

Context as a Tool: Context Management for Long-Horizon SWE-Agents

Anonymous ACL submission

Abstract

Agents based on large language models have recently shown strong potential on real-world software engineering (SWE) tasks that require long-horizon interaction with repository-scale codebases. However, most existing agents rely on append-only context maintenance or passively triggered compression heuristics, which often lead to context explosion, semantic drift, and degraded reasoning in long-running interactions. We propose CAT, a new context management paradigm that elevates context maintenance to a callable tool integrated into the decision-making process of agents. CAT formalizes a structured context workspace consisting of stable task semantics, condensed long-term memory, and high-fidelity short-term interactions, and enables agents to proactively compress historical trajectories into actionable summaries at appropriate milestones. To support context management for SWE-agents, we propose a trajectory-level supervision framework, CAT-GENERATOR, based on an offline data construction pipeline that injects context-management actions into complete interaction trajectories. Using this framework, we train a context-aware model, SWE-Compressor. Experiments on SWE-Bench-Verified demonstrate that SWE-Compressor reaches a 57.6% solved rate and significantly outperforms ReAct-based agents and static compression baselines, while maintaining stable and scalable long-horizon reasoning under a bounded context budget.

1 Introduction

Large language models (LLMs) (Yang et al., 2025a; Touvron et al., 2023; Anthropic, 2025; Achiam et al., 2023) have achieved remarkable progress on tasks such as code generation and bug fixing (Chen, 2021; Austin et al., 2021; Liu et al., 2024b; Chai et al., 2024). As research attention increasingly shifts toward real-world software engineering (SWE) scenarios, maintaining stable and

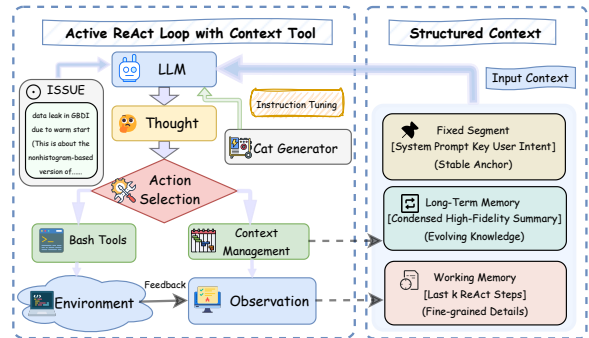


Figure 1: Overview of CAT with a structured context workspace for long-horizon reasoning.

effective reasoning in complex and long-horizon interactive tasks has emerged as a primary challenge in agent research. Most SWE tasks (e.g., repository-level issue resolution (Jimenez et al., 2023)) often require LLMs to continuously interpret environment feedback, execute actions, and revise strategies over hundreds of interaction rounds, which poses greater demands on context management.

Most existing code agents relying on paradigms such as ReAct (Yao et al., 2022), which adopts an append-only context maintenance strategy, where all past interactions are continuously concatenated into the messages. While this method may suffice for short-horizon tasks, it often leads to rapid context expansion in long-horizon scenarios, resulting in information redundancy, semantic drift, and even reasoning collapse. To mitigate these challenges, prior work (Jiang et al., 2023; Shinn et al., 2023; Wang et al., 2025b; Packer et al., 2023) has explored context compression, summarization, or multi-level memory mechanisms to constrain context size. Nevertheless, previous approaches treat context management as a passively triggered heuristic mechanism, so they lack the flexibility to adapt compression timing and content across different task phases, limiting their adaptability and scalability in complex environments. In this paper, we propose a new context management paradigm, CAT. We argue that in long-horizon interactive

tasks, context management should be internalized as a model capability rather than enforced through external constraints. Motivated by this perspective, CAT treats context management as a callable and plannable tool, on par with environment-interaction tools such as code editing and command execution, which integrates context maintenance into the action process as an active and learnable component.

In Figure 1, CAT organizes the context into a structured workspace consisting of stable task-semantic anchors, an evolvable long-term memory, and a short-term working memory. It further enables the agent to proactively trigger context folding at stage boundaries, compressing redundant histories into high-fidelity and actionable long-term memory representations. Building on this design, we introduce a trajectory-level supervision framework, CAT-GENERATOR, which injects context-management behaviors into complete interaction trajectories via offline reconstruction. Using CAT-GENERATOR, we train a context-aware model, SWE-Compressor, allowing it to learn when to compress context, how to generate effective summaries, and how to reuse compressed representations during subsequent reasoning.

Experimental results on SWE-Bench show that CAT substantially outperforms ReAct without context management and baselines with static compression strategies under the same model scale and interaction budget, demonstrating the strong context scalability in long-horizon interaction settings. The contribution of his paper is summarized as:

- We propose CAT, a new context management paradigm that internalizes context maintenance as a learnable, tool-based capability for long-horizon interactive reasoning.
- We design a structured context workspace with proactive context folding, enabling effective memory construction and sustained reasoning under bounded context budgets.
- We introduce CAT-GENERATOR, a trajectory-level supervision framework for learning context-management behaviors, and train a context-aware model, SWE-Compressor, that effectively compresses and reuses context during extended interactions.

2 SWE-Compressor

Unlike prior agents treating context management as a passive heuristic or post-processing step, CAT

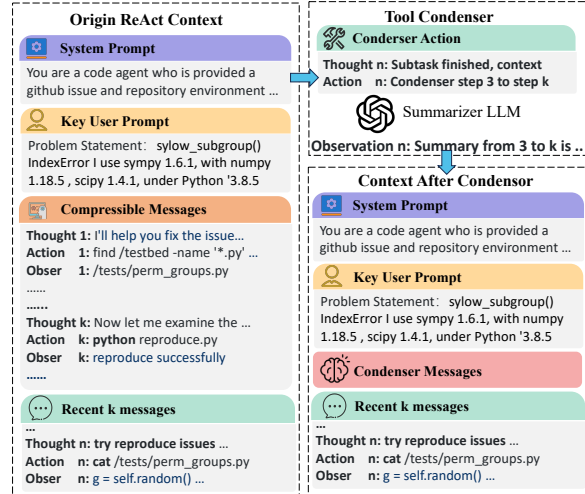


Figure 2: Example of structured context condensation in CAT during long-horizon tasks.

elevates it to a callable and plannable capability on par with environment-interaction tools, making context maintenance an active and learnable component of the agent’s decision-making process. Specifically, we formalize a structured context workspace, introduce tool-based context management with segmented compression for ReAct, and present a trajectory-level supervision framework (Figure 3) with an offline data pipeline that enables effective context compression and reuse in subsequent reasoning and decision-making.

2.1 Structured Context Workspace

When performing long-horizon ReAct reasoning in complex environments, agent performance largely depends on how its working context is organized. Motivated by this observation, CAT models context as a controllable and dynamically updated cognitive workspace composed of three functional segments, as illustrated in Figure 2. The first is a fixed segment that preserves the system prompt and key user intent. The second is a long-term memory segment that stores a condensed, high-fidelity summary of historical trajectories. The third is a high-fidelity working memory segment that retains the most recent k ReAct interaction steps. This design provides a stable semantic anchor for the task while preserving fine-grained information from recent environment feedback, thereby supporting precise contextualized actions. Formally, the working context at step t is represented as

$$C(t) = (Q, M(t), I^{(k)}(t)) \quad (1)$$

where Q denotes the non-compressible component consisting of the system prompt and key user objec-

154 tives, $M(t)$ represents the high-fidelity summary
 155 of historical trajectories, and $I^{(k)}(t)$ denotes the
 156 complete records of the most recent k ReAct inter-
 157 actions. At initialization, $C(1) = (Q, \emptyset, \emptyset)$. As rea-
 158 soning progresses, the agent continuously updates
 159 recent interactions and long-term summaries, and
 160 adjusts the content and structure of $M(t)$ through
 161 the context-management tool when needed. In this
 162 way, stable goals, condensed knowledge, and fine-
 163 grained working memory are coordinated within a
 164 unified framework, mitigating both semantic drift
 165 and uncontrolled context expansion.

166 2.2 Context Management as a First-Class Tool

167 The key innovation of CAT is to explicitly model
 168 context management as a callable tool operation
 169 and to place it at the same decision level as
 170 environment-interaction tools such as file editing
 171 or command execution. Consequently, when gen-
 172 erating a response at step t , the agent evaluates
 173 whether to invoke the context-management tool in
 174 the same manner as it selects other external actions.
 175 In practice, the agent tends to proactively trigger
 176 context management when a subtask has been com-
 177 pleted and requires a stage-wise summary, when
 178 the trajectory has grown sufficiently large that his-
 179 torical compression is necessary to maintain oper-
 180 ational efficiency, or when subsequent reasoning
 181 benefits more from a concise, structured summary
 182 than from verbose raw logs. Through this design,
 183 context management is transformed from a post
 184 hoc procedure into a self-regulating reasoning step,
 185 becoming an integral part of the agent’s strategy
 186 rather than an external constraint.

187 2.3 Structured Memory Generation

188 For the compressible historical segment, CAT ap-
 189 plies structured summarization to condense infor-
 190 mation accumulated during long-horizon reasoning
 191 into a compact and actionable long-term memory.
 192 The goal is to reduce context length while pre-
 193 serving information that remains causally relevant
 194 for subsequent decisions. The resulting long-term
 195 memory summarizes key aspects of task progress,
 196 including intermediate goals, adopted strategies
 197 and their outcomes, salient environment feedback,
 198 and persistent constraints that continue to shape fu-
 199 ture reasoning. After summarization, the agent
 200 reconstructs its working context by combining
 201 the fixed task semantics, the condensed long-term
 202 memory, and the most recent high-fidelity interac-
 203 tion steps, enabling stable and consistent reasoning

under a bounded context budget. 204

205 2.4 Supervised Trajectory Construction

Data Generation Pipeline To internalize 206
 207 toolized context management capability in LLM,
 208 we propose CAT-GENERATOR, a two-stage
 209 retrospective trajectory construction pipeline,
 210 as illustrated in Figure 3. CAT-GENERATOR
 211 first generates complete base ReAct reasoning
 212 trajectories without introducing any context com-
 213 pression operations, to maximize the naturalness
 214 and completeness of task-solving behaviors. Then,
 215 these trajectories are minimally and structurally re-
 216 constructed: context management tool invocations
 217 are injected at appropriate steps without altering
 218 the original sequence of environment interactions,
 219 yielding SFT data that are consistent with the
 220 reasoning paradigm of CAT.

221 Phase I: Base ReAct Trajectory Generation.

222 Given a task instance, we first deploy a standard Re-
 223 Act agent to execute a complete interaction process
 224 in a controlled environment, while explicitly dis-
 225 abling the context management tool and allowing
 226 only environment-related actions (e.g., code edit-
 227 ing, command execution, or information retrieval).
 228 The resulting raw trajectory is denoted as

$$229 \mathcal{T}_{\text{base}} = \{(C_{\text{base}}(1), R_{\text{base}}(1)), \dots, (C_{\text{base}}(T), R_{\text{base}}(T))\} \quad (2)$$

230 where $R_{\text{base}}(t)$ follows the standard ReAct re-
 231 sponse structure, consisting of Thought and Ac-
 232 tion. The objective of this stage is to preserve as
 233 much intermediate reasoning, failure patterns, and
 234 environment feedback as possible, thereby provid-
 235 ing sufficient causal evidence for learning context
 236 compression decisions in later stages.

237 Phase II: Trajectory Refactoring by injecting

Compression Operation. We perform an offline 238
 239 reconstruction of the base trajectory $\mathcal{T}_{\text{base}}$ to obtain
 240 an augmented trajectory $\mathcal{T}_{\text{retro}}$ that includes con-
 241 text management tool invocations. This stage is
 242 to transform context compression from an implicit
 243 side effect during generation into an explicitly mod-
 244 eled, controllably injected tool operation.

(1) Condensor Position Generation. We begin 245
 246 by conducting a structured analysis of the base tra-
 247 jectory to identify a set of candidate time steps sui-
 248 table for inserting context management tool calls,
 249 denoted as $\mathcal{A} = \{a_1, \dots, a_m\}$. The selection of
 250 insertion position jointly considers multiple signals,
 251 including: (i) context expansion signals, such as

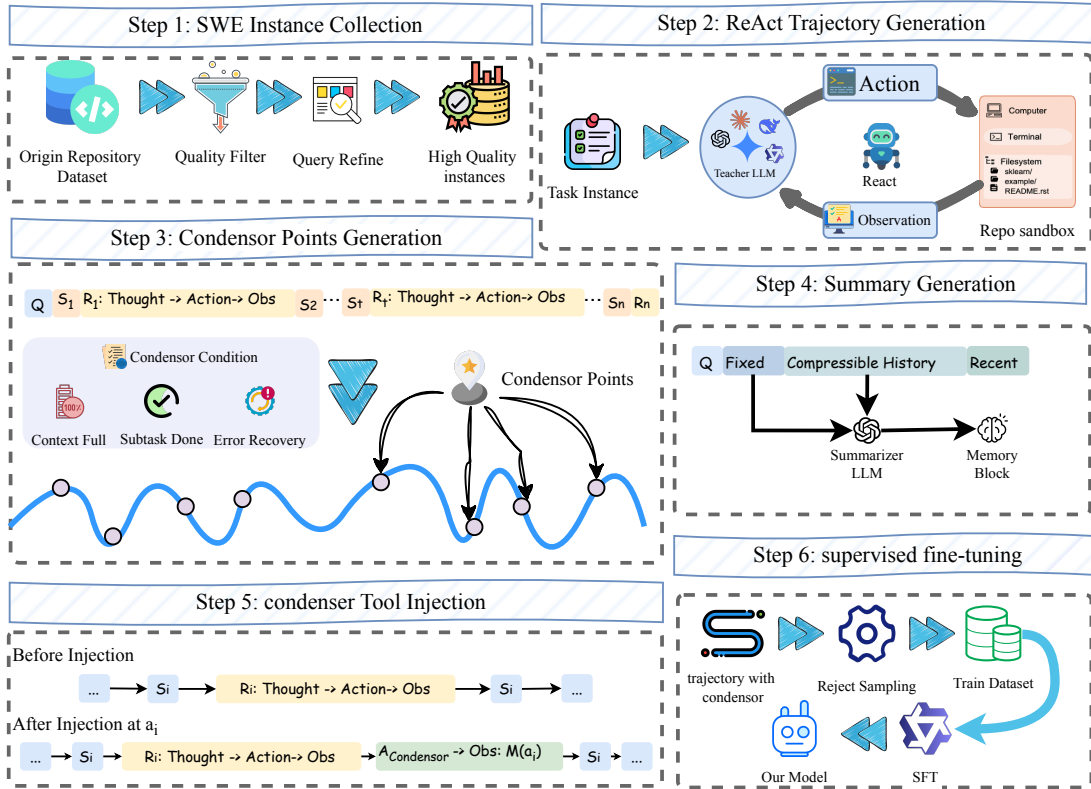


Figure 3: Overview of the data construction and training pipeline for CAT. The process includes SWE instance collection, base ReAct trajectory generation, condenser point identification, structured summary generation, tool-based context injection, and supervised fine-tuning with rejection sampling.

Statistic	CAT-Instruct
<i>Trajectory length (steps)</i>	
Average	87.4
Median	77.5
Max	500
<i>Context size (tokens)</i>	
Avg tokens per step	13,044
Median tokens per step	10,875
Max tokens per step	65,536
<i>Context-mgmt actions / traj</i>	
Average	4.22
Median	4.00
Max	26
<i>Compression</i>	
Avg tokens before	15,585
Avg tokens after	4,676
Avg ratio (%)	30

Table 1: Statistics of SFT data for CAT-Instruct.

sustained growth in context length or decreasing token utilization efficiency; (ii) structural boundary signals, such as subtask completion, strategy switching, or intermediate milestones; and (iii) error-correction signals, where new feasible directions emerge after repeated failures. We treat these signals as heuristic triggers rather than fixed thresholds, aligning insertion points with natural moments for stage-wise summarization in long-

horizon reasoning and improving the learnability of compression behaviors.

(2) Segmented Context Construction and Compression Input Preparation. For each candidate insertion position a_i , we construct a segmented representation of the context at that step by partitioning the visible context into a fixed segment Q , a recent high-fidelity working memory segment $I^{(k)}(a_i)$, and a compressible historical trajectory segment. The fixed segment and the recent trajectory segment are preserved verbatim, while the historical segment is provided as input to model for long-term memory summarization. This procedure enforces the structured context workspace constraints of CAT during the offline stage.

(3) Long-Term Memory Block Generation. Based on the segmented context, we invoke a high-capacity language model to generate a structured long-term memory block $M(a_i)$. The summarizer uses the same backbone as the reasoning model (SWE-Compressor), keeping summarization aligned with the agent’s internal reasoning style. This memory block aims to faithfully summarize critical information from the compressible history,

including completed subtasks, attempted strategies and their outcomes, important environment state changes, facts that continue to constrain subsequent decisions, and key information that remains useful for future steps. The generated $M(a_i)$ serves as the Observation of a context management tool invocation and is written into the long-term memory segment for subsequent reasoning.

(4) Trajectory Stitching and Minimal-Intrusion Compression Injection. Finally, we adopt a minimal-intrusion trajectory stitching strategy to inject context management behaviors into the original ReAct trajectory. Specifically, for each insertion point a_i , we explicitly insert a context management tool invocation at that step as an independent Action, whose corresponding Observation is the generated long-term memory block $M(a_i)$. $\mathcal{T} = \{(C(1), R(1)), (C(2), R(2)), \dots, (C(T), R(T))\}$ where T denotes the total number of steps required to complete the task. Each trajectory fully captures the process from the initial task specification, through multiple rounds of environment interaction and stage-wise context compression, to final task completion. The advantage of trajectory-level supervision lies in preserving the temporal continuity of context evolution, allowing the model not only to observe the immediate summary produced by a context-management tool invocation but also to learn how that summary influences reasoning and decision-making across subsequent steps. Moreover, it enables the model to learn cross-step strategic judgments, such as how different invocation timings affect downstream efficiency and stability, thereby better reflecting the decision structure of real long-horizon tasks.

(5) Rejection Sampling Fine-Tuning. To construct high-quality trajectory-level SFT data, we apply a rejection sampling strategy during data curation to filter interaction trajectories using trajectory-level and step-level criteria. At the trajectory level, samples that fail to complete the task or enter unrecoverable error states are discarded. At the step level, we further remove trajectories exhibiting unreasonable context-management behaviors, such as excessively frequent tool invocations with minimal information gain, severe semantic drift, or internal state inconsistencies. The resulting curated set of trajectories constitutes our SFT dataset (CAT-Instruct). Table 1 summarizes key statistics of CAT-Instruct, including trajectory length, context size, and the frequency and effec-

tiveness of context-management actions. Using CAT-Instruct for supervised fine-tuning, we obtain SWE-Compressor, which internalizes context management as a learned model capability.

3 Experiments

Datasets. We evaluate the proposed method on the SWE-Bench-Verified (Hou et al., 2024) subset. SWE-Bench is a benchmark designed to assess the ability of LLMs to solve real-world software engineering tasks collected from 12 real-world GitHub repositories. The Verified split is a high-quality subset of SWE-Bench consisting of 500 instances, which are manually curated to provide clearer problem descriptions and more reliable evaluation criteria. We report solved rate as the primary evaluation metric, defined as the proportion of instances that are successfully resolved.

Training Data Construction. For training, we collect a large number of instances from two open-source datasets, SWE-smith (Yang et al., 2025b) and SWE-ReBench (Badertdinov et al., 2025). We first employ CAT-GENERATOR to automatically generate agent interaction trajectories and apply a rejection sampling strategy to filter high-quality samples. This process yields a curated set of 20k supervised fine-tuning instances, referred to as CAT-Instruct, which effectively enhance the model’s context-management capability. In addition, following the data construction protocol of SWE-smith, we collect an additional 20k high-quality supervised fine-tuning instances, denoted as BASE-INSTRUCT, which do not involve context-management skills. These data are used to train baseline models, ensuring fair and comparable evaluation against the proposed method.

Agent Post-training. We adopt Qwen2.5-Coder-32B (Hui et al., 2024) as the base model and perform post-training on the CAT-Instruct dataset to obtain the final model, SWE-Compressor. The model is trained for up to three epochs using the AdamW (Loshchilov and Hutter, 2017) optimizer with a weight decay of 0.01. We employ a cosine learning rate schedule with a warm-up ratio of 0.1 and a peak learning rate of 5×10^{-5} . During inference, we use the OpenHands (Wang et al., 2024) framework, where the agent can invoke tools including execute_bash, str_replace_editor, submit, and context. For all experiments, the temperature is fixed to 0.0. The model is trained

Method		SWE-Bench Verified	
Model	Model Size	Scaffold	Pass@1
ReAct Agent with 100B+ LLM			
GPT-5.1 (OpenAI, 2025)	🔒	OpenHands	76.3
GPT-4o (OpenAI, 2023)	🔒	Agentless	38.8
Claude-3.5-Sonnet (Anthropic, 2024)	🔒	OpenHands	53.0
Claude-4.5-Sonnet (Anthropic, 2025)	🔒	OpenHands	77.2
Gemini-2.5-Pro (Google Cloud, 2025)	🔒	OpenHands	59.6
Gemini-3-Pro (Google DeepMind, 2025)	🔒	OpenHands	76.2
ReAct Agent			
R2E-Gym-32B (Jain et al., 2025)	32B	OpenHands	34.4
SWE-Gym-32B (Pan et al., 2024a)	32B	OpenHands	20.6
SWE-agent-LM-32B (Yang et al., 2025b)	32B	SWE-agent	40.2
DeepSWE-32B-Preview (Luo et al., 2025)	32B	OpenHands	42.2
SWE-Mirror-LM-32B (Wang et al., 2025a)	32B	OpenHands	52.2
FrogBoss-32B (Sonwane et al., 2025)	32B	OpenHands	54.6
Seed-OSS-36B (Team, 2025)	36B	OpenHands	55.2
Llama3-SWE-RL-70B (Wei et al., 2025)	70B	OpenHands	41.0
Lingma-SWE-GPT-72B (Ma et al., 2024)	72B	SWE-SynInfer	28.8
SWE-Fixer-72B (Xie et al., 2025)	72B	SWE-Fixer	32.8
GLM-4.5-Air	12/106B	OpenHands	57.6
Qwen3-235B-A22B (Yang et al., 2025a)	22/235B	OpenHands	34.4
Qwen3-Coder-480B-A35B (Yang et al., 2025a)	35/480B	OpenHands	69.6
DeepSeek-V3.1 (Liu et al., 2024a)	37/671B	OpenHands	61.0
DeepSeek-R1-0528 (Guo et al., 2025a)	37/671B	OpenHands	45.6
Summary Agent			
ReAct Agent	32B	OpenHands	49.8
Threshold-Compression Agent	32B	OpenHands	53.8
CAT			
SWE-Compressor	32B	OpenHands	57.6

Table 2: Performance comparison on SWE-Bench Verified (N=500). We report Pass@1 results for different agent systems under a unified evaluation setting, grouped by model scale and agent framework.

with a context length of 65,536 tokens; for the evaluation reported in Table 2, we allow the agent to perform up to 500 interaction rounds.

Baselines. We compare the proposed method with the following baselines: (1) ReAct (Yao et al., 2022): This baseline follows the ReAct framework and does not employ any explicit context management. Once the context window is exhausted, the dialogue terminates early. (2) Threshold-Compression (OpenHands, 2025): This agent applies context compression only when the context length exceeds a predefined threshold. Upon triggering, it follows the same compression scheme as CAT: the system prompt and key user intent are preserved verbatim, together with the most recent k interaction messages, while all remaining earlier messages are summarized into a compact representation. For all baselines, we use the same base model as SWE-Compressor and use the same summarizer backbone for any compression operation. SFT is performed on the BASE-INSTRUCT dataset, which consists of 20k instances without context-management capabilities. Besides, we also compare our method with existing closed-

source and open-source systems, such as GPT-5 and DeepSeek-R1.

Main Results. Table 2 presents the main experimental results, demonstrating the effectiveness of CAT. On the challenging SWE-Bench-Verified benchmark, SWE-Compressor achieves a 57.6% solved rate, reaching state-of-the-art performance under the setting of agent post-training on a 32B model. Under the same fine-tuning data budget, SWE-Compressor significantly outperforms both the ReAct Agent and the Threshold-Compression Agent baselines. Moreover, its performance is comparable to that of substantially larger models, and in some settings even surpasses them, under the same agent framework. These results indicate that CAT and CAT-GENERATOR are particularly effective for long-horizon interactive software engineering tasks such as those in SWE-Bench.

4 Further Analysis

Token Usage Analysis To evaluate the context management capability of CAT, we analyze 500 interaction trajectories from SWE-Bench. Specifically, we report the number of surviving trajectories

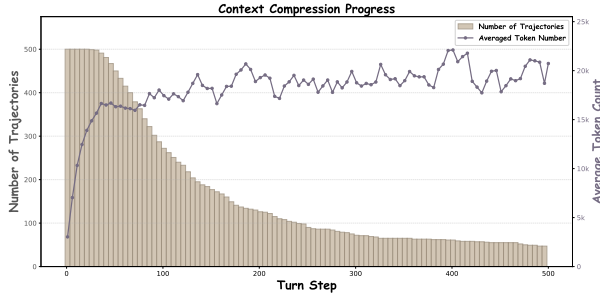


Figure 4: Context token usage and trajectory survival of CAT over interaction rounds on SWE-Bench-Verified.

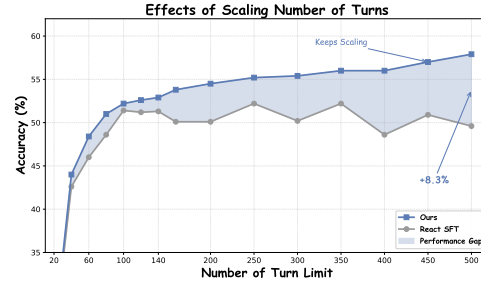
at each interaction round ($|\mathcal{T}(t)|$) and the average number of context tokens over the same set of trajectories at that round ($A(t)$). The average context token count $A(t)$ is formally defined as follows:

$$A(t) = \frac{1}{|\mathcal{T}(t)|} \sum_{j \in \mathcal{T}(t)} \text{TokenCount}(C_t^{(j)}) \quad (3)$$

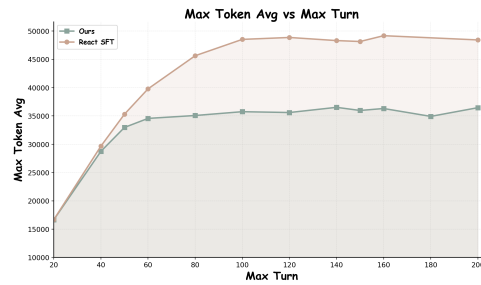
where $\mathcal{T}(t)$ denotes the set of surviving trajectories whose lengths exceed t interaction rounds, and $C_t^{(j)}$ represents the maintained context of trajectory j at round t . In Figure 4, CAT maintains a highly compact context. As the interaction progresses, the average token count quickly stabilizes after approximately 100 rounds and remains below 32k tokens, without exhibiting continuous growth over time, indicating the effectiveness of CAT in preventing context explosion. Moreover, the trajectory survival curve indicates that, in our experiments, more than 40% of tasks remain interactive after 100 rounds. This observation further suggests that CAT equips the model with stable and extensible long-horizon interaction capabilities, highlighting its strong potential for addressing highly complex and long-running software engineering tasks.

Context Comparison between CAT and ReAct. Figure 5 evaluates CAT and ReAct on SWE-Bench under identical maximum interaction round budgets. Two methods exhibit different behaviors across varying interaction budgets. Across all comparable interaction budgets, the model equipped with CAT consistently outperforms the SFT-based ReAct baseline, indicating more efficient utilization of historical information under the same reasoning budget. As the interaction budget increases, ReAct performance saturates after around 60 rounds and subsequently degrades, primarily due to its append-only context strategy: once the context window is filled, additional interactions fail to introduce effective information and instead

hinder further reasoning. In contrast, CAT continues to improve steadily with increasing interaction budgets, maintaining a clear upward trend even at 500 rounds, suggesting that CAT can continuously integrate salient information within a bounded context budget while compressing redundant history. In Figure 5, the context token usage of CAT remains stable at approximately 35k tokens, whereas the ReAct baseline rapidly exhausts the available context window.



(a) Performance of varying interaction budgets.



(b) Context token usage comparison.

Figure 5: Comparison of scalability and efficiency between CAT and ReAct on SWE-Bench. (a) CAT exhibits an upward trend in performance as the interaction budget increases to 500 rounds, whereas ReAct saturates and degrades after 60 rounds. (b) CAT maintains stable context usage (35k tokens) via condensation, while ReAct rapidly exhausts the context window.

Performance by Task Difficulty. We partition SWE-Bench into difficulty levels based on the original dataset’s reported human solution time. Specifically, instances are categorized as easy (≤ 15 minutes, 194 instances), medium (15 minutes–1 hour, 261 instances), and hard (≥ 1 hour, 45 instances). Figure 6 presents agent performance stratified by task difficulty, comparing the scores obtained under different context management strategies, showing that CAT delivers stable and consistent performance improvements across easy, medium, and hard instances. The performance gains are substantially larger on the medium and hard subsets than on the easy subset. This observation suggests that when tasks require more complex reasoning processes

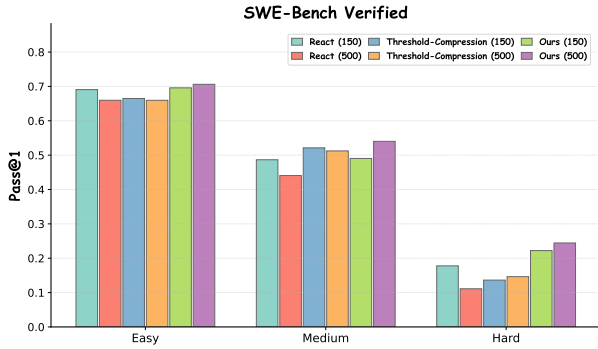


Figure 6: Performance classified by task difficulty on SWE-Bench Verified. CAT outperforms baselines, with notably larger gains on medium and hard tasks that demand complex, long-horizon reasoning.

Method	Max Steps	Tokens	Pass (%)
ReAct	150	1.96M	53.2
ReAct	500	2.54M	48.8
Threshold-Compression	150	2.49M	54.2
Threshold-Compression	500	5.18M	53.8
CAT (Base SFT)	150	1.95M	53.0
CAT (Base SFT)	500	5.07M	55.0
CAT	150	1.89M	54.8
CAT	500	2.75M	57.8

Table 3: Performance and token usage under different maximum interaction steps on SWE-Bench-Verified.

and longer-horizon context maintenance, the training signal introduced by CAT-Instruct becomes more effective. These findings further highlight the advantages of CAT-GENERATOR and CAT in addressing difficult tasks.

Effect of Interaction Budget and Token Efficiency. Table 3 summarizes the performance and token usage of different methods under varying maximum interaction budgets. When a larger number of interactions is allowed (500 steps), CAT achieves the highest pass rate among all methods, indicating stronger long-horizon reasoning capability. Under a smaller interaction budget (150 steps), CAT both maintains competitive performance and uses the fewest tokens across all methods. This favorable performance-efficiency trade-off is primarily attributed to the proposed context management mechanism, which prioritizes salient information while compressing redundant history, enabling efficient reasoning within a constrained context budget. The comparison between CAT and its Base SFT variant isolates the effect of CAT-GENERATOR. Results show that SWE-Compressor consistently outperforms Base SFT model, especially under larger interaction budgets, confirming the importance of CAT-GENERATOR.

5 Related Work

Code Agents. As LLMs plateau on standalone code generation, recent work shifts toward agentic systems for real-world software engineering, typically evaluated on SWE-bench and Multi-SWE-bench (Jimenez et al., 2023; Zan et al., 2025). Existing approaches improve performance either by enhancing agent designs with interactive tools and test-time scaling (Wang et al., 2024; Yang et al., 2024; Jain et al., 2025; Lin et al., 2025; Gao et al., 2025), or by strengthening model-level agentic capability through large executable environments, synthetic supervision, and reinforcement learning (Pan et al., 2024a; Badertdinov et al., 2025; Guo et al., 2025b; Yang et al., 2025b; Wang et al., 2025a; Sonwane et al., 2025; Wei et al., 2025; He et al., 2025; Luo et al., 2025).

Context Management. To support long-horizon decision making, prior work explores context compression and memory mechanisms, including saliency-based filtering, hierarchical summarization, and multi-level memory architectures (Li, 2023; Jiang et al., 2023; Pan et al., 2024b; Ye et al., 2025; Sun et al., 2025; Packer et al., 2023; Wang et al., 2023; Hu et al., 2025; Xiao et al., 2024). However, most rely on static compression or fixed memory policies, whereas our Tool Condensor enables dynamic, execution-driven context management that actively preserves decision-critical information over extended horizons.

6 Conclusion

In this work, we propose CAT, a context management paradigm that treats context maintenance as a first-class, toolized capability in long-horizon agents. By integrating context management into the decision process of agent, CAT enables proactive and structured condensation of interaction history, overcoming the limitations of append-only contexts and passive compression. We further introduce a trajectory-level supervision framework with an offline retrofitting pipeline to inject context-management actions into full interaction trajectories. Experiments on SWE-Bench demonstrate that CAT consistently outperforms ReAct agents and static compression baselines, while maintaining stable context usage and scalability under extended interaction budgets, underscoring the importance of modeling context evolution as an active and learnable component of agent behavior.

569 Limitations

570 **Heuristic supervision for tool invocation.** CAT-
571 GENERATOR injects context-management actions
572 by identifying condenser positions from recon-
573 structed trajectories using heuristic signals (e.g.,
574 context growth and stage boundaries). As a result,
575 the learned invocation behavior may not be globally
576 optimal under different interaction styles, reposi-
577 tories, or context budgets, and we do not provide
578 a principled objective for minimizing tool calls or
579 maximizing information gain.

580 Summary fidelity and error accumulation.

581 CAT relies on LLM-generated long-term memory
582 blocks. We do not include an explicit verification
583 mechanism to detect hallucinations, omissions, or
584 inconsistencies in compressed memories, nor do
585 we quantify how summary errors accumulate over
586 hundreds of rounds. This can potentially propa-
587 gate incorrect beliefs to later steps when the agent
588 reuses compressed context.

589 Ethical Statement

590 **Potential Risks.** Our work strengthens long-
591 horizon SWE agents by improving context man-
592 agement, which may increase the effectiveness of
593 automated code changes. Potential risks include (i)
594 over-reliance on agent outputs in safety- or security-
595 critical settings, where erroneous patches or halluci-
596 nated long-term memories could cause failures; (ii)
597 misuse to generate or scale unwanted code changes;
598 and (iii) privacy risks when agents process reposi-
599 tory artifacts that may contain personal identifiers.
600 We recommend sandboxed execution, human re-
601 view, and conservative deployment practices.

602 References

603 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama
604 Ahmad, Ilge Akkaya, Florencia Leoni Aleman,
605 Diogo Almeida, Janko Altenschmidt, Sam Altman,
606 Shyamal Anadkat, et al. 2023. Gpt-4 technical report.
607 *arXiv preprint arXiv:2303.08774*.

608 Anthropic. 2024. **Introducing claude 3.5 son-**
609 **net.** [https://www.anthropic.com/news/claude-](https://www.anthropic.com/news/claude-3-5-sonnet)
610 [3-5-sonnet](https://www.anthropic.com/news/claude-3-5-sonnet). Accessed: 2025-12-22.

611 Anthropic. 2025. **Introducing claude sonnet**
612 **4.5.** [https://www.anthropic.com/news/claude-](https://www.anthropic.com/news/claude-sonnet-4-5)
613 [sonnet-4-5](https://www.anthropic.com/news/claude-sonnet-4-5). Accessed: 2025-12-22.

614 Jacob Austin, Augustus Odena, Maxwell Nye, Maarten
615 Bosma, Henryk Michalewski, David Dohan, Ellen
616 Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. 2021.

Program synthesis with large language models. *arXiv preprint arXiv:2108.07732*. 617 618

Ibragim Badertdinov, Alexander Golubev, Maksim
Nekrashevich, Anton Shevtsov, Simon Karasik, An-
drei Andriushchenko, Maria Trofimova, Daria Litv-
intseva, and Boris Yangel. 2025. Swe-rebench: An
automated pipeline for task collection and decon-
taminated evaluation of software engineering agents.
arXiv preprint arXiv:2505.20411. 619 620 621 622 623 624 625

Linzhen Chai, Shukai Liu, Jian Yang, Yuwei Yin,
Ke Jin, Jiaheng Liu, Tao Sun, Ge Zhang, Changyu
Ren, Hongcheng Guo, et al. 2024. Mceval: Mas-
sively multilingual code evaluation. *arXiv preprint*
arXiv:2406.07436. 626 627 628 629 630

Mark Chen. 2021. Evaluating large language models
trained on code. *arXiv preprint arXiv:2107.03374*. 631 632

Pengfei Gao, Zhao Tian, Xiangxin Meng, Xinchun
Wang, Ruida Hu, Yuanan Xiao, Yizhou Liu, Zhao
Zhang, Junjie Chen, Cuiyun Gao, et al. 2025.
Trae agent: An llm-based agent for software en-
gineering with test-time scaling. *arXiv preprint*
arXiv:2507.23370. 633 634 635 636 637 638

Google Cloud. 2025. **Gemini 2.5 pro | generative ai**
on vertex ai. [https://docs.cloud.google.com/](https://docs.cloud.google.com/vertex-ai/generative-ai/docs/models/gemini/2-5-pro)
[vertex-ai/generative-ai/docs/models/](https://docs.cloud.google.com/vertex-ai/generative-ai/docs/models/gemini/2-5-pro)
[gemini/2-5-pro](https://docs.cloud.google.com/vertex-ai/generative-ai/docs/models/gemini/2-5-pro). Accessed: 2025-12-22. 639 640 641 642

Google DeepMind. 2025. **Gemini 3 pro — our most**
intelligent ai model. [https://deepmind.google/](https://deepmind.google/models/gemini/pro/)
[models/gemini/pro/](https://deepmind.google/models/gemini/pro/). Accessed: 2025-12-22. 643 644 645

Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song,
Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma,
Peiyi Wang, Xiao Bi, et al. 2025a. Deepseek-r1: In-
centivizing reasoning capability in llms via reinforce-
ment learning. *arXiv preprint arXiv:2501.12948*. 646 647 648 649 650

Lianghong Guo, Yanlin Wang, Caihua Li, Pengyu Yang,
Jiachi Chen, Wei Tao, Yingtian Zou, Duyu Tang, and
Zibin Zheng. 2025b. Swe-factory: Your automated
factory for issue resolution training data and evalua-
tion benchmarks. *arXiv preprint arXiv:2506.10954*. 651 652 653 654 655

Zhenyu He, Qingping Yang, Wei Sheng, Xiaojian
Zhong, Kechi Zhang, Chenxin An, Wenlei Shi,
Tianle Cai, Di He, Jiase Chen, and Jingjing
Xu. 2025. Swe-swiss: A multi-task fine-tuning
and rl recipe for high-performance issue reso-
lution. [https://www.notion.so/SWE-Swiss-](https://www.notion.so/SWE-Swiss-A-Multi-Task-Fine-Tuning-and-RL-Recipe-for-High-Performance-Issue-Resolution-21e174dedd4880ea829ed4c861c44f88)
[A-Multi-Task-Fine-Tuning-and-RL-Recipe-](https://www.notion.so/SWE-Swiss-A-Multi-Task-Fine-Tuning-and-RL-Recipe-for-High-Performance-Issue-Resolution-21e174dedd4880ea829ed4c861c44f88)
[for-High-Performance-Issue-Resolution-](https://www.notion.so/SWE-Swiss-A-Multi-Task-Fine-Tuning-and-RL-Recipe-for-High-Performance-Issue-Resolution-21e174dedd4880ea829ed4c861c44f88)
[21e174dedd4880ea829ed4c861c44f88](https://www.notion.so/SWE-Swiss-A-Multi-Task-Fine-Tuning-and-RL-Recipe-for-High-Performance-Issue-Resolution-21e174dedd4880ea829ed4c861c44f88). Notion
Blog. 656 657 658 659 660 661 662 663 664 665

Xinyi Hou, Yanjie Zhao, Yue Liu, Zhou Yang, Kailong
Wang, Li Li, Xiapu Luo, David Lo, John Grundy,
and Haoyu Wang. 2024. Large language models
for software engineering: A systematic literature re-
view. *ACM Transactions on Software Engineering*
and Methodology, 33(8):1–79. 666 667 668 669 670 671

672	Mengkang Hu, Tianxing Chen, Qiguang Chen, Yao Mu,	Shukai Liu, Linzheng Chai, Jian Yang, Jiajun Shi,	729
673	Wenqi Shao, and Ping Luo. 2025. Hiagent: Hier-	He Zhu, Liran Wang, Ke Jin, Wei Zhang, Hualei	730
674	archical working memory management for solving	Zhu, Shuyue Guo, et al. 2024b. Mdeval: Mas-	731
675	long-horizon agent tasks with large language model.	sively multilingual code debugging. <i>arXiv preprint</i>	732
676	In <i>Proceedings of the 63rd Annual Meeting of the</i>	<i>arXiv:2411.02310</i> .	733
677	<i>Association for Computational Linguistics (Volume</i>		
678	<i>1: Long Papers)</i> , pages 32779–32798.	Ilya Loshchilov and Frank Hutter. 2017. Decou-	734
		pled weight decay regularization . <i>arXiv preprint</i>	735
679	Chengying Huan, Ziheng Meng, Yongchao Liu,	<i>arXiv:1711.05101</i> .	736
680	Zhengyi Yang, Yun Zhu, Yue Yun, Shipeng Li, Rong		
681	Gu, Xiabao Wu, Haitao Zhang, et al. 2025. Scal-	Michael Luo, Naman Jain, Jaskirat Singh, Sijun	737
682	ing graph chain-of-thought reasoning: A multi-agent	Tan, Ameen Patel, Qingyang Wu, Alpaya Ariyak,	738
683	framework with efficient llm serving. <i>arXiv preprint</i>	Colin Cai, Shang Zhu Tarun Venkat, Ben Athi-	739
684	<i>arXiv:2511.01633</i> .	waratkun, Manan Roongta, Ce Zhang, Li Erran	740
		Li, Raluca Ada Popa, Koushik Sen, and Ion	741
685	Binyuan Hui, Jian Yang, Zeyu Cui, Jiayi Yang, Day-	Stoica. 2025. Deepswe: Training a state-of-	742
686	iheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang,	the-art coding agent from scratch by scaling	743
687	Bowen Yu, Keming Lu, et al. 2024. Qwen2. 5-coder	rl. https://pretty-radio-b75.notion.site/	744
688	technical report. <i>arXiv preprint arXiv:2409.12186</i> .	DeepSWE-Training-a-Fully-Open-sourced-	745
		State-of-the-Art-Coding-Agent-by-Scaling-	746
689	Naman Jain, Jaskirat Singh, Manish Shetty, Liang	RL-22281902c1468193aabb9a8c59bbe33 . No-	747
690	Zheng, Koushik Sen, and Ion Stoica. 2025. R2e-	tion Blog.	748
691	gym: Procedural environments and hybrid verifiers		
692	for scaling open-weights swe agents. <i>arXiv preprint</i>	Yingwei Ma, Rongyu Cao, Yongchang Cao, Yue Zhang,	749
693	<i>arXiv:2504.07164</i> .	Jue Chen, Yibo Liu, Yuchen Liu, Binhua Li, Fei	750
		Huang, and Yongbin Li. 2024. Lingma swe-gpt: An	751
694	Huiqiang Jiang, Qianhui Wu, Chin-Yew Lin, Yuqing	open development-process-centric language model	752
695	Yang, and Lili Qiu. 2023. Llmlingua: Compressing	for automated software improvement. <i>arXiv preprint</i>	753
696	prompts for accelerated inference of large language	<i>arXiv:2411.00622</i> .	754
697	models. <i>arXiv preprint arXiv:2310.05736</i> .		
		OpenAI. 2023. Gpt-4 technical report . <i>arXiv preprint</i>	755
698	Carlos E Jimenez, John Yang, Alexander Wettig,	<i>arXiv:2303.08774</i> .	756
699	Shunyu Yao, Kexin Pei, Ofir Press, and Karthik		
700	Narasimhan. 2023. Swe-bench: Can language mod-	OpenAI. 2025. Gpt-5.1: A smarter, more conversational	757
701	els resolve real-world github issues? <i>arXiv preprint</i>	chatgpt . https://openai.com/index/gpt-5-1/ .	758
702	<i>arXiv:2310.06770</i> .	Accessed: 2025-12-22.	759
		OpenHands. 2025. Openhands context condensensation	760
703	Minki Kang, Wei-Ning Chen, Dongge Han, Huseyin A	for more efficient ai agents .	761
704	Inan, Lukas Wutschitz, Yanzhi Chen, Robert Sim,		
705	and Saravan Rajmohan. 2025. Acon: Optimizing	Charles Packer, Vivian Fang, Shishir_G Patil, Kevin	762
706	context compression for long-horizon llm agents.	Lin, Sarah Wooders, and Joseph_E Gonzalez. 2023.	763
707	<i>arXiv preprint arXiv:2510.00615</i> .	Memgpt: Towards llms as operating systems.	764
		Jiayi Pan, Xingyao Wang, Graham Neubig, Navdeep	765
708	Yucheng Li. 2023. Unlocking context constraints of	Jaitly, Heng Ji, Alane Suhr, and Yizhe Zhang. 2024a.	766
709	llms: Enhancing context efficiency of llms with self-	Training software engineering agents and verifiers	767
710	information-based content filtering. <i>arXiv preprint</i>	with swe-gym. <i>arXiv preprint arXiv:2412.21139</i> .	768
711	<i>arXiv:2304.12102</i> .		
712	Jiaye Lin, Yifu Guo, Yuzhen Han, Sen Hu, Ziyi Ni,	Zhuoshi Pan, Qianhui Wu, Huiqiang Jiang, Menglin	769
713	Licheng Wang, Mingguang Chen, Hongzhang Liu,	Xia, Xufang Luo, Jue Zhang, Qingwei Lin, Victor	770
714	Ronghao Chen, Yangfan He, et al. 2025. Se-agent:	Rühle, Yuqing Yang, Chin-Yew Lin, et al. 2024b.	771
715	Self-evolution trajectory optimization in multi-step	Llmlingua-2: Data distillation for efficient and faith-	772
716	reasoning with llm-based agents. <i>arXiv preprint</i>	ful task-agnostic prompt compression. <i>arXiv preprint</i>	773
717	<i>arXiv:2508.02085</i> .	<i>arXiv:2403.12968</i> .	774
		Noah Shinn, Federico Cassano, Ashwin Gopinath,	775
718	Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang,	Karthik Narasimhan, and Shunyu Yao. 2023. Re-	776
719	Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi	flexion: Language agents with verbal reinforcement	777
720	Deng, Chenyu Zhang, Chong Ruan, et al. 2024a.	learning . <i>Advances in Neural Information Process-</i>	778
721	Deepseek-v3 technical report. <i>arXiv preprint</i>	<i>ing Systems</i> , 36:8634–8652.	779
722	<i>arXiv:2412.19437</i> .		
723	Jun Liu, Zhenglun Kong, Changdi Yang, Fan Yang,	Atharv Sonwane, Isadora White, Hyunji Lee, Matheus	780
724	Tianqi Li, Peiyan Dong, Joannah Nanjeyke, Hao	Pereira, Lucas Caccia, Minseon Kim, Zhengyan	781
725	Tang, Geng Yuan, Wei Niu, et al. 2025. Rcr-router:	Shi, Chinmay Singh, Alessandro Sordoni, Marc-	782
726	Efficient role-aware context routing for multi-agent	Alexandre Côté, et al. 2025. Bugpilot: Complex	783
727	llm systems with structured memory. <i>arXiv preprint</i>		
728	<i>arXiv:2508.04903</i> .		

784	bug generation for efficient learning of swe skills.	Yuan-An Xiao, Pengfei Gao, Chao Peng, and Yingfei Xiong. 2025. Improving the efficiency of llm agent systems through trajectory reduction. <i>arXiv preprint arXiv:2510.19898</i> .	838
785			839
786	Weiwei Sun, Miao Lu, Zhan Ling, Kang Liu, Xuesong Yao, Yiming Yang, and Jiecao Chen. 2025. Scaling long-horizon llm agent via context-folding. <i>arXiv preprint arXiv:2510.11967</i> .	Chengxing Xie, Bowen Li, Chang Gao, He Du, Wai Lam, Difan Zou, and Kai Chen. 2025. Swe-fixer: Training open-source llms for effective and efficient github issue resolution. <i>arXiv preprint arXiv:2509.23586</i> .	840
787			841
788			842
789			843
790	Valentin Tablan, Scott Taylor, Gabriel Hurtado, Kristofer Bernhem, Anders Uhrenholt, Gabriele Farei, and Karo Moilanen. 2025. Smarter together: Creating agentic communities of practice through shared experiential learning. <i>arXiv preprint arXiv:2511.08301</i> .		844
791			845
792			846
793		An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, et al. 2025a. Qwen3 technical report. <i>arXiv preprint arXiv:2505.09388</i> .	847
794			848
795	ByteDance Seed Team. 2025. Seed-oss-36b-instruct: A 36b instruction-tuned open-source large language model. https://huggingface.co/ByteDance-Seed/Seed-OSS-36B-Instruct . Accessed: 2025-12-22.		849
796			850
797		John Yang, Carlos E Jimenez, Alexander Wettig, Kilian Lieret, Shunyu Yao, Karthik Narasimhan, and Ofir Press. 2024. Swe-agent: Agent-computer interfaces enable automated software engineering. <i>Advances in Neural Information Processing Systems</i> , 37:50528–50652.	851
798			852
799			853
800	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. <i>arXiv preprint arXiv:2302.13971</i> .		854
801			855
802			856
803		John Yang, Kilian Lieret, Carlos E Jimenez, Alexander Wettig, Kabir Khandpur, Yanzhe Zhang, Binyuan Hui, Ofir Press, Ludwig Schmidt, and Diyi Yang. 2025b. Swe-smith: Scaling data for software engineering agents. <i>arXiv preprint arXiv:2504.21798</i> .	857
804			858
805			859
806	Bing Wang, Xinnian Liang, Jian Yang, Hui Huang, Shuangzhi Wu, Peihao Wu, Lu Lu, Zejun Ma, and Zhoujun Li. 2023. Scm: Enhancing large language model with self-controlled memory framework. <i>arXiv e-prints</i> , pages arXiv–2304.		860
807			861
808		Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R Narasimhan, and Yuan Cao. 2022. React: Synergizing reasoning and acting in language models. In <i>The eleventh international conference on learning representations</i> .	862
809			863
810			864
811	Junhao Wang, Daoguang Zan, Shulin Xin, Siyao Liu, Yurong Wu, and Kai Shen. 2025a. Swe-mirror: Scaling issue-resolving datasets by mirroring issues across repositories. <i>arXiv preprint arXiv:2509.08724</i> .		865
812			866
813		Rui Ye, Zhongwang Zhang, Kuan Li, Huifeng Yin, Zhengwei Tao, Yida Zhao, Liangcai Su, Liwen Zhang, Zile Qiao, Xinyu Wang, et al. 2025. Agent-fold: Long-horizon web agents with proactive context management. <i>arXiv preprint arXiv:2510.24699</i> .	867
814			868
815			869
816	Qingyue Wang, Yanhe Fu, Yanan Cao, Shuai Wang, Zhiliang Tian, and Liang Ding. 2025b. Recursively summarizing enables long-term dialogue memory in large language models. <i>Neurocomputing</i> , 639:130193.		870
817			871
818		Daoguang Zan, Zhirong Huang, Wei Liu, Hanwu Chen, Linhao Zhang, Shulin Xin, Lu Chen, Qi Liu, Xiaojian Zhong, Aoyan Li, et al. 2025. Multi-swe-bench: A multilingual benchmark for issue resolving. <i>arXiv preprint arXiv:2504.02605</i> .	872
819			873
820			874
821	Xingyao Wang, Boxuan Li, Yufan Song, Frank F Xu, Xiangru Tang, Mingchen Zhuge, Jiayi Pan, Yueqi Song, Bowen Li, Jaskirat Singh, et al. 2024. Openhands: An open platform for ai software developers as generalist agents. <i>arXiv preprint arXiv:2407.16741</i> .		875
822			876
823			
824			
825			
826	Yuxiang Wei, Olivier Duchenne, Jade Copet, Quentin Carbonneaux, Lingming Zhang, Daniel Fried, Gabriel Synnaeve, Rishabh Singh, and Sida I Wang. 2025. Swe-rl: Advancing llm reasoning via reinforcement learning on open software evolution. <i>arXiv preprint arXiv:2502.18449</i> .		
827			
828			
829			
830			
831			
832	Chaojun Xiao, Pengle Zhang, Xu Han, Guangxuan Xiao, Yankai Lin, Zhengyan Zhang, Zhiyuan Liu, and Maosong Sun. 2024. Inllm: Training-free long-context extrapolation for llms with an efficient context memory. <i>Advances in Neural Information Processing Systems</i> , 37:119638–119661.		
833			
834			
835			
836			
837			

877
878
879
880
881
882
883
884
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
918
919
920
921
922
923
924
925
926
927

A Related Work

Code Agents. As LLMs approach saturation on traditional code generation tasks, research has shifted toward their agentic capabilities in real-world codebases, typically evaluated on SWE-bench and Multi-SWE-bench (Jimenez et al., 2023; Zan et al., 2025). Existing efforts fall into two main directions. The first focuses on agent design. OpenHands and SWE-Agent (Wang et al., 2024; Yang et al., 2024) equip LLMs with interfaces such as editors and shells, enabling iterative file edits and command execution. Building on this setup, R2E-Gym, SE-Agent, and TraeAgent (Jain et al., 2025; Lin et al., 2025; Gao et al., 2025) pursue test-time scaling by increasing sampling and decision steps to approximate upper-bound performance. The second direction seeks to improve model-level agentic capability. SWE-Gym, SWE-rebench, and SWE-Factory (Pan et al., 2024a; Badertdinov et al., 2025; Guo et al., 2025b) build large executable environments for agent training, while SWE-Smith, SWE-Mirror, and BugPilot (Yang et al., 2025b; Wang et al., 2025a; Sonwane et al., 2025) generate synthetic tasks and trajectories for supervised fine-tuning. Reinforcement learning has also been explored: SWE-RL (Wei et al., 2025) uses patch similarity as rewards, and SWE-Swiss and Deep-SWE (He et al., 2025; Luo et al., 2025) investigate execution-based rewards as a promising alternative.

Context Management. As agent technologies become crucial for long-horizon tasks, recent work aims to help LLMs sustain stable decision-making in complex environments. One line develops efficient context compression and memory management: methods like Selective Context, AgentFold, Context-folding, and LLMingua (Li, 2023; Jiang et al., 2023; Pan et al., 2024b; Ye et al., 2025; Sun et al., 2025) shorten inputs via salient-information filtering or hierarchical summaries, while AgentDiet and ACON (Xiao et al., 2025; Kang et al., 2025) enhance long-trajectory fidelity through reflection, contrastive learning, or latent-space compression. Another direction builds multi-level memory systems—MemGPT, SCM, HIAGENT, and InfLLM (Packer et al., 2023; Wang et al., 2023; Hu et al., 2025; Xiao et al., 2024)—coordinating short- and long-term memory with retrieval modules for consistent extended reasoning. Multi-agent studies explore role-aware routing, shared memory, and divide-and-conquer frameworks to mitigate context explosion (Liu et al., 2025; Tablan et al.,

2025; Huan et al., 2025). Yet, these approaches rely on static compression, heuristic retrieval, or fixed memory, limiting adaptability and long-term coherence. In contrast, Tool Condensor enables dynamic, execution-driven context management, flexibly removing redundancy while preserving critical information—essential for complex long-horizon tasks.

B Discussion

This subsection discusses CAT from two complementary perspectives: its generality beyond software engineering tasks and the underlying mechanisms revealed by the further analysis.

B.1 Generality beyond SWE tasks

Although this work focuses on software engineering scenarios, the core idea of CAT—explicitly modeling context management as a callable tool within the agent’s action space—is not inherently specific to SWE. Instead, it targets a broader class of agentic tasks characterized by long-horizon interaction, persistent dependence on historical decisions, and continuously growing contexts. Beyond repository-level code repair, this paradigm naturally extends to tasks such as deep research, complex document analysis, long-term planning, and multi-stage decision making. In these settings, agents must maintain stable representations of key facts, intermediate conclusions, and long-term objectives across many interaction rounds. Similar to SWE, append-only context strategies or passively triggered compression heuristics are prone to semantic drift and information overload. By internalizing context management as part of the agent’s decision process, CAT offers a unified and general approach to addressing these challenges.

B.2 Interpreting the further analysis results

The experiments presented in the Further Analysis section are designed to provide complementary insights into why CAT is effective, rather than serving as isolated empirical observations. First, the token usage analysis shows that CAT rapidly stabilizes context length within a bounded range and maintains this stability over hundreds of interaction rounds, demonstrating its ability to prevent uncontrolled context growth in practice. Second, the comparison with ReAct under identical interaction budgets reveals that proactive and structured context folding not only reduces redundancy but

928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975

976	also improves the effective utilization of historical	Accordingly, we do not obtain individual consent	1025
977	information, leading to sustained performance	beyond what is implied by the original public post-	1026
978	gains as interaction horizons increase. Third, the	ing and the terms of use of the underlying platforms	1027
979	difficulty-based analysis indicates that the bene-	and datasets.	1028
980	fits of CAT become more pronounced on medium	Ethics review board approval. Because this	1029
981	and hard tasks, where deeper reasoning chains and	work does not involve human-subject experiments	1030
982	longer temporal dependencies are required. Finally,	or new data collection from human participants, we	1031
983	the joint analysis of interaction budgets and token	did not seek ethics review board (IRB) approval;	1032
984	efficiency highlights the role of CAT-GENERATOR	the study is typically considered exempt under stan-	1033
985	in teaching the model to balance when to compress	dard IRB criteria for research using publicly avail-	1034
986	context and how to preserve decision-critical infor-	able information.	1035
987	mation.	AI assistants. AI-assisted tools were used solely	1036
988	Taken together, these analyses support a central	to polish language and improve clarity during	1037
989	conclusion: the strength of CAT does not stem	manuscript preparation. All technical content, anal-	1038
990	from compression alone, but from treating context	yses, and conclusions were conceived, verified, and	1039
991	evolution as an explicit, learnable decision process.	approved by the authors.	1040
992	This design enables agents to maintain accurate		
993	task understanding under bounded context budgets,		
994	providing a principled foundation for stable long-		
995	horizon reasoning.		
996	C Checklist		
997	Personally identifying or offensive content.		
998	Our data sources (SWE-Bench-Verified and train-		
999	ing instances from SWE-smith and SWE-ReBench)		
1000	are derived from public GitHub repositories/issues,		
1001	which may contain usernames, commit metadata,		
1002	or other identifiers, and may include offensive		
1003	content in natural-language discussions. In this		
1004	work, we do not claim a complete, formal audit		
1005	for PII/offensive content beyond relying on the		
1006	upstream dataset construction and filtering proto-		
1007	cols described in the corresponding dataset papers.		
1008	When releasing any additional artifacts, we recom-		
1009	mend (and plan) to apply automated scrubbing for		
1010	common PII patterns (e.g., emails, access tokens)		
1011	and to avoid redistributing raw issue threads when		
1012	not necessary.		
1013	D: Human subjects / annotators. We did not re-		
1014	cruit or compensate human participants/annotators		
1015	for this work. We use publicly released bench-		
1016	marks/datasets; any human curation (e.g., SWE-		
1017	Bench-Verified) is described in the corresponding		
1018	dataset papers.		
1019	Data consent. We do not newly collect data di-		
1020	rectly from individuals. Our benchmarks and train-		
1021	ing instances are derived from publicly available		
1022	software engineering artifacts (e.g., GitHub reposi-		
1023	tories and issue threads) released as part of SWE-		
1024	Bench-Verified, SWE-smith, and SWE-ReBench.		