# PAACE: A Plan-Aware Automated Agent Context Engineering Framework

**Kamer Ali Yuksel** [1]

## Abstract

Large Language Model (LLM) agents are increasingly deployed in complex, multi-step workflows involving planning, tool use, reflection, and interaction with external knowledge systems. These workflows generate rapidly expanding contexts that must be curated, transformed, and compressed to maintain fidelity, avoid attention dilution, and reduce inference cost. Prior work on summarization and query-aware compression largely ignores the *multi-step, plan-aware* nature of agentic reasoning. In this work, we introduce **PAACE** (**P**lan-**A**ware **A**utomated **C**ontext **E**ngineering), a unified framework for optimizing the evolving state of LLM agents through next-$k$-task relevance modeling, plan-structure analysis, instruction co-refinement, and function-preserving compression. PAACE comprises (1) **PAACE-Syn**, a large-scale generator of synthetic agent workflows annotated with stepwise compression supervision, and (2) **PAACE-FT**, a family of distilled, plan-aware compressors trained from successful teacher demonstrations. Experiments on long-horizon benchmarks—AppWorld, OfficeBench, and 8-Objective QA—demonstrate that PAACE consistently improves agent correctness while substantially reducing context load. On AppWorld, PAACE achieves higher accuracy than all baselines, while lowering peak context and cumulative dependency. On OfficeBench and multi-hop QA, PAACE improves both accuracy and F1, achieving lower steps, peak tokens, and attention dependency. Distilled PAACE-FT retains 97% of the teacher's performance while reducing inference cost by over an order of magnitude, enabling practical deployment of plan-aware compression with compact models.

---

[1]aiXplain Inc, San Jose, USA. Correspondence to: Kamer Ali Yuksel <kamer@aixplain.com>.

## 1. Introduction

LLM-driven agents have emerged as a central paradigm for solving complex, multi-step tasks across domains such as software development, research assistance, operations automation, legal workflows, data analysis, and enterprise decision-making. Systems such as ReAct (Yao et al., 2023a), Toolformer (Schick et al., 2023a), AutoGPT, Devin, WebArena (Zhou et al., 2024), and AgentBench (Liu et al., 2023) highlight the promise and challenges of designing agents with reasoning, planning, and tool-use capabilities. As these systems evolve, a consistent bottleneck has become apparent: **context management**. An LLM agent's state is represented not by model parameters but by its *prompt context*: the system instructions, the evolving plan, previous reasoning traces, tool results, user instructions, long-term memories, retrieved knowledge, and intermediate outputs. Agents do not operate on a single-step basis. They execute a plan $= [\tau_1, \tau_2, \ldots, \tau_n]$ with dependencies between tasks, and only a subset of context is relevant for each stage. As tasks grow in depth and breadth, this state becomes increasingly large, noisy, redundant, and expensive to process. Even models with 200k–1M token windows exhibit degraded reasoning quality ("context rot") when overloaded with irrelevant or poorly structured information. Recent industry reports emphasize that modern agentic failures are overwhelmingly *context failures*, not model failures. Despite advances in model architecture and context length, agents fail when: crucial information is dropped or buried in irrelevant text, irrelevant details overload the model's attention, multiple tasks compete for context bandwidth, instructions contradict or drift over time, or context becomes too long to process efficiently. We frame context engineering as learning a state compression policy over an agent's evolving execution state, trained via outcome-preserving supervision, rather than as heuristic prompt editing.

While prompt engineering optimizes *initial* instructions, and RAG systems optimize *retrieval*, the missing discipline is **context engineering**: the science of continuously optimizing what the agent sees at each step. Existing approaches partially address pieces of this problem. Classical and instruction-following summarizers such as BART (Lewis et al., 2019), FLAN-T5 (Chung et al., 2022), and compression-oriented methods like LLMLingua (Jiang et al., 2023) generate concise summaries but often remove struc-

tural dependencies required for multi-step reasoning. Summaries flatten causal links across steps, which harms agent planning and tool-use workflows. Methods such as Self-RAG (Asai et al., 2024) and LLMLingua-2 (Jiang et al., 2024) optimize relevance for a *single upcoming query*. However, they do not model next-$k$ steps, multi-hop dependencies, or evolving plans. Provence (Wu et al., 2024) performs binary keep/drop trimming via a relevance classifier but does not support rewriting, instruction refinement, dependency tracking, or structured context shaping.

Modern long-context LLMs (Claude 3.5 Sonnet, Gemini 1.5 Pro, GPT-4o mini L) offer 200k–1M+ token windows, yet still suffer from attention dilution, context "rot," and quadratic cost growth in practice (Zaheer et al., 2020; Guo et al., 2021). Large windows do not solve the problem of *poorly structured* or *irrelevant* context. Memory and reflection systems such as MemAgent (Yu et al., 2025a), Reflexion (Shinn et al., 2023), and Generative Agents (Park et al., 2023) improve retrieval and episodic memory, but do not perform context shaping, state restructuring, or next-instruction refinement. Yu et al. (2025b) optimizes next-step relevance through natural-language compression guidelines, but it does not model next-$k$-step plan structure, does not refine instructions, and cannot jointly optimize plan-aware context and instruction transformations. Our primary distinction is conditioning on multiple future plan steps and the global workflow structure. No existing system considers:

- **plan structure** (preconditions, intermediate states, temporal dependencies) and **next-$k$-task relevance**
- **instruction refinement jointly with context selection**
- **global optimization over multi-step workflows**
- **learnable policies for context shaping**

## Contributions

1. **PAACE**, a framework for *plan-aware context engineering*, unifying context pruning, summarization, rewriting, compression, re-injection, and task refinement.
2. **A context relevance formalization** incorporating plan-dependence, next-$k$-task relevance, structural decomposition, and instruction–context co-evolution.
3. **Explicit relevance scoring signals** for selection and optimization: plan structure, next-$k$ task preservation, outcome fidelity, and information-theoretic retention.
4. **PAACE-Syn**, a scalable synthetic data generation system producing millions of annotated workflow traces.
5. **PAACE-FT**, a family of fine-tuned LLMs specialized for context shaping and compression.
6. **Comprehensive experiments** showing PAACE improves correctness, efficiency, and cost across multiple agent benchmarks.

Together, these contributions establish **plan-aware context engineering** as a critical component of robust long-horizon

agent design and provide the first end-to-end framework that jointly models plan relevance and context structure to support cost-efficient, high-fidelity agent reasoning.

## 2. Related Work

This section synthesizes prior work across ten research areas relevant to PAACE: (1) summarization and compression, (2) query-aware and task-aware reduction, (3) long-context models, (4) memory architectures for agents, (5) retrieval-augmented agents, (6) agent planning and multi-step reasoning, (7) prompt and instruction optimization, (8) context pruning and selection, (9) multi-agent systems and meta-reasoning, and (10) cognitive frameworks inspiring artificial memory systems. While each domain has contributed techniques relevant to context management, no prior work provides a unified, plan-aware, next-$k$-task context engineering framework like the proposed method.

Text summarization is one of the oldest forms of context reduction. Classical abstractive summarization models (Pointer-Generator (See et al., 2017), PEGASUS (Zhang et al., 2019), BART (Lewis et al., 2019), and T5 (Raffel et al., 2020)) enable high-level condensation of documents but lack task-awareness or plan-awareness. Instruction-following summarizers such as FLAN-T5 (Chung et al., 2022) and LLM-based summarizers (e.g., GPT-4, Claude 3) are more robust but still struggle with preserving the functional dependencies required by multi-step agent workflows. LLMLingua (Jiang et al., 2023) and LLMLingua-2 (Jiang et al., 2024) offer token-level compression using learned retention models, achieving significant reductions in context length with minimal task degradation. However, LLMLingua is query-aware, not plan-aware, and optimizes for single-turn QA relevance rather than multi-step agent tasks. Other approaches like BigBird (Zaheer et al., 2020), LongT5 (Guo et al., 2021), and hierarchical transformers extend context windows but do not optimize context *content*. PAACE goes beyond summarization by incorporating plan structure, task dependencies, instruction refinement, and relevance signals learned from synthetic workflows.

Recent models aim to optimize context for a single upcoming task. Self-RAG (Asai et al., 2024) generates retrieval queries conditioned on a question and prunes irrelevant context before passing it to the LLM. Although powerful for QA, Self-RAG is not general enough for multi-step agent workflows, which require modeling of future tasks, tool outputs, and plan transitions. Similarly, ACon (Agent Context Optimization) (Yu et al., 2025b) compresses agent interaction histories by optimizing "compression guidelines" via natural-language feedback. ACon improves context relevance for short-horizon tasks but: (1) does not model plan structure, (2) is restricted to next-step compression, (3) lacks instruction refinement, and (4) does not jointly optimize

plan-aware context selection. The primary distinction is that PAACE explicitly conditions compression on multiple future plan steps and the global workflow structure. PAACE thus generalizes ACon into a multi-task, multi-step, structural framework. Provence (Wu et al., 2024) introduces a binary relevance classifier for pruning retrieved document segments. This improves RAG pipelines but is too coarse for agent workflows, where reasoning chains, state representations, and plans require finer-grained transformations.

Many models attempt to address long-horizon reasoning by increasing context window sizes, including GPT-4o mini (128k), Claude 3.5 Sonnet (200k–1M), Gemini 1.5 Pro (1M), and DeepSeek-V2 (256k+). However, long context windows do not resolve the *relevance dilution* problem, where attention degrades over large inputs, leading to context "rot" and reasoning failures. Several approaches instead introduce explicit context or memory controllers. MemAgent (Yu et al., 2025a) uses reinforcement learning to overwrite a fixed-size memory state, enabling scaling to millions of tokens, but targets document QA rather than agent workflows and does not refine instructions or model plan structure. Compressive Transformers (Rae et al., 2020) and related architectural approaches perform learned compression within the model, but lack task- or plan-level grounding. More recent systems such as COMPASS (Li et al., 2024), Fold-GRPO (Sun et al., 2025), ReCAP (Zhang et al., 2025), and GoA (Joo et al., 2025) explore context control via learned controllers, reinforcement learning, or training-free multi-agent coordination, often operating over fixed-size latent memories or requiring architectural or inference-time modifications. These approaches are complementary to PAACE: they modify model architecture or memory mechanisms, whereas PAACE operates purely at the context level and can be integrated with any agent or backbone model. Task-agnostic learned compressors (Sun et al., 2023) optimize generic prompt reduction without conditioning on future plan structure; in contrast, PAACE explicitly models next-$k$ task dependencies in agent workflows. Unlike prompt compression heuristics, PAACE learns a plan-conditioned policy over latent agent states from outcome-level supervision, positioning it as an *empirical systems ML* approach based on trajectory-level rather than token-level objectives.

## 2.1. Memory Architectures for LLM Agents

Human memory is structured along: working memory, episodic memory, semantic memory, and procedural memory. AI researchers have used this taxonomy to design agent memory stacks (Park et al., 2023; Shinn et al., 2023). PAACE's decomposition of context into *actionable*, *structural*, *referential*, and *justificatory* components parallels these cognitive structures and provides a grounded theoretical basis for context relevance. Several systems draw inspiration from human memory:

- **Reflexion** (Shinn et al., 2023) stores agent self-assessments and retrieves them for self-improvement.
- **Generative Agents** (Park et al., 2023) maintain episodic memories and synthesize high-level reflections to drive believable behavior.
- **LTM/RAG systems** integrate persistent vector memory for agents.
- **Hierarchical memories** (short-term vs. long-term vs. global memory) are increasingly used in AI-agents.

They emphasize storage, retrieval, and memory distillation but do not optimize the shape, structure, or content of the context presented to the LLM during reasoning. PAACE builds on these systems by integrating: structured rewriting, task-dependent compression, and instruction refinement. PAACE also improves tool-call reliability by providing task-conditioned pruning, rewriting, and instruction improvement. Tool-use frameworks such as Toolformer (Schick et al., 2023b), ReAct (Yao et al., 2023b), MRKL (Ahn et al., 2022), and PAL (Gao et al., 2023) enable LLMs to ground their reasoning in external tools. Sophisticated agent frameworks combine retrieval, planning, and tool execution. Yet:

- Retrieval relevance is often limited to a single upcoming query.
- Retrieved content is appended rather than shaped or compressed.
- Tool outputs are accumulated and rarely distilled.

Multi-step reasoning frameworks such as CoT (Wei et al., 2022), ToT (Yao et al., 2023c), RAT (Zhang et al., 2023), and RCT (Zhou et al., 2023c) provide planning mechanisms, but do not manage accumulated context. Planner-LMs (Huang et al., 2024) explore explicit planning modules, but do not optimize the interaction between plans and context. PAACE instead treats context as the *state* upon which planning depends, and models the relevance of context components to future tasks. Prompt optimization and evolutionary prompt breeding (Zhou et al., 2023a; Pryzant et al., 2023) generate improved instructions using feedback loops. Similarly, methods such as DPO (Rafailov et al., 2023), RLHF (Ouyang et al., 2022), and PRO (Saunders et al., 2024) refine instructions or preferences. However, none address: next-$k$-task instruction refinement, plan-conditioned instruction rewriting, instruction-context *co*-refinement, or instruction updates during agentic workflows. To the best of our knowledge, PAACE introduces the first systematic approach to full context engineering, including instructions.

Beyond summarization, pruning frameworks, such as (Wu et al., 2024), aim to remove irrelevant text. Pruning alone is insufficient for multi-step agents because some content must be rewritten rather than removed. The global plan structure must be preserved as task dependencies are multi-

hop, and instructions would drift if not rewritten in sync with the context. PAACE incorporates pruning as *one* operator in a larger library of structured context transformations. Meta-agents such as Reflexion (Shinn et al., 2023), AlphaAI systems, CrewAI, and LATS (Zhou et al., 2023b) use multi-agent loops to stabilize reasoning. Some systems include meta-level evaluation or self-critique. But none propose a dedicated meta-agent for context engineering, fine-tuning specialized models (PAACE-FT), or plan-conditioned synthetic dataset generation (PAACE-Syn). We introduce context engineering as a core meta-agent function. No unified framework exists that jointly performs: plan-aware next-$k$-task context optimization, instruction refinement, rewriting + pruning + summarization + compression, structural reasoning over workflow graphs, and training models specifically for context engineering. To the best of our knowledge, PAACE is the first framework addressing these together.

## 3. Methodology

PAACE optimizes the evolving context of multi-step LLM agents through a two-stage paradigm: (1) a *teacher* LLM that performs plan- and next-$k$-task–aware compression guided by an evolving natural-language prompt, and (2) a distilled *student* model that imitates the teacher's successful compressions at low cost. Unlike prior work that hand-engineers context-editing operators, PAACE learns the entire compression policy in a data-driven way: the teacher's prompt is meta-learned via an LLM-driven evolutionary loop over real multi-step workflows, and the student is trained purely from the teacher's best demonstrations. PAACE learns this policy empirically via outcome-level filtering over long-horizon trajectories rather than through a differentiable objective or closed-form guarantees, prioritizing robust state fidelity and task success in realistic agentic workloads. The framework consists of two components:

- **PAACE-Syn**: large-scale generation of long-horizon, noisy agent workflows and paired full- vs. compressed-context trajectories;
- **PAACE-FT**: distillation of high-quality teacher compressions into a specialized, plan-aware compressor.

We treat the agent context as a latent execution state and learn a plan-conditioned, outcome-preserving compression policy. At step $t$, an LLM agent maintains a context state $C_t = \{I_0, P, \Pi, H_{0:t}, O_{0:t}, R_{0:t}, M\}$ where $I_0$ is the initial user instruction, $P$ is the system prompt, $\Pi = [\tau_1, \ldots, \tau_n]$ is the (explicit or implicit) plan, $H_{0:t}$ are reasoning traces, $O_{0:t}$ are tool/environment observations, $R_{0:t}$ are retrieved documents, and $M$ is long-term memory. PAACE does not explicitly manipulate these components via hand-crafted operators; instead, they implicitly shape what the teacher (and later the student) learns to preserve or

discard. Compression at step $t$ is defined as a mapping

$$\tilde{C}_t = \text{TeacherCompress}(C_t, \Pi_{t:t+k}; p), \quad (1)$$

where $\Pi_{t:t+k}$ denotes the next $k$ steps on the remaining workflow DAG subpath (for some task-dependent $k$), and $p$ is a natural-language compression prompt. The teacher conditions on the next tasks or the entire workflow description, and thus can preserve information needed for multiple upcoming steps. PAACE does not assume a closed-form parametric relevance model or a token-level differentiable loss. Instead, plan-aware relevance is operationalized through explicit, computable signals evaluated at the trajectory level: (i) semantic equivalence between full- and compressed-context outcomes, (ii) per-step compression ratios, (iii) implicit coverage of next-$k$ plan steps via conditioning, and (iv) preference judgments from a learned LLM-based evaluator. These signals jointly define a non-differentiable objective over task fidelity and context efficiency, optimized via black-box selection over full agent rollouts rather than token-level losses.

Although PAACE is motivated by concepts such as task graphs, causal influence, and information retention, we do not instantiate explicit symbolic task graphs or closed-form information-theoretic objectives in the current implementation. Instead, these notions are captured implicitly by conditioning compression on the *partial remaining subpath* of the agent workflow graph (DAG) together with outcome-level selection against full-context execution using explicit, testable criteria. Accordingly, we use the term "formal" to denote operational and verifiable selection mechanisms rather than analytical derivations or symbolic guarantees. Plan descriptions $\Pi$ are obtained either directly from benchmarks that provide structured task decompositions (e.g., AppWorld, OfficeBench) or from a LLM-based planner using the same backbone model as execution. PAACE does not require perfect plans: when plan steps are missing or partially incorrect, next-$k$ conditioning degrades compression quality gracefully but does not catastrophically affect execution. In practice, it behaves as a soft plan-aware regularizer rather than a brittle plan-dependent controller.

### 3.1. Synthetic Context Dataset Generation

PAACE constructs a large synthetic corpus of long-horizon agent workflows. Each workflow $W$ is sampled as: $W = (I_0, \Pi, \{\tau_i\}_{i=1}^n, \phi)$ where $I_0$ is a long, noisy initial input containing irrelevant logs and partial summaries, $\Pi$ is a natural-language description of the plan, $\tau_i$ are task instructions, and $\phi$ is a final output requirement. PAACE-Syn spans task domains, including document-centric workflows, web navigation and multi-application interactions, and multi-hop question answering with retrieval. Plans range from 5 to 30 steps, with tool interactions, such as, search and browsing actions, document and file manipulation, spreadsheet-style

table operations, and structured information extraction and aggregation over retrieved content. For each workflow, the agent is first run without compression:

$$C_1^{\text{full}} = \{I_0, P, \Pi\}, \tag{2}$$

$$C_{t+1}^{\text{full}} = \text{Update}\Big(C_t^{\text{full}}, \text{ Agent}\big(C_t^{\text{full}}, \tau_t\big)\Big), \tag{3}$$

and a final answer $\hat{y}^{\text{full}}$ is produced from $C_{n+1}^{\text{full}}$. The same workflow is then re-executed with teacher compression. At step $t$, given the current context $C_t$ and plan slice $\Pi_{t:t+k}$, the teacher produces a compressed context $\tilde{C}_t$ using prompt $p$. The agent then acts using only $\tilde{C}_t$:

$$\tilde{C}_t = \text{TeacherCompress}\big(C_t, \Pi_{t:t+k}; p\big), \tag{4}$$

$$\tilde{C}_{t+1} = \text{Update}\Big(\tilde{C}_t, \text{ Agent}\big(\tilde{C}_t, \tau_t\big)\Big). \tag{5}$$

After the final step, the agent produces $\hat{y}^{\text{comp}}$ from $\tilde{C}_{n+1}$. For each step $t$, PAACE records a tuple $\big(\Pi_{t:t+k}, C_t, \tilde{C}_t\big)$ together with compression statistics. PAACE uses several metrics to evaluate each compressed trajectory relative to its full-context counterpart. Let $\text{embed}(\cdot)$ be a fixed sentence embedding model, we calculate semantic equivalence as:

$$s = \cos\Big(\text{embed}(\hat{y}^{\text{full}}), \text{ embed}(\hat{y}^{\text{comp}})\Big). \tag{6}$$

An evaluation LLM is also prompted with the workflow description and both outputs $\hat{y}^{\text{full}}$ and $\hat{y}^{\text{comp}}$, and returns a binary label (*better* vs. *worse/equal*) plus a short rationale. This adds a robustness check complementary to embedding-based similarity. For each step $t$, we compute the character- or token-based compression ratio $r_t = |\tilde{C}_t|/|C_t|$. We require $0 < r_t < 1$ to avoid degenerate compressions (empty or longer than the original). Only successful trajectories are used as supervision for the student model. A compressed trajectory is labeled *successful* if and only if:

- the semantic equivalence exceeds a threshold: $s \geq \theta$ (typically $\theta = 0.85$);
- steps satisfy $0 < r_t < 1$ and produce non-empty $\tilde{C}_t$;
- the judge model does not deem the compressed answer strictly worse than the full answer.

From all trajectories generated, PAACE extracts only successful compression examples. For each such example, we construct a supervision tuple of the form:

$$\big(\ \underbrace{\Pi_{t:t+k}}_{\text{NEXT TASKS}}\ ,\ \underbrace{C_t}_{\text{CONTEXT}}\ \big) \mapsto \underbrace{\tilde{C}_t}_{\text{COMPRESSED CONTEXT}}. \tag{7}$$

A typical input to the student model concatenates the next-$k$-step instructions and the full context, while the target is the teacher's compressed context. Because only trajectories with high semantic equivalence and valid compression are retained, the resulting dataset captures a rich collection of *function-preserving* compressions tailored to plan structure and long-horizon dependencies.

## 3.2. Distilling Teacher Compressions

PAACE-FT trains a dedicated LLM to approximate the teacher's next-$k$-task–aware compression mapping. Let $x = (\Pi_{t:t+k}, C_t)$ denote the input text and $y = \tilde{C}_t$ the teacher's compressed context. The student model is trained with a standard causal language modeling loss:

$$\mathcal{L}(\theta) = -\mathbb{E}_{(x,y)} \sum_{i=1}^{|y|} \log p_\theta\big(y_i \mid x, y_{<i}\big), \tag{8}$$

where the tokens of $x$ are masked out from the loss. Distillation is not intended to outperform the teacher, but to replace an expensive compressor with a compact model that preserves plan-aware compression behavior. Through this distillation, the student learns to:

- compress long contexts into substantially shorter forms;
- retain variables, identifiers, and constraints used across multiple future tasks;
- discard irrelevant logs, noise, and dead-end reasoning;
- implicitly respect plan structure without explicit symbolic representations.

At inference time, the student can be called repeatedly at each step, providing fast, plan-aware compressions without running the evolutionary procedure or invoking a large teacher model. Because the student's behavior is distilled from a teacher across thousands of synthetic workflows, PAACE yields agents that maintain compact, relevant contexts over long horizons while preserving correctness, improving robustness to context dilution, and reducing inference cost. At deployment:

1. The agent obtains or constructs a plan $\Pi = [\tau_1, \ldots, \tau_n]$.
2. At step $t$, given the current context $C_t$ and a slice $\Pi_{t:t+k}$, the student compressor produces $\tilde{C}_t$. The agent executes task $\tau_t$ conditioned on $\tilde{C}_t$, and yields an updated context $C_{t+1}$.
3. Step 2 repeats until the plan completes, at which point a final answer is produced.

The teacher compression behavior is controlled entirely by a natural-language prompt $p$. Instead of hand-tuning $p$, PAACE meta-learns a population of prompts by treating context compression as a black-box policy and optimizing it via LLM-driven evolution over many workflows. Each prompt variant $p_i$ in the population is associated with statistics estimated from its evaluations:

- success rate: fraction of compressed trajectories labeled successful;
- mean semantic equivalence on successful trajectories;
- mean compression ratio over successful trajectories;
- a scalar reward that trades off fidelity and compression.

PAACE combines rankings by reward, success rate, mean equivalence, and mean ratio into a composite score used for selection. Evolution is run in a steady-state, asynchronous fashion. At any time, each worker process is assigned one prompt variant and evaluates it on several independent workflows, producing full- and compressed-context trajectories, trajectory labels (success/failure), and updated statistics for the variant. Prompt variants are updated online as new results arrive, and global reward statistics are continually recomputed. The teacher's compression policy improves over time as new prompts are proposed, evaluated, and selected based on real multi-step performance.

### 3.3. Supervision Quality, Scope, and Robustness

Because supervision is generated using an LLM teacher and evaluator, PAACE-Syn may inherit biases from the teacher's compression preferences. To mitigate this, we enforce strict outcome-level agreement with the full-context execution and discard any compression that degrades task success, even if favored by the judge. Empirically, this filtering reduces mode collapse and prevents overfitting to stylistic or heuristic compression artifacts. In addition, we emphasize that instruction co-refinement in PAACE is *implicit* rather than a separate optimization objective: instructions are rewritten only insofar as they are embedded in and reshaped by the compressed context. PAACE does not perform explicit instruction-only optimization; isolating instruction refinement as a standalone future-work objective.

While embedding similarity and LLM-based judging provide scalable and empirically robust proxies for task fidelity, they do not guarantee formal semantic equivalence or safety-critical correctness. PAACE thus targets functional preservation with respect to benchmark task outcomes rather than formal verification. Incorporating symbolic checks or domain-specific validators is an important direction for future work. Despite these limitations, using both embedding similarity and LLM-based judging substantially reduces degenerate compressions. Ablation on synthetic validation workflows shows that removing either filter increases failure rates by approximately 8–12%, primarily due to overly aggressive deletions or instruction drift that preserves surface similarity but breaks downstream task execution.

## 4. Experiments

We evaluate PAACE on two long-horizon agentic benchmarks that require multi-step reasoning, tool-use, and cross-application state tracking. These benchmarks provide diverse evaluation signal across planning, observation integration, memory demands, and long-context robustness. For each benchmark and model setting, we compare PAACE against baselines under a fixed execution backbone, planner, tool interfaces, and inference configuration. Within

each comparison group, the only difference between methods is the context management strategy. We emphasize that PAACE does not rely on any task-specific heuristics. All compression is learned from synthetic plan-conditioned demonstrations generated by PAACE-Syn and distilled into PAACE-FT. The **teacher compressor** is implemented using the OpenAI **GPT-OSS-120B** model with a 65,536-token context window. This model provides consistent multi-hop reasoning quality needed for next-$k$ compression. The **student compressor** for deployment, PAACE-FT, is distilled into **Qwen3–4B-Instruct**, due to its sufficiently large context window and fine summarization quality. To assess qualitative reproducibility without reliance on proprietary APIs, we conducted auxiliary sanity-check runs using the open-weight **GPT-OSS-20B** model. Due to its substantially smaller scale, these runs are not included in the quantitative comparison tables; however, they confirm that the qualitative effects of plan-aware next-$k$ compression and the relative ordering of methods persist across benchmarks.

PAACE-Syn consists of approximately 1.2M synthetic workflows comprising roughly 9.5B tokens before compression, generated entirely offline and amortized across all downstream tasks within benchmarks. We plan to release the full codebase, prompts, and a representative subset of the synthetic dataset upon publication to support reproducibility and follow-on research. Using the synthetic supervision produced by PAACE-Syn, the distilled compressor PAACE-FT preserves 97–98% of the teacher's performance while reducing inference cost by more than an order of magnitude. When plugged into the agent loop, PAACE-FT achieves nearly identical accuracy to the teacher-compressed version, confirming that next-$k$ compression behavior transfers well to smaller models. The teacher compressor incurs substantial overhead and is used only during offline data generation. Once distilled, PAACE-FT incurs no further optimization or teacher cost and can be reused across tasks within the same environment without additional prompt evolution or rollouts. At inference time, PAACE-FT adds less than 8% latency per step while reducing total input tokens by 35–60%, yielding net reductions in agent cost.

The primary goal of our experiments was to assess:

- **(RQ1)** Whether plan-aware next-$k$ compression improves agent correctness under long-context pressure.
- **(RQ2)** Whether PAACE reduces peak context length while retaining task fidelity.
- **(RQ3)** Whether the teacher–student PAACE-FT models preserve compression quality.

PAACE does not explicitly decompose pruning, rewriting, and summarization into independently parameterized operators; instead, these behaviors emerge jointly from the learned compression policy. As a result, operator-level ablations are ill-defined in our setting because pruning, rewriting,

and summarization emerge jointly from a single learned policy. However, we evaluate a constrained variant in which instruction rewriting is disabled and the compressor is restricted to extractive deletion only. This restriction leads to noticeably higher failure rates on multi-step tasks, indicating that implicit instruction reshaping contributes materially to long-horizon robustness. We evaluate the proposed method (PAACE) on three categories of long-horizon agentic tasks:

- **AppWorld** (Trivedi et al., 2024): multi-application tasks with heterogeneous observations.
- **OfficeBench** (Wang et al., 2024): document-centric tool chains and structured operations.
- **Multi-Objective QA** (Zhou et al., 2025): multi-hop retrieval QA with tool-based search.

Each benchmark measures both task success and interaction-dependent efficiency metrics (throughout this paper, *efficiency* refers to *token efficiency*—reductions in effective context length and cumulative attention dependency):

- **Acc / EM / F1**: benchmark-defined task performance.
- **Peak**: maximum input context length in the trajectory.
- **Dependency**: cumulative attention load measured as $\text{Dep} = \sum_{t=1}^{T} |C_t|$ where $|C_t|$ is the number of input tokens at step $t$. Reported in millions of tokens, this metric approximates total transformer attention cost and correlates with latency and quadratic compute growth. Because all methods are evaluated with identical agent backbones and step counts, Dependency is directly comparable across compression strategies, even when peak context lengths differ.

We compare PAACE against Acon (Yu et al., 2025b) and the following context-reduction baselines:

- **No Compression**: full interaction history preserved.
- **FIFO**: keep most recent $k$ turns, discard older.
- **Retrieval**: embedding-based past interaction selection.
- **LLMLingua**: extractive long-context compression.
- **Prompting**: an heuristic summarization instruction.

Tables 1, 2 and 3 summarize results on AppWorld, OfficeBench, and 8-Objective QA. PAACE consistently matches or outperforms all baselines in performance while substantially reducing context costs. PAACE surpasses the no-compression baseline on several benchmarks, suggesting that plan-aware compression can act as a form of regularization, reducing distraction from irrelevant or stale context. Across two long-horizon benchmarks, PAACE delivers consistently higher task performance, substantially lower context consumption, and more stable multi-step reasoning than all baselines. The results demonstrate the importance of incorporating plan structure, next-$k$ relevance, and workflow-level synthetic supervision into context optimization. PAACE enables long-horizon agents to maintain

*Table 1.* Average results on AppWorld with standard deviations reflecting Easy, Medium, and Hard. Across all benchmarks, PAACE maintains state fidelity and improves accuracy across multi-hop workflows while reducing context size and lowering cost.

| Method | Acc↑ | Steps↓ | Peak↓ | Dep↓ |
|---|---|---|---|---|
| No Compression | 56.00 | 16.14 | 9.93 | 5.96 |
| FIFO | 45.80 | 28.48 | 6.73 | 5.69 |
| Retrieval | 27.40 | 33.17 | 8.39 | 6.68 |
| LLMLingua | 39.30 | 24.42 | 7.50 | 6.37 |
| Prompting | 43.50 | 24.01 | 6.93 | 5.29 |
| ACON UT | 51.20 | 20.92 | 7.17 | 4.49 |
| ACON UTCO | 56.50 | 22.82 | 7.33 | 4.69 |
| **PAACE (ours)** | **59.00** | **19.20** | **6.23** | **3.75** |
| **Std.** | **± 9.68** | **± 5.92** | **± 0.89** | **± 1.60** |

*Table 2.* Results on OfficeBench.

| Method | Acc↑ | Steps↓ | Peak↓ | Dep↓ |
|---|---|---|---|---|
| No Compression | 76.84 | 11.52 | 7.27 | 4.43 |
| FIFO | 67.37 | 12.26 | 4.02 | 2.64 |
| Retrieval | 65.26 | 16.20 | 4.33 | 2.06 |
| LLMLingua | 70.53 | 10.89 | 4.65 | 1.85 |
| Prompting | 71.58 | 10.13 | 4.40 | 1.10 |
| ACON UT | 74.74 | 13.13 | 4.93 | 3.85 |
| ACON UTCO | 72.63 | 11.54 | 4.54 | 1.91 |
| **PAACE (ours)** | **78.10** | **10.48** | **4.29** | **1.64** |

coherent state representations without exceeding practical context budgets, even when deployed with compact compressors. To sum up, experimental results across AppWorld, OfficeBench, and Multi-Objective QA illustrate that *plan-aware context engineering* is a critical ingredient for robust long-horizon agents. PAACE consistently improves accuracy, lowers context cost, and stabilizes multi-step reasoning, indicating that context optimization should be treated as a core architectural module rather than an auxiliary utility. We note that improvements over the strongest baseline (ACON UTCO) are sometimes modest relative to the observed variance. These results should thus be interpreted as evidence that PAACE does not harm task performance while substantially reducing context cost, rather than as demonstrating a large absolute accuracy gain on this benchmark. Table 4 shows that moderate lookahead ($k = 2$) is sufficient for tool-centric benchmarks, while multi-hop QA benefits from a longer relevance horizon ($k = 3$) because retrieved evidence is often consumed several steps after acquisition. Larger $k$ values yield diminishing returns due to increased compression difficulty, indicating that next-$k$ conditioning should be selected based on task dependency structure.

A central finding is that agents benefit not only from retaining relevant information but also from *removing* distracting and stale information. While long-context architectures

Table 3. Average results on 8-Objective QA.

| Method | EM↑ | F1↑ | Steps↓ | Peak↓ | Dep↓ |
|---|---|---|---|---|---|
| No Compression | 0.366 | 0.488 | 15.78 | 10.35 | 3.32 |
| FIFO | 0.293 | 0.388 | 19.26 | 5.09 | 2.51 |
| Retrieval | 0.331 | 0.438 | 20.06 | 5.11 | 2.62 |
| LLMLingua | 0.363 | 0.481 | 17.68 | 5.68 | 2.24 |
| Prompting | 0.376 | 0.478 | 18.70 | 4.73 | 1.66 |
| ACON UT | 0.373 | 0.494 | 17.14 | 4.71 | 1.57 |
| ACON UTCO | 0.335 | 0.458 | 17.79 | 4.65 | 1.50 |
| **PAACE (ours)** | **0.402** | **0.512** | **16.86** | **4.41** | **1.41** |

Table 4. Ablation on next-$k$ conditioning for PAACE.

| Benchmark | $k$ | Acc / EM↑ | F1↑ | Peak↓ | Dep↓ |
|---|---|---|---|---|---|
| AppWorld | 1 | 56.5 | – | 6.05 | 3.95 |
| | **2** | **59.0** | – | **6.23** | **3.75** |
| | 3 | 58.6 | – | 6.48 | 3.82 |
| OfficeBench | 1 | 76.3 | – | 4.15 | 1.72 |
| | **2** | **78.1** | – | **4.29** | **1.64** |
| | 3 | 77.6 | – | 4.51 | 1.69 |
| 8-Objective QA | 1 | 0.381 | 0.497 | 4.18 | 1.36 |
| | 2 | 0.394 | 0.505 | 4.30 | 1.39 |
| | **3** | **0.402** | **0.512** | **4.41** | **1.41** |

push token windows into the hundreds of thousands or millions, our experiments show that model performance still degrades when exposed to ill-structured or semantically dilute histories. This aligns with emerging empirical evidence on context "rot," where transformers deteriorate as they attend over long, heterogeneous sequences. PAACE demonstrates that explicit modeling of *plan structure* and *next-$k$-step relevance* yields substantial benefits. Unlike query-aware compression methods that optimize for a single upcoming question, PAACE preserves information required for multi-step chains, such as tool-call sequences, document references, and causal dependencies. PAACE'd contexts serve as coherent state representations, avoiding the loss of global task structure that often impairs baselines.

A notable contribution of PAACE is the effectiveness of synthetic workflows produced by PAACE-Syn. Existing datasets for agentic research either lack plan annotations or are too limited in size. In contrast, PAACE-Syn generates millions of diverse traces with explicit plan structures, multi-hop dependencies, variable noise, and rich tool-use patterns. These workflows provide dense supervision for training next-$k$–aware compressors without requiring expensive human annotation. The strong performance of the distilled PAACE-FT—recovering up to $97\%$ of the teacher's behavior—demonstrates that *synthetic supervision can be reliably transferred to compact models*. For practical deployment of PAACE, this result is significant: the student compressor can be executed cheaply at every step of an agent loop, removing dependence on a large teacher LLM.

Another interesting observation is that PAACE often surpasses the "no compression" baseline. This suggests that compression is not merely an efficiency tool but also a form of *regularization*. Removing irrelevant or contradictory information forces the agent to reason over a more coherent, structured state, reducing distraction and mitigating the risk of attending to spurious details. This may explain the improved task success on OfficeBench and 8-Objective QA, where unfiltered histories contain verbose logs and redundant tool outputs. Our results indicate that PAACE significantly mitigates instruction drift—a common failure mode in long-running workflows. By conditioning compression on the next-$k$ tasks, PAACE preserves alignment between the agent's plan and the evolving state. This ensures that rewritten instructions, retained constraints, and contextual details remain consistent across steps, reducing the accumulation of cascading errors.

## 5. Conclusion

We introduced PAACE, a framework for plan-aware automated context engineering in long-horizon LLM agents. By modeling context as a structured, plan-dependent state and optimizing it via outcome-level supervision—through relevance scoring, rewriting, summarization, pruning, and instruction co-refinement—PAACE enables agents to reason more robustly under tight context budgets. Across multiple benchmarks, PAACE consistently improves task correctness, reasoning stability, and effective context utilization while substantially reducing peak context length. These results show that learned, plan-aware context shaping can serve as a first-class component of agent architectures, akin to retrieval augmentation in knowledge-intensive systems.

PAACE targets robustness and context efficiency, not guaranteed reductions in total API cost, which depend on deployment and compression frequency. We make no claims of convergence, optimality, or formal guarantees; rather, our results demonstrate that outcome-driven learning suffices to induce robust and transferable compression behavior in realistic multi-step workflows. PAACE intentionally learns environment- and workflow-specific compression policies. While cross-domain generalization is an important direction for future work, our findings suggest that plan-aware relevance is inherently task-dependent in long-horizon settings, and that such specialization is a feature—not a limitation—of effective context engineering.

Overall, PAACE frames context engineering as a learnable, plan-aware optimization problem and provides a practical and empirical foundation for building scalable, reliable agents that operate coherently over extended interactions.

## Broader Impact Statement

From a broader impact perspective, improved context efficiency can enable the deployment of long-horizon agents under tighter resource, latency, and cost constraints, potentially widening access to complex agentic systems in practical settings. At the same time, aggressive or poorly validated compression policies could omit safety-relevant instructions, constraints, or provenance information if applied outside the conditions under which they were trained. PAACE mitigates this risk by requiring outcome-level agreement with full-context execution and by discarding compressions that degrade task success. Nevertheless, PAACE is not designed as a safety or alignment mechanism, and its use in safety-critical or high-stakes domains would require additional validation, domain-specific checks, or human oversight.

## References

Ahn, M. et al. Mrkl systems: Modular reasoning, knowledge and language. *arXiv preprint arXiv:2205.00445*, 2022.

Asai, A., Chen, X., and Hajishirzi, H. Self-rag: Learning to retrieve, generate, and critique through self-reflection. *arXiv preprint arXiv:2403.14580*, 2024.

Chung, H. W., Hou, L., Longpre, S., et al. Scaling instruction-finetuned language models. *arXiv preprint arXiv:2210.11416*, 2022.

Gao, L., Madaan, A., et al. Pal: Program-aided language models. In *ICML*, 2023.

Guo, M. et al. Longt5: Efficient text-to-text transformer for long sequences. In *arXiv preprint arXiv:2112.07916*, 2021.

Huang, J. et al. Planning with large language models via iterative refinement. In *NeurIPS*, 2024.

Jiang, H., Ye, F., et al. Llmlingua: Efficient context compression for large language models. *arXiv preprint arXiv:2310.07041*, 2023.

Jiang, H. et al. Llmlingua-2: Efficient and robust context compression for large language models. *arXiv preprint arXiv:2404.17762*, 2024.

Joo, T., Ishida, S., Sosnovik, I., Lim, B., Rezaei-Shoshtari, S., Gaier, A., and Giaquinto, R. Graph of agents: Principled long-context modeling by emergent multi-agent collaboration. *arXiv preprint arXiv:2509.21848*, 2025. URL https://arxiv.org/abs/2509.21848.

Lewis, M. et al. Bart: Denoising sequence-to-sequence pre-training. In *ACL*, 2019.

Li, J., Zhou, X., Zhang, Y., Zhao, R., and Chen, Y. Compass: Controllable memory policy for long-context language models. *arXiv preprint arXiv:2402.05618*, 2024.

Liu, X., Guo, J., et al. Agentbench: Evaluating llm agents via simulated environments. In *NeurIPS Datasets and Benchmarks Track*, 2023. URL https://arxiv.org/abs/2308.03688.

Ouyang, L. et al. Training language models to follow instructions with human feedback. *arXiv preprint arXiv:2203.02155*, 2022.

Park, J. et al. Generative agents: Interactive simulacra of human behavior. *arXiv preprint arXiv:2304.03442*, 2023.

Pryzant, R. et al. Automated prompt engineering. In *NeurIPS*, 2023.

Rae, J. et al. Compressive transformers for long-range sequence modeling. *ICLR*, 2020.

Rafailov, R. et al. Direct preference optimization: Optimizing language models without reinforcement learning. *arXiv preprint arXiv:2305.18290*, 2023.

Raffel, C. et al. Exploring the limits of transfer learning with a unified text-to-text transformer. In *Journal of Machine Learning Research*, 2020.

Saunders, W. et al. Preference ranking optimization. *arXiv preprint arXiv:2402.07872*, 2024.

Schick, T., Dwivedi-Yu, J., et al. Toolformer: Language models can teach themselves to use tools. In *International Conference on Learning Representations (ICLR)*, 2023a. URL https://arxiv.org/abs/2302.04761.

Schick, T. et al. Toolformer: Language models can teach themselves to use tools. In *ICLR*, 2023b.

See, A. et al. Get to the point: Summarization with pointer-generator networks. *ACL*, 2017.

Shinn, N. et al. Reflexion: Language agents with verbal reinforcement learning. *arXiv preprint arXiv:2303.11366*, 2023.

Sun, S., Liu, Z., Chen, X., and Li, J. Prompt compression for large language models. *arXiv preprint arXiv:2310.06201*, 2023.

Sun, W., Lu, M., Ling, Z., Liu, K., Yao, X., Yang, Y., and Chen, J. Scaling long-horizon llm agents via context-folding. *arXiv preprint arXiv:2510.11967*, 2025. URL https://arxiv.org/abs/2510.11967.

Trivedi, H., Khot, T., Hartmann, M., Manku, R., Dong, V., Li, E., Gupta, S., Sabharwal, A., and Balasubramanian, N. Appworld: A controllable world of apps and people for benchmarking interactive coding agents. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 16022–16076, Bangkok, Thailand, August 2024. Association for Computational Linguistics. URL https://doi.org/10.18653/v1/2024.acl-long.850.

Wang, Z., Cui, Y., Zhong, L., Zhang, Z., Yin, D., Lin, B. Y., and Shang, J. Officebench: Benchmarking language agents across multiple applications for office automation. *arXiv preprint arXiv:2407.19056*, 2024. URL https://doi.org/10.48550/arXiv.2407.19056.

Wei, J. et al. Chain-of-thought prompting elicits reasoning in large language models. In *NeurIPS*, 2022.

Wu, Y. et al. Provence: Learned context pruning for retrieval-augmented generation. *arXiv preprint arXiv:2403.16975*, 2024.

Yao, S., Bosma, D., Zhao, J., et al. React: Synergizing reasoning and acting in language models. In *International Conference on Learning Representations (ICLR)*, 2023a. URL https://arxiv.org/abs/2210.03629.

Yao, S. et al. React: Synergizing reasoning and acting in language models. In *ICLR*, 2023b.

Yao, S. et al. Tree of thoughts: Deliberate problem solving with large language models. In *arXiv preprint arXiv:2305.10601*, 2023c.

Yu, B. et al. Memagent: Learned memory mechanisms for long-context language models. *arXiv preprint arXiv:2501.01234*, 2025a.

Yu, Y., Li, H., Zheng, J., et al. Agent context optimization: Compressing agent interaction histories via natural language guidance. *arXiv preprint arXiv:2501.00001*, 2025b.

Zaheer, M. et al. Big bird: Transformers for longer sequences. *NeurIPS*, 2020.

Zhang, J. et al. Pegasus: Pre-training with extracted gap sentences for abstractive summarization. *ICML*, 2019.

Zhang, T. et al. React in the wild: Reasoning-augmented tool use. *arXiv preprint arXiv:2306.03631*, 2023.

Zhang, Z. et al. Recap: Recursive context-aware reasoning and planning for large language model agents. *arXiv preprint arXiv:2510.23822*, 2025. URL https://arxiv.org/abs/2510.23822.

Zhou, J. et al. Large language model prompt optimization. *arXiv preprint arXiv:2309.03409*, 2023a.

Zhou, S. et al. Large language models as tool makers. *arXiv preprint arXiv:2304.03439*, 2023b.

Zhou, Y. et al. Recursive chain of thought: Multi-step reasoning with structured steps. *arXiv preprint arXiv:2311.11543*, 2023c.

Zhou, Z., Zhang, Z., et al. Webarena: A realistic web environment for building autonomous agents. In *International Conference on Learning Representations (ICLR)*, 2024. URL https://arxiv.org/abs/2307.15238.

Zhou, Z., Qu, A., Wu, Z., Kim, S., Prakash, A., Rus, D., Zhao, J., Low, B. K. H., and Liang, P. P. MEM1: Learning to synergize memory and reasoning for efficient long-horizon agents. *arXiv preprint arXiv:2506.15841*, 2025. URL https://doi.org/10.48550/arXiv.2506.15841.