680–686. IEEE.	791
737		792
738		793
739		794
740		
741		
742		
743	Jielin Qiu, Zuxin Liu, Zhiwei Liu, Rithesh Murthy, Jianguo Zhang, Haolin Chen, Shiyu Wang, Ming Zhu, Liangwei Yang, Juntao Tan, and 1 others. 2025. Locobench-agent: An interactive benchmark for llm agents in long-context software engineering. <i>arXiv preprint arXiv:2511.13998</i> .	795
744		796
745		797
746		798
747		799
748		800
749	Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Romain Sauvestre, Tal Remez, and 1 others. 2023. Code llama: Open foundation models for code. <i>arXiv preprint arXiv:2308.12950</i> .	801
750		802
751		803
752		804
753		805
754	Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Eric Hambro, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. 2023. Toolformer: Language models can teach themselves to use tools. <i>Advances in Neural Information Processing Systems</i> , 36:68539–68551.	806
755		807
756		808
757		809
758		810
759		
760	Arun James Thirunavukarasu, Darren Shu Jeng Ting, Kabilan Elangovan, Laura Gutierrez, Ting Fang Tan, and Daniel Shu Wei Ting. 2023. Large language models in medicine. <i>Nature medicine</i> , 29(8):1930–1940.	811
761		812
762		813
763		814
764		815
765		816
766	Sizhe Wang, Zhengren Wang, Dongsheng Ma, Yongan Yu, Rui Ling, Zhiyu Li, Feiyu Xiong, and Wentao Zhang. 2025. Codeflowbench: A multi-turn, iterative benchmark for complex code generation. <i>arXiv preprint arXiv:2504.21751</i> .	
767		
768		
769		
770	Jie JW Wu and Fatemeh H Fard. 2025. Humaneval-comm: Benchmarking the communication competence of code generation for llms and llm agents. <i>ACM Transactions on Software Engineering and Methodology</i> , 34(7):1–42.	
771		
772		
773		
774		
775	Qinyun Wu, Chao Peng, Pengfei Gao, Ruida Hu, Haoyu Gan, Bo Jiang, Jinhe Tang, Zhiwen Deng, Zhanming Guan, Cuiyun Gao, and 1 others. 2024. Repomas-tereval: Evaluating code completion via real-world repositories. <i>arXiv preprint arXiv:2408.03519</i> .	
776		
777		
778		
779		
780	An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, and 1 others. 2025. Qwen3 technical report. <i>arXiv preprint arXiv:2505.09388</i> .	
781		
782		
783		
784		
785	John Yang, Carlos E Jimenez, Alexander Wettig, Kilian Lieret, Shunyu Yao, Karthik Narasimhan, and Ofir	
786		
	Press. 2024. Swe-agent: Agent-computer interfaces enable automated software engineering. <i>Advances in Neural Information Processing Systems</i> , 37:50528–50652.	
	Zexun Zhan, Shuzheng Gao, Ruida Hu, and Cuiyun Gao. 2025. Sr-eval: Evaluating llms on code generation under stepwise requirement refinement. <i>arXiv preprint arXiv:2509.18808</i> .	
	Fengji Zhang, Linqun Wu, Huiyu Bai, Guancheng Lin, Xiao Li, Xiao Yu, Yue Wang, Bei Chen, and Jacky Keung. 2024a. Humaneval-v: Evaluating visual understanding and reasoning abilities of large multimodal models through coding tasks. <i>arXiv preprint arXiv:2410.12381</i> .	
	Kechi Zhang, Jia Li, Ge Li, Xianjie Shi, and Zhi Jin. 2024b. Codeagent: Enhancing code generation with tool-integrated agent systems for real-world repo-level coding challenges. <i>arXiv preprint arXiv:2401.07339</i> .	
	Linghao Zhang, Shilin He, Chaoyun Zhang, Yu Kang, Bowen Li, Chengxing Xie, Junhao Wang, Maoquan Wang, Yufan Huang, Shengyu Fu, and 1 others. 2025. Swe-bench goes live! <i>arXiv preprint arXiv:2505.23419</i> .	
	Terry Yue Zhuo, Minh Chien Vu, Jenny Chim, Han Hu, Wenhao Yu, Ratnadira Widayarsi, Imam Nur Bani Yusuf, Haolan Zhan, Junda He, Indraneil Paul, and 1 others. 2024. Bigcodebench: Benchmarking code generation with diverse function calls and complex instructions. <i>arXiv preprint arXiv:2406.15877</i> .	

A Prompts and Implementation Details

To ensure the reproducibility of ClarEval, we provide the detailed prompts used in our construction pipeline and evaluation framework. These details are extracted directly from our data synthesis engine and evaluation scripts.

A.0.1 Implementation Details

All experiments use standardized prompting strategies tailored to each model’s interface, with a unified instruction encouraging proactive clarification. We set temperature=0.1 for deterministic generation and limit output length to 512 tokens. Evaluation is fully automated via the ClarEval framework, which parses model responses and computes metrics against human-annotated gold standards. All experiments were conducted on a cluster of NVIDIA A100 80GB GPUs.

A.1 Ambiguity Injection Prompts

We employed GPT-4o to transform standard HumanEval tasks into ambiguous variants. Table 5 presents the specific prompts used for each ambiguity type.

A.2 Dialogue Script Construction

To enable multi-turn evaluation, we generated a “User Simulator Script” using the prompt in Table 6.

A.3 Evaluation Implementation

In our evaluation, all agents shared a unified system instruction (Table 7). The User Simulator operated deterministically based on the generated scripts.

B Efficiency-Adjusted Recall (EAR)

B.1 Metric Definition

While *Missing Premises Recall* (MPR) and *Average Turns to Clarify* (ATC) evaluate clarification completeness and efficiency separately, a unified metric is desirable to penalize agents that achieve high recall only through excessive dialogue turns. Drawing inspiration from Discounted Cumulative Gain (DCG) in information retrieval, we introduce Efficiency-Adjusted Recall (EAR).

EAR measures the proportion of missing premises successfully clarified, with each clarified premise discounted by a logarithmic factor of the turn number in which it was resolved. Formally, let \mathcal{P} be the set of all missing premises for a task, and $\mathcal{R} \subseteq \mathcal{P}$ be the set of clarified premises. For

each clarified premise $p \in \mathcal{R}$, let $Turn(p)$ denote the dialogue turn index (starting from 1) where the premise was first resolved. The EAR score is defined as:

$$EAR = \frac{1}{|\mathcal{P}|} \sum_{p \in \mathcal{R}} \frac{1}{\log_2(Turn(p) + 1)} \quad (7)$$

Where:

- A premise clarified in the **1st turn** contributes $\frac{1}{\log_2(1+1)} = 1.0$ (full value).
- A premise clarified in the **2nd turn** contributes $\frac{1}{\log_2(3)} \approx 0.63$.
- A premise clarified in the **3rd turn** contributes $\frac{1}{\log_2(4)} = 0.5$.
- Unclearified premises ($p \notin \mathcal{R}$) contribute 0.

Interpretation: EAR creates a soft penalty for delayed clarification. An agent that recovers all information immediately ($EAR = 1.0$) is superior to one that recovers all information over multiple turns ($EAR < 1.0$), which is in turn superior to an agent that fails to recover information ($EAR \approx 0$). This metric effectively addresses the trade-off between the “Lazy Agent” bias (low ATC but low MPR) and the “Inefficient Agent” bias (high MPR but high ATC).

B.2 Quantitative Analysis

Model	MPR	ATC	EAR (Approx)	Rank
GPT5-Coder	0.951	1.493	0.722	1
Claude-Opus 4.1	0.930	1.900	0.605	2
Claude-Sonnet 4.5	0.823	1.719	0.570	3
DeepSeek-Coder	0.796	1.772	0.541	4
O3	0.791	1.771	0.538	5
CodeX-GPT5	0.748	1.681	0.526	6
O4	0.712	1.576	0.522	7
Qwen3-Coder	0.640	1.645	0.456	8
Qwen2.5-Coder	0.643	2.396	0.365	9
Aider-GPT5	0.396	1.576	0.290	10
CodeLlama-70B	0.385	1.602	0.279	11

Table 3: Efficiency-Adjusted Recall (EAR) rankings. Higher EAR indicates better balance between identifying missing premises and minimizing user dialogue turns.

Table 3 presents the ranked EAR scores. This metric reveals significant performance distinctions that are less visible when viewing MPR and ATC in isolation:

- 891 • **Dominance of GPT5-Coder:** GPT5-Coder 938
892 achieves a dominant EAR of **0.722**, outper- 939
893 forming the second-best model (Claude-Opus 940
894 4.1) by a substantial margin ($\sim 19\%$). This 941
895 confirms that GPT5-Coder is not only thor- 942
896 ough but also exceptionally efficient, often 943
897 resolving ambiguities in fewer turns. 944
- 898 • **The Efficiency Penalty:** While Qwen2.5- 945
899 Coder achieves a comparable MPR (0.643) to 946
900 Qwen3-Coder (0.640), its EAR drops signifi- 947
901 cantly (0.365 vs. 0.456). This penalty reflects 948
902 Qwen2.5-Coder’s high ATC (2.396), indicat- 949
903 ing a tendency to ask fragmented questions 950
904 over many turns, which degrades the user ex- 951
905 perience. 952

906 B.3 The Efficiency-Cognitive Load Trade-off 953

907 A critical finding of our study is the mismatch 954
908 between capability and interaction efficiency in 955
909 modern code agents. While open-weights mod- 956
910 els like Qwen2.5-Coder achieve competitive intent 957
911 detection scores (KQC: 0.855), their interaction 958
912 efficiency lags significantly behind frontier mod- 959
913 els like GPT-5-Coder (ATC: 2.396 vs. 1.493). We in- 960
914 terpret this not merely as a performance gap, but as 961
915 a trade-off in clarification planning. 962

916 **Batching vs. Serial Strategies.** Our qualitative 963
917 audit reveals distinct behavioral patterns: 964

- 918 • **The Batching Strategy (GPT-5-Coder):** 965
919 Frontier models tend to employ a “Parallel” 966
920 approach. When faced with composite ambi- 967
921 guities (e.g., missing both sort order and data 968
922 type), these agents consolidate inquiries into 969
923 a single, structured turn. While this increases 970
924 the per-turn reading load for the user, it sig- 971
925 nificantly minimizes interaction latency (Low 972
926 ATC) and accelerates the time-to-solution. 973
- 927 • **The Fragmented Strategy (Qwen2.5- 974
928 Coder):** Conversely, other models often 975
929 exhibit a “Serial Chain” behavior. While 976
930 a serial approach can theoretically reduce 977
931 per-turn cognitive load, our analysis shows 978
932 that models like Qwen2.5 do not employ this 979
933 strategically. Instead, they exhibit *fragmented* 980
934 *reasoning*—asking about a data type in 981
935 Turn 1, waiting for feedback, and only then 982
936 realizing they also need the sort order in Turn 983
937 3. 984

This distinction is crucial: whereas strategic se- 938
rial questioning can be a valid user-centric choice, 939
the observed fragmented questioning indicates a 940
failure in global planning. The high ATC in these 941
models stems from an inability to anticipate de- 942
pendencies, imposing a severe “Cognitive Tax” of 943
friction on the user. Thus, ClarEval’s ATC met- 944
ric serves as a proxy for an agent’s ability to plan 945
clarifications efficiently. 946

947 C Dialogue Script Construction 948

(Multi-Turn)

949 To evaluate the agent’s ability to recover this miss- 950
951 ing information, we constructed a “User Simulator 952
953 Script” \mathcal{S}_i for each task. The script \mathcal{S}_i is derived 954
955 from the original ground truth T_i and acts as a con- 956
957 ditional lookup table for the simulator. 958

Each script is structured as a set of trigger- 959
response pairs: 960

$$961 \mathcal{S}_i = \{(I_{trigger}, R_{content})\} 962$$

963 where $I_{trigger}$ represents the intent the agent asks 964
965 about (e.g., “Input Format”, “Edge Case”), and 966
967 $R_{content}$ is the pre-defined user answer restored 968
969 from the ground truth. 970

971 For example, in a task where the *Missing* 972
973 *Premises* strategy removed the constraint “The list 974
975 must be sorted descending”, the script records this 976
977 as a conditional response: *If Agent asks about or- 978
979 der \rightarrow Reply “Sort it descending.”* This ensures 980
981 that the evaluation of multi-turn capabilities is de- 982
983 terministic and consistent across all models. 984

965 D Robustness Validation of the User 966

Simulator

967 To address concerns regarding the potential rigidity 968
969 of the rule-based matching mechanism used in our 970
971 User Simulator, we conducted a rigorous validation 972
973 experiment. Our goal was to determine if the rule- 974
975 based keywords (Triggers) sufficiently capture the 976
977 semantic intent of clarification questions generated 978
979 by advanced agents. 980

974 D.1 Methodology: Semantic Consistency 975

Check

976 We randomly sampled 200 interaction turns from 977
978 the evaluation logs of the top-performing models 979
980 (GPT-5 Coder and Claude-Opus 4.1). We then em- 981
982 ployed a “LLM-as-a-Judge” approach using GPT- 983
984 4o to semantically evaluate the agent’s questions 985

981	against the ground truth intent, independent of the	up” always means dropna. This is the most	1026
982	keyword triggers.	common failure mode in single-turn settings.	1027
983	The LLM Judge was instructed to classify each		
984	agent question into one of two categories:		
985	• Hit (Semantic): The question effectively asks	• Generic Querying: The agent recognizes	1028
986	for the missing information required by the	something is wrong but asks a non-actionable	1029
987	ground truth.	question, such as “ <i>Is there anything else you</i>	1030
988	• Miss (Semantic): The question is irrelevant,	<i>need?”</i> or “ <i>Can you provide more details?”</i>	1031
989	too vague, or hallucinates constraints.	without specifying <i>what</i> details are missing.	1032
990		These do not trigger the User Simulator’s spe-	1033
991	We compared these semantic labels with the de-	specific responses.	1034
992	terministic labels (Hit/Miss) generated by our rule-		
993	based User Simulator.	• Hallucinated Constraints: Instead of clarify-	1035
994		ing, the agent invents its own constraints (e.g.,	1036
995	D.2 Results	“ <i>I will assume the data is sorted</i> ”), effectively	1037
996	Table 8 presents the confusion matrix between	bypassing the clarification process.	1038
997	the Rule-Based Simulator and the LLM Semantic		
998	Judge.	E.2 Case Study: “Ambiguous Terminology”	1039
999	The analysis yields the following insights:	Table 9 demonstrates a concrete example of how	1040
1000	• Agreement Rate: The overall agreement be-	different agents handle the ambiguous term “Orga-	1041
1001	tween our rule-based script and the semantic	nize” in the task: “ <i>Organize the list of items based</i>	1042
1002	judge is 96.5% .	<i>on the 'id'.</i> ”	1043
1003	• False Negatives (2.5%): In only 5 out of		
1004	200 cases did the rule-based simulator fail to	E.3 Implications	1044
1005	recognize a valid question. This typically oc-	This analysis suggests that future model training	1045
1006	curred when agents used highly metaphorical	should focus on Semantic Sensitivity —teaching	1046
1007	language (e.g., asking about “data hygiene”	models to flag non-technical descriptors (like	1047
1008	instead of “null values”).	“fast”, “best”, “organize”) as triggers for mandatory	1048
1009		clarification, rather than treating them as stylistic	1049
1010	• Conclusion: The high alignment confirms	noise.	1050
1011	that our trigger keyword sets are sufficiently		
1012	diverse and robust. The rule-based approach	F Fine-Grained Analysis by Ambiguity	1051
1013	provides a computationally efficient and deter-	Type	1052
1014	ministic evaluation without significantly pe-	F.1 KQC Breakdown	1053
1015	nalizing the semantic flexibility of SOTA mod-	Appendix Table 10 and Table 11 shows KQC per-	1054
1016	els.	formance across different ambiguity types. In	1055
1017		single-turn settings, all models perform best on	1056
1018	E Qualitative Analysis of Ambiguity	Missing Goal scenarios, with GPT5-Coder achiev-	1057
1019	Resolution	ing 0.945. Performance on Ambiguous Termi-	1058
1020	To provide deeper insights into why models strug-	nology is consistently lowest across models (e.g.,	1059
1021	gle significantly with Ambiguous Terminology	GPT5-Coder: 0.918; Claude-Opus: 0.817), indicat-	1060
1022	(A_T), we conducted a qualitative analysis of failure	ing this is the most challenging ambiguity type for	1061
1023	cases. We identified three primary error modes that	intent inference.	1062
1024	hinder clarification performance.	In multi-turn settings, most models suffer sig-	1063
1025		nificant KQC degradation. Notably, Aider-GPT5	1064
		drops from 0.857 to 0.364 on Missing Goal, and	1065
		CodeLlama-70B falls from 0.854 to 0.434, reveal-	1066
		ing severe limitations in context management and	1067
		iterative intent refinement.	1068
	E.1 Taxonomy of Clarification Failures		
	• Assumptive Generation (The “Silent Fail-	F.2 MPR Breakdown	1069
	ure”): The agent ignores the ambiguity enti-	As shown in Appendix Table 12 and Table 13,	1070
	rely and generates code based on a high-	GPT5-Coder consistently achieves high MPR	1071
	probability assumption (e.g., assuming “clean		

1072 across all ambiguity types, particularly excelling
1073 in multi-turn Missing Premises scenarios (0.943).
1074 Claude-Opus 4.1 also shows robust multi-turn per-
1075 formance (0.930–0.935 across types), indicating
1076 reliable technical clarification capabilities.

1077 Most models show improved MPR in multi-turn
1078 vs. single-turn settings, validating the value of ex-
1079 tended interaction for uncovering implementation
1080 details. However, Aider-GPT5 and CodeLlama-
1081 70B show minimal improvement, suggesting their
1082 multi-turn reasoning mechanisms are underdevel-
1083 oped.

1084 **F.3 Efficiency Analysis**

1085 Appendix Table 14 reveals that GPT5-Coder
1086 achieves the most consistent efficiency across am-
1087 biguity types (ATC: 1.422–1.602), often asking
1088 multiple critical questions in a single turn. Con-
1089 versely, Qwen2.5-Coder shows the highest ATC
1090 values (2.338–2.486), reflecting a sequential, less
1091 efficient questioning strategy. Interestingly, ATC
1092 patterns are relatively stable across ambiguity types
1093 for most models, suggesting that questioning effi-
1094 ciency is a model-intrinsic trait, independent of
1095 ambiguity nature.

1096 **G The Value of ClarEval**

1097 ClarEval successfully moves beyond functional cor-
1098 rectness, providing a standardized framework to
1099 assess and, ultimately, improve the collaborative
1100 intelligence of code agents. By quantifying capabil-
1101 ities across intent inference, technical clarification,
1102 and interaction efficiency, our benchmark enables
1103 targeted improvements in agent design. This shift
1104 from correctness-only to collaboration-aware eval-
1105 uation is essential for developing the next gener-
1106 ation of usable, trustworthy, and human-aligned
1107 programming assistants.

Ambiguity Type	Injection Prompt Strategy	Transformation Example
Missing Goal (A_G)	Strategy: Logic Omission The prompt simulates a user who provides data context but forgets the objective. <ul style="list-style-type: none"> Identify core algorithmic goal ("what to do"). Remove the specific logic entirely. Retain input context descriptions. 	<i>Original:</i> "Given a list of integers, <u>return the sum of all even numbers.</u> " ↓ <i>Ambiguous:</i> "I have a list of integers. Write a function to <u>process them.</u> "
Missing Premises (A_P)	Strategy: Constraint Removal The prompt simulates a user who has a goal but omits critical boundary conditions. <ul style="list-style-type: none"> Keep the high-level goal. Identify and remove specific constraints (e.g., "shortest", "starts with"). 	<i>Original:</i> "Find the <u>shortest</u> palindrome that <u>starts with</u> the given string." ↓ <i>Ambiguous:</i> "Write a function that finds a <u>palindrome</u> using the given string."
Ambiguous Terms (A_T)	Strategy: Vague Replacement The prompt simulates a non-technical user employing colloquial language. <ul style="list-style-type: none"> Identify precise technical terms (e.g., "sort", "dict"). Replace with vague synonyms (e.g., "organize", "structure") without removing context. 	<i>Original:</i> " <u>Sort</u> the list of <u>dictionaries</u> by the 'id' key." ↓ <i>Ambiguous:</i> " <u>Organize</u> the list of <u>items</u> based on the 'id'."

Table 4: The Ambiguity Injection strategies used to construct ClarEval. We apply distinct transformation rules—Omission (A_G), Constraint Removal (A_P), and Vague Replacement (A_T)—to systematically generate diverse types of underspecification from the original HumanEval tasks.

Type	Injection Prompt Template
System	You are a software requirement engineer. Your task is to rewrite clear coding problems into ambiguous or underspecified versions to test an AI assistant’s clarification capabilities.
Missing Goal	Task: Rewrite the following problem description to introduce ambiguity by removing the explicit goal . Instructions: <ol style="list-style-type: none"> Keep the context (input data description, setup). Remove the specific instruction of what to calculate or return. Replace the goal with a vague request (e.g., "process the data", "handle the input"). Ensure the code cannot be written correctly without asking "What should I do with this data?".
Missing Premises	Task: Rewrite the following problem description to introduce ambiguity by omitting necessary constraints . Instructions: <ol style="list-style-type: none"> Keep the high-level goal clear (e.g., "sort the list"). Remove specific constraints required for a unique solution (e.g., sort order, handling of duplicates, edge cases like empty lists). The resulting description should look plausible but lead to multiple possible valid implementations.
Ambiguous Terms	Task: Rewrite the following problem description to introduce ambiguity by using vague terminology . Instructions: <ol style="list-style-type: none"> Identify precise technical terms in the original description (e.g., "dictionary", "prime number", "recursion"). Replace them with non-technical, colloquial, or ambiguous synonyms (e.g., "lookup list", "special number", "repeated logic"). Do not remove the logic, but make the <i>language</i> imprecise.

Table 5: The prompt templates used to inject three types of ambiguity into coding tasks.

Component	Content
Script Gen. Prompt	<p>Task: Compare the Original Task and the Ambiguous Task below. Identify the missing information that an agent would need to ask to recover the Original Task’s requirements.</p> <p>Input: Original Task: {original_prompt} Ambiguous Task: {ambiguous_prompt}</p> <p>Output Format (JSON): Generate a list of Key-Value pairs where:</p> <ul style="list-style-type: none"> - Trigger: A short keyword or phrase representing the missing intent (e.g., "Sort Order"). - Response: The specific detail from the Original Task (e.g., "Sort in descending order").

Table 6: The prompt used to generate the User Simulator’s lookup table (Ground Truth intents).

Component	Implementation Logic / Prompt
Agent System Prompt	<p>You are a helpful and professional coding assistant. When you receive a coding task from the user, please follow these steps:</p> <ol style="list-style-type: none"> 1. Carefully analyze the user’s instruction. 2. If the instruction is ambiguous, underspecified, or lacks critical details, you should ask clarifying questions to the user to resolve the ambiguity. 3. Do not guess the user’s intent if it is unclear. 4. Once you have sufficient information, generate the code to solve the problem.
User Simulator	<p>Mechanism: Deterministic, Rule-Based Simulation using Pre-defined Scripts.</p> <p>Role Setup: The simulator follows a pre-defined script S_i containing trigger-response pairs derived from the ground truth T_{clear}.</p> <p>Interaction Logic:</p> <ol style="list-style-type: none"> 1. Input Matching: The agent’s question is matched against pre-defined trigger keywords/phrases in S_i. 2. Response Retrieval: <ul style="list-style-type: none"> – If a trigger match is found → Return the corresponding pre-written response from the script. – If no match is found → Return fallback: "I don’t have specific requirements for that; please follow standard best practices." <p>Termination: Upon code generation or reaching the maximum number of turns ($T = 5$).</p>

Table 7: Details of the Agent System Prompt and the deterministic, rule-based User Simulator logic used in `eval.py`. The simulator operates based on pre-defined scripts to ensure reproducibility and fairness across all model evaluations.

	LLM Judge: Hit	LLM Judge: Miss
Rule-Based: Hit	84 (True Positive)	2 (False Positive)
Rule-Based: Miss	5 (False Negative)	109 (True Negative)

Table 8: Confusion Matrix of Simulator Decisions vs. LLM Semantic Judge. The high agreement rate confirms the reliability of the rule-based mechanism.

Model	Response Strategy	Outcome	Analysis
GPT-5 Coder (Strong)	Clarification Question: “By ‘organize’, do you mean sorting the list? If so, should it be in ascending or descending order?”	Success (Hit)	The agent correctly mapped the vague term “organize” to the technical concept “sort” and identified the missing parameter (order).
Qwen2.5-Coder (Inefficient)	Clarification Question: “What is the structure of the items?” (Next Turn): “How should I handle the ID?”	Partial Success (High ATC)	The agent asked tangential questions first. It took 3 turns to finally ask about the sorting order, resulting in a high ATC score.
CodeLlama-70B (Weak)	Direct Action: “Here is the python code to sort the list of dictionaries by ID in ascending order...”	Failure (Miss)	The model exhibited Assumptive Generation . It failed to recognize “organize” as ambiguous and arbitrarily chose “ascending sort,” leading to potential misalignment with user intent.

Table 9: Comparison of Agent Responses to Ambiguous Terminology. While GPT-5 Coder correctly identifies the need to clarify the sort order, weaker models either ask inefficient questions or make unsafe assumptions.

Model	Missing Goal	Missing Premises	Ambiguous Terminology
Aider-GPT5	0.600	0.593	0.571
CodeX-GPT5	0.551	0.585	0.536
GPT5-Coder	0.830	0.661	0.810
CodeLlama-70B	0.539	0.555	0.512
Claude-opus4.1	0.607	0.650	0.631
Claude-sonnet-4.5	0.625	0.636	0.633
O3	0.572	0.588	0.556
O4	0.536	0.554	0.542
DeepSeek-Coder	0.586	0.589	0.588
Qwen2.5-Coder	0.615	0.618	0.633
Qwen3-Coder	0.590	0.621	0.565

Table 10: Fine-grained results for KQC in the Single-Turn setting. We report the agent’s performance across three distinct ambiguity types: Missing Goal, Missing Premises, and Ambiguous Terminology. (Note: Higher values indicate a stronger capability to infer the correct user intent immediately from the initial underspecified instruction.)

Model	Missing Goal	Missing Premises	Ambiguous Terminology
Aider-GPT5	0.364	0.373	0.392
CodeX-GPT5	0.639	0.623	0.615
GPT5-Coder	0.505	0.554	0.529
CodeLlama-70B	0.434	0.389	0.459
Claude-opus4.1	0.753	0.768	0.742
Claude-sonnet-4.5	0.618	0.649	0.631
O3	0.569	0.586	0.562
O4	0.535	0.562	0.517
DeepSeek-Coder	0.651	0.702	0.671
Qwen2.5-Coder	0.531	0.549	0.551
Qwen3-Coder	0.681	0.683	0.639

Table 11: Fine-grained results for KQC in the Multi-Turn setting. We report the agent’s performance across three distinct ambiguity types: Missing Goal, Missing Premises, and Ambiguous Terminology. (Note: Higher values indicate a stronger capability to infer the correct user intent immediately from the initial underspecified instruction.)

Model	Missing Goal	Missing Premises	Ambiguous Terminology
Aider-GPT5	0.600	0.593	0.571
CodeX-GPT5	0.551	0.585	0.536
GPT5-Coder	0.830	0.661	0.810
CodeLlama-70B	0.539	0.555	0.512
Claude-opus4.1	0.607	0.650	0.631
Claude-sonnet-4.5	0.625	0.636	0.633
O3	0.572	0.588	0.556
O4	0.536	0.554	0.542
DeepSeek-Coder	0.586	0.589	0.588
Qwen2.5-Coder	0.615	0.618	0.633
Qwen3-Coder	0.590	0.621	0.565

Table 12: Fine-grained results for PIR in the Single-Turn setting. We report the agent’s performance across three distinct ambiguity types: Missing Goal, Missing Premises, and Ambiguous Terminology. (Note: Higher values indicate a stronger capability to infer the correct user intent immediately from the initial underspecified instruction.)

Model	Missing Goal	Missing Premises	Ambiguous Terminology
Aider-GPT5	0.393	0.389	0.405
CodeX-GPT5	0.729	0.747	0.768
GPT5-Coder	0.954	0.943	0.956
CodeLlama-70B	0.367	0.360	0.430
Claude-opus4.1	0.930	0.924	0.867
Claude-sonnet-4.5	0.809	0.822	0.837
O3	0.822	0.777	0.773
O4	0.715	0.718	0.703
DeepSeek-Coder	0.799	0.803	0.785
Qwen2.5-Coder	0.657	0.623	0.649
Qwen3-Coder	0.644	0.635	0.641

Table 13: Fine-grained results for MPR in the Multi-Turn setting. We report the agent’s performance across three distinct ambiguity types: Missing Goal, Missing Premises, and Ambiguous Terminology. (Note: Higher values indicate a stronger capability to infer the correct user intent immediately from the initial underspecified instruction.)

Model	Missing Goal	Missing Premises	Ambiguous Terminology
Aider-GPT5	1.603	1.548	1.577
CodeX-GPT5	1.625	1.766	1.651
GPT5-Coder	1.422	1.455	1.602
CodeLlama-70B	1.533	1.466	1.807
Claude-opus4.1	1.898	1.869	1.933
Claude-sonnet-4.5	1.539	1.885	1.732
O3	1.840	1.701	1.772
O4	1.458	1.629	1.640
DeepSeek-Coder	1.599	1.849	1.868
Qwen2.5-Coder	2.363	2.486	2.338
Qwen3-Coder	1.689	1.599	1.646

Table 14: Analysis of interaction efficiency using ATC across the three ambiguity types. The metric measures the average number of dialogue turns required for an agent to successfully recover a missing premise. (Note: Lower values for this metric indicate better efficiency, signifying that the agent secures necessary clarifications with fewer user interactions.)