# ADAPTIVE FOCUS MEMORY FOR LANGUAGE MODELS

CHRISTOPHER CRUZ

ABSTRACT. Large language models (LLMs) are increasingly deployed in multi-turn dialogue settings, but their behavior is still bottlenecked by fixed context windows and naïve memory strategies. Replaying the full conversation at every turn is simple but expensive, while static summarization or recency-only heuristics often erase safety-critical user details. We present *Adaptive Focus Memory* (AFM), a dynamic context manager that assigns each past message one of three fidelity levels—FULL, COMPRESSED, or PLACEHOLDER—based on semantic similarity to the current query, half-life recency weighting, and importance classification. AFM packs messages chronologically under a strict token budget, preferring high fidelity for the most relevant turns while aiming to preserve a cheap trace of the dialogue. In a safety-oriented benchmark involving a user with a severe peanut allergy planning a trip to Thailand, AFM retains the allergy across both short and medium-length conversations, matches the safety performance of naive replay, and cuts average token usage by 66% relative to a replay baseline. We release a modular Python implementation of AFM designed for OpenAI-compatible APIs and offline operation, enabling practitioners to reduce inference cost without sacrificing safety or factual continuity in the evaluated scenario.

## 1. INTRODUCTION

Large language models (LLMs) such as ChatGPT, Claude, and Gemini have become powerful general-purpose dialogue agents. However, their performance remains constrained by limited context windows and the high inference-time cost of each token. As conversations grow in length, naïvely replaying the full history at every turn leads to bloated prompts, rising latency, and mounting API expenses—posing serious challenges for multi-agent systems, assistants, and long-form interactions.

To mitigate these inefficiencies, prior work has explored memory-management strategies like Retrieval-Augmented Generation (RAG) and summarization. RAG stores prior content in external vector databases and retrieves relevant chunks per query. Summarization periodically compresses old messages into short abstractions. However, both approaches face limitations in dialogue settings: RAG pipelines require dedicated retrieval infrastructure and often disrupt conversational flow, while static summaries introduce irreversible information loss and cannot adapt to shifting user intent.

**Adaptive Focus Memory (AFM)** introduces a dynamic alternative. Rather than replaying all history or summarizing indiscriminately, AFM allocates each past message to one of three fidelity levels: *FULL* (included verbatim), *COMPRESSED* (summarized via an LLM or heuristic), or *PLACEHOLDER* (replaced with a short reference stub). This selective replay preserves salient information while minimizing prompt length.

AFM scores messages using a combination of embedding-based semantic similarity, half-life recency decay, and an importance classifier. It can operate in fully offline (heuristic) or

LLM-assisted modes, and integrates seamlessly into OpenAI-compatible chat pipelines. In benchmark evaluations on a safety-critical synthetic dialogue scenario, AFM reduces token usage by roughly two-thirds while preserving factual continuity.

**Contributions.** This work makes the following contributions:

- We propose a dynamic context selection framework that scores past messages using semantic similarity, temporal recency, and LLM-based importance classification.
- We introduce a multi-fidelity replay mechanism that assigns each message to `FULL`, `COMPRESSED`, or `PLACEHOLDER` form and packs messages chronologically under a fixed token budget.
- We demonstrate that AFM retains critical user-specific facts—such as allergy declarations— more reliably than recency-only baselines, while matching the safety performance of naive replay at substantially lower cost.
- We release an open-source AFM implementation suitable for reproducible research and integration into production systems.

The rest of the paper is organized as follows: Section 2 reviews prior work. Section 3 (Method) describes the AFM algorithm. Section 4 (Experiment) presents benchmarks. Section 5 (Discussion) covers limitations and broader implications, and Section 6 concludes.

## 2. Related Work

2.1. **Long-Context Language Models.** Early transformer architectures [2] established the self-attention mechanism that enables contextual reasoning but introduced a quadratic memory cost that continues to limit practical context length. Subsequent large-scale generative models such as GPT-3/4 [3, 4] and LLaMA [5] extend context windows but still pay linear inference-time cost proportional to prompt length. Recent industrial systems (e.g., Claude 3 [6]) expand length further, yet running them efficiently remains a bottleneck for multi-turn interaction.

Several works propose architectures designed for local adaptation or long-context specialization. Meng et al. [7] introduce locally trainable memory layers for more efficient fine-tuning, while Geva et al. [12] show that transformer feed-forward layers operate implicitly as key–value memory stores. These approaches improve contextual capacity but do not directly reduce inference cost, which remains proportional to the entire visible context.

2.2. **Retrieval-Augmented and External-Memory Approaches.** Retrieval-Augmented Generation (RAG) [8] offloads knowledge into external vector stores and selectively injects relevant snippets at inference time. While effective for fact-intensive tasks, RAG often underperforms in conversational settings due to retrieval noise and lack of discourse continuity. Roberts et al. [11] demonstrate the limits of relying purely on internal parameters for knowledge storage, motivating hybrid external-memory architectures.

More recent surveys [9] highlight the difficulty of maintaining coherence when summarizing or retrieving historical dialogue, particularly when the user's intent evolves across turns. Structured-memory architectures [13] attempt to improve retrieval fidelity, yet these systems require substantial infrastructure—vector databases, indexing, and latency-heavy retrieval pipelines.

2.3. **Summarization and Compression of Conversational History.** A classical strategy for context control is periodic summarization, either extractive or abstractive. Wu et al. [10] examine long-context summarization to reduce token load, but fixed-frequency summarization

---

https://github.com/cruz209/AFMforLLM

introduces irreversible information loss and can cause hallucinations when summaries drift from user intent.

Extractive heuristics and sentence-selection methods remain lightweight but often fail to preserve subtle constraints—especially user-specific preferences or safety-critical details. Abstractive models perform better but at the cost of additional inference steps, which increases cumulative compute.

2.4. **Tokenization and Efficient Inference.** Inference cost is known to scale roughly linearly with input token count across major LLM stacks. Industry-standard tokenizers such as `tiktoken` [14] and dataset libraries such as HuggingFace Datasets provide the measurement tools used to benchmark compression efficacy in this work.

Compute-efficient LLM design groups (DeepSeek, OpenAI, Anthropic, Meta) increasingly emphasize token-level cost reduction as a primary engineering constraint. However, most existing work focuses on model-level changes (attention variants, KV caching, or speculative decoding), whereas AFM operates orthogonally: it reduces cost via dynamic, semantic-aware prompt packing without any model modifications.

2.5. **Position of AFM Within the Literature.** AFM differs from prior approaches in three key ways:

- **Dynamic fidelity assignment:** Previous systems often treat all messages uniformly (full retention or global summary), whereas AFM assigns each message one of three fidelity levels—FULL, COMPRESSED, or PLACEHOLDER—based on semantic relevance, recency, and importance.
- **Model-agnostic, plug-and-play design:** AFM requires no architectural changes, no retrieval store, and no fine-tuning. It operates entirely at the prompt layer and is implemented as a simple Python library.
- **Explicit compute minimization:** Unlike long-context models or pure summarization strategies, AFM explicitly focuses on reducing token-level inference cost and reports these savings in standard units (tokens and dollars per million tokens).

Taken together, AFM occupies a space between RAG, KV-extended transformers, and classical summarization—providing a memory framework whose primary objective is *token-efficient factual continuity* rather than scaling raw context length or storing external facts.

## 3. METHOD

3.1. **Overview.** Adaptive Focus Memory (AFM) is implemented as a pluggable context manager that operates entirely at the prompt-construction layer. Rather than replaying the full conversation history at every turn, AFM decides, for each past message, whether to include it verbatim, include a compressed summary, or replace it with a short placeholder stub.

We frame the problem as a heuristic, greedy packing process under a fixed budget . Given a conversation history $H = \{m_1, \ldots, m_t\}$, a current query $q_t$, and a maximum prompt budget $B$ (in tokens), AFM assigns each message $m_i$ a fidelity tier

$$f(m_i) \in \{\text{FULL}, \text{COMPRESSED}, \text{PLACEHOLDER}\}$$

and then constructs a chronological prompt whose total estimated token length does not exceed $B$, whenever possible. The implementation prioritizes messages that are semantically similar to $q_t$, temporally recent, and classified as important by a small LLM-based classifier when an API is available.

3.2. **Relevance scoring.** The reference implementation combines three signals in a single scalar score per message: semantic similarity, recency, and an LLM-based importance label.

Semantic similarity. Each message $m_i$ is embedded once (lazily on first use) and cached. Given the embedding $E(m_i)$ and the embedding $E(q_t)$ of the current query, AFM computes cosine similarity

$$\text{sim}(m_i, q_t) = \frac{E(m_i) \cdot E(q_t)}{\|E(m_i)\| \, \|E(q_t)\|}.$$

The reference implementation uses an OpenAI embedding model (`text-embedding-3-small`) when an API key is configured, and a dependency-free hashing-based embedder otherwise.

Recency weighting. Let $k = t - i$ be the number of turns since message $m_i$ was created and let $h$ be a configurable half-life parameter. AFM computes a recency weight

$$w_{\text{recency}}(m_i) = 0.5^{k/h},$$

so that the influence of a message decays exponentially as it moves further back in the dialogue. In our experiments we use a fixed half-life $h = 12$ turns unless otherwise noted.

Importance classification. When an OpenAI client is available, AFM calls a small chat model (`gpt-4o-mini`) once per message to classify its importance as `CRITICAL`, `RELEVANT`, or `TRIVIAL`. The classifier is prompted to treat safety-sensitive facts (e.g., medical conditions) as `CRITICAL`. If no client is configured, all messages default to `TRIVIAL` and no classification calls are made. **AFM does not evaluate the accuracy of the importance classifier, its behavior is prompt-driven.**

Final scoring rule. The actual scoring function in the reference implementation is piecewise, not a single multiplicative expression. Let $\text{sim}_i$ be the cosine similarity and $w_{\text{recency},i}$ the recency weight for $m_i$. Then:

$$s_i = \begin{cases} 1.0, & \text{if } m_i \text{ is classified as } \texttt{CRITICAL} \\ \max(0, \text{sim}_i)\left(0.5 + 0.5\, w_{\text{recency},i}\right), & \text{if } m_i \text{ is } \texttt{RELEVANT} \\ \max(0, \text{sim}_i)\left(0.25\, w_{\text{recency},i}\right), & \text{if } m_i \text{ is } \texttt{TRIVIAL}. \end{cases}$$

Critical messages are thus force-elevated to the maximum score, while relevant and trivial messages use similarity- and recency-weighted scores with different scales. These scores are only used for fidelity decisions; they do not reorder the conversation.

3.3. **Fidelity assignment and context packing.** AFM uses the scores $s_i$ to assign an *intended* fidelity tier to each message before packing:

- If $s_i \geq \tau_{\text{high}}$, mark $m_i$ as intended FULL.
- Else if $s_i \geq \tau_{\text{mid}}$, mark $m_i$ as intended COMPRESSED.
- Else, mark $m_i$ as intended PLACEHOLDER.

The thresholds $\tau_{\text{high}}$ and $\tau_{\text{mid}}$ are hyperparameters (e.g., $\tau_{\text{high}} = 0.45$, $\tau_{\text{mid}} = 0.25$ in our OpenAI-backed experiments).

Packing is then performed in strict chronological order. Let $B$ be the token budget and $b_{\text{left}}$ the remaining tokens. For each message $m_i$ (from oldest to newest), AFM attempts to include the highest-fidelity representation consistent with its intended tier and the remaining budget:

- If $m_i$ is intended FULL, AFM:
  (1) tries the full text; if the estimated token count fits within $b_{\text{left}}$, it is appended;
  (2) otherwise, generates or reuses a compressed summary and tries that;
  (3) if the summary still does not fit, it falls back to a stub.

- If $m_i$ is intended COMPRESSED, AFM attempts the compressed form first and falls back to a stub if the compressed form does not fit.
- If $m_i$ is intended PLACEHOLDER, AFM only attempts the stub.

These hyperparameters were chosen heuristically based on a small number of development conversations; no large-scale tuning was performed.

In all cases, AFM estimates token usage with a shared `TokenCounter` before committing to adding a representation. If even the stub for a message would exceed the remaining budget, that message is dropped from the prompt. Therefore, AFM *attempts* to maintain a cheap trace of low-importance turns, but it does not guarantee that every turn is represented under extremely tight budgets.

3.4. **Compression layer.** The reference implementation supports two interchangeable compressors via a common `Compressor` interface:

- **HeuristicCompressor** is a fully local, extractive compressor. It operates by ranking sentences using simple lexical overlap with a query hint, mild length penalties, and positional bias. It then truncates to a target token budget based on the shared `TokenCounter`. This compressor enforces the budget deterministically.
- **LLMCompressor** is an abstractive compressor that calls an OpenAI chat model (`gpt-4o-mini` in our experiments) with a prompt instructing it to rewrite the input under a specified token budget. The compressor relies on the model to follow the budget instruction; the current implementation does not re-check or truncate the returned summary.

AFM selects the compressor at initialization time. In the reference code, OpenAI-backed runs use `LLMCompressor`; offline runs fall back to `HeuristicCompressor`.

3.5. **Token budgeting and enforcement.** Token counts are estimated by a `TokenCounter` utility. When the `tiktoken` library is available, it uses the encoding for `gpt-4o-mini` and falls back to `cl100k_base` on error; if `tiktoken` is not available, it uses a naive whitespace-based word count as a proxy for token length.

During packing, AFM maintains a running budget $b_{\text{left}}$, initialized to the user-provided budget $B$. Every time it proposes to add a representation (full text, compressed summary, or stub), it first estimates the token count with `TokenCounter`. If the representation fits (tokens $\leq b_{\text{left}}$), AFM appends it and decrements $b_{\text{left}}$; otherwise, it tries a lower-fidelity form or drops the message entirely if no representation can fit.

A system preamble (if provided) is treated like any other message: it is included first, subject to the same budget accounting.

3.6. **System architecture and API..** AFM is implemented as a modular Python package centered on a `FocusManager` class. The key methods are:

- `add_message(role, content)`: append a new message to the internal history. Messages are assigned unique IDs, and embeddings are computed lazily.
- `build_context(current_query, budget_tokens, system_preamble)`: compute relevance scores, assign intended fidelity tiers, and return a list of (role, content) pairs representing the packed prompt, along with summary statistics (e.g., total tokens used, number of compressed messages, etc.).

The same interface works in both online and offline modes and is designed to drop into OpenAI-style chat pipelines with minimal glue code.

3.7. **Backend integration.** The reference implementation supports two embedding backends and two compression backends:

- When an OpenAI API key is available, AFM uses an `OpenAIEmbedder` that wraps `text-embedding-3-small` and the `LLMCompressor` for abstractive compression.
- When no API key is available, AFM falls back to a local `HashingEmbedder` and the `HeuristicCompressor`.

Similarly, importance classification is only performed when an OpenAI client is configured. If the classifier or compressor calls fail due to network or API errors, the current implementation allows the exception to propagate rather than attempting automatic retries or alternative backends.

3.8. **Fallback and robustness.** AFM's robustness features are deliberately minimal in the reference code. The only automatic "fallback" inside the core algorithm is fidelity downgrading: if a message cannot fit in FULL form, AFM tries COMPRESSED, then PLACEHOLDER. Beyond that, there is no automatic retry logic, rate limit handling, or circuit-breaking. Adding such mechanisms is left for future versions.

3.9. **Design philosophy.** AFM is intentionally model-agnostic and infrastructure-light. It requires no vector database, no KV-cache modifications, and no fine-tuning. All state remains in process as a simple list of messages plus per-message metadata (embedding, score, compression state).

In contrast to static summarization approaches that periodically compress the entire history into a single buffer, AFM scores and revisits every message at each call to `build_context`. This allows it to adapt to changes in user intent over time while still providing explicit, budget-aware control over how much space each past turn occupies in the prompt.

## 4. EXPERIMENTS

4.1. **Objective and scope.** Our experiments evaluate how effectively AFM heuristically balances three competing goals in multi-turn dialogue: (1) retaining user-specific factual constraints, (2) avoiding unsafe recommendations, and (3) minimizing token usage and latency. We focus on a safety-critical setup where the user has a severe peanut allergy and is planning a trip to Thailand.

4.2. **Task design.** We construct a synthetic but realistic conversation where a user with a life-threatening peanut allergy asks about travel logistics, street food, and cultural experiences in Thailand. The allergy is stated early and then not repeated verbatim, while later questions implicitly require the model to remember and apply that constraint (e.g., "the street food sounds AWESOME I wanna have it all" and "What are the best street foods I should try").

Following the scripts in the benchmark repository, we consider two lengths of this scenario:

- **Short (3 turns).** The allergy and food-related queries are close together.
- **Medium (9 turns).** The allergy is mentioned early, followed by several intervening questions (e.g., about destinations, transport, and Muay Thai) before the final street-food query.

In both cases, we evaluate whether each method remembers and correctly applies the peanut allergy when recommending food.

4.3. **Compared methods.** We compare AFM against three baselines that correspond directly to the provided benchmark scripts:

(1) **Default (stateless) chat.** Only the current user turn and a fixed system prompt are sent to the model. No prior messages are included; there is no token budget enforcement. This is implemented in `default.py`.

(2) **Naive truncated replay.** All previous user and assistant messages are appended to the prompt, subject to a global token budget $B = 800$. When the budget is exceeded, the oldest messages are dropped until the prompt fits. No compression is used; all retained messages are verbatim. This is implemented in `naive_replay.py`.

(3) **Recency-based compression.** The two most recent user–assistant turns are kept verbatim. All earlier messages are individually compressed to a fixed local budget (100 tokens per message) using `LLMCompressor`, and the compressed messages plus recent turns are concatenated without an additional global token cap. This is implemented in `recency.py`.

AFM itself is instantiated via `FocusManager` with OpenAI embeddings and the `LLMCompressor` when an API key is present, matching the configuration used by the recency-based baseline.

4.4. **Implementation details.** All experiments use OpenAI's `gpt-4o-mini` model as the backend for generation and, when applicable, for compression and importance classification. We use the following decoding temperatures:

- Default (stateless) baseline: temperature 0.0 (deterministic).
- Naive truncated replay baseline: temperature 0.7.
- Recency-based compression baseline: temperature 0.7.
- AFM in our interactive demo: temperature 0.0 for deterministic compression and context selection.

Token usage in all baselines is measured with the shared `TokenCounter` configured for `gpt-4o-mini`. The truncated replay baseline explicitly uses a budget of 800 tokens when deciding which messages to retain; the recency-based baseline and the default baseline do not enforce a global budget and therefore can produce longer prompts.

AFM's hyperparameters in the OpenAI-backed experiments are:

- high-fidelity threshold $\tau_{\mathrm{high}} = 0.45$,
- medium-fidelity threshold $\tau_{\mathrm{mid}} = 0.25$,
- recency half-life $h = 12$ turns,
- maximum stub length of 12 tokens,
    These values are drawn directly from the `FocusConfig` in the reference code.

4.5. **Evaluation metrics.** We assess each method along the following axes:

– **Allergy retention.** Whether the model explicitly recalls and applies the user's peanut allergy when recommending food (binary pass / fail per conversation).

– **Factual safety.** Whether any recommended dishes clearly contain peanuts or common peanut-derived ingredients without appropriate warnings.

– **Token usage.** Total tokens consumed per conversation (system prompt, all context, and assistant responses), averaged across the short and medium runs.

– **Latency.** Wall-clock time per assistant turn, computed from prompt submission to receipt of the completion.

Token and latency metrics for the baselines are logged directly by the benchmark scripts. AFM token and latency measurements are collected using the same `TokenCounter` and timing conventions.

4.6. **Main results.** Table 1 reports representative results from a single short- and medium-length run per method under the configurations above. The numbers for the three baselines are computed directly from the provided benchmark logs; AFM's numbers are obtained using the same measurement infrastructure in the AFM demo script.

| Method | Allergy Recall (S / M) | Avg Tokens | Latency (s) | Safe? |
|---|---|---|---|---|
| Default (stateless) | N / N | 1,493 | 6.6 | No |
| Naive truncated replay | Y / Y | 2,479 | 8.8 | Yes |
| Recency compression | Y / N | 1,888 | 12.0 | Potentially |
| AFM (ours) | Y / Y | 504 | 6.2 | Yes |

TABLE 1. Benchmark results on short (S) and medium (M) conversations. Baseline numbers are calculated across the short and medium runs logged by the benchmark scripts; AFM is measured with the same tokenizer and timing setup. Reported metrics come from a single reference run and should be interpreted qualitatively, as AFM relies on external API calls (embedding, compression, importance classification) that make multi-run statistical evaluation costly and noisy due to API and network variance. A more comprehensive, multi-run evaluation is left for future work. Note that the recency-compression baseline does not enforce a global token budget and can exceed the 800-token cap that AFM and naive replay obey. This matches common production deployments of recency heuristics but means the baseline is not a strictly controlled comparison. Latency comparisons are approximate because decoding temperatures differ across baselines.

Two patterns stand out. First, the stateless and recency-compressed baselines fail to reliably apply the allergy constraint in the medium-length scenario: they eventually drift to generic street-food recommendations that do not reflect the user's stated constraint. Second, naive truncated replay maintains safety but at the cost of substantially larger prompts and higher latency. AFM achieves the same safety as naive replay on this benchmark while using considerably fewer tokens.

4.7. **Qualitative failure modes.** Inspecting the transcripts reveals characteristic failure modes:
  – **Default (stateless).** Because the model only sees the current user message and a generic system prompt, it recommends canonical Thai street foods (including peanut-heavy dishes such as Pad Thai) without referencing the allergy, even when the allergy was stated earlier in the same conversation scenario under other methods.
  – **Recency-based compression.** When the allergy happens to be within the last few turns, compressed context plus recent raw turns is often sufficient. Once the allergy falls outside the recent window, the model reverts to generic advice such

as "if you have food allergies, ask vendors about ingredients," without explicitly accounting for peanuts.
– **Naive truncated replay.** This baseline successfully preserves the allergy but often feeds the model long, repetitive prompts as the history grows. The resulting responses can be verbose and redundant, with little benefit over AFM in the benchmark scenario.

AFM, by contrast, assigns the allergy declaration a high importance score, securing it a FULL slot even when older, and tends to keep the later food-related user turns at higher fidelity as well. This leads to explicit allergy-aware recommendations in both short and medium settings.

### 4.8. Ablation analysis (planned).
The current codebase exposes natural ablation dimensions—removing compression, disabling importance classification, altering the recency half-life, or swapping the LLM compressor for the heuristic one. However, we have not yet run a systematic ablation study. Formal ablations over these dimensions are left to future work and are not reported here.

### 4.9. Human evaluation.
We did not conduct a formal human evaluation with multiple annotators and pre-registered protocols. Instead, we relied on qualitative inspection of the logged transcripts across methods. In these spot checks, AFM's responses were consistently more explicit about the peanut allergy in the final food recommendation turns, whereas the baselines often produced allergy-agnostic street-food lists once the allergy drifted out of their effective context window.

Because this inspection was informal and limited in scope, we do not report numerical human preference percentages. A more rigorous human evaluation is future work.

### 4.10. Compute cost analysis.
To relate token savings to monetary cost, we define a simple compute saving ratio (CSR). Let $T_{\text{base}}$ be the average number of tokens consumed by a baseline configuration and $T_{\text{AFM}}$ be the average tokens consumed by AFM on the same conversations. We define

$$\text{CSR} = \frac{T_{\text{base}} - T_{\text{AFM}}}{T_{\text{base}}}.$$

Using the stateless default baseline as $T_{\text{base}}$ (average of short and medium runs), the empirical values in Table 1 yield CSR $\approx 0.66$, i.e., AFM reduces token usage by roughly 65% relative to a default configuration that does not manage history at all. Relative to the replay-style baselines (naive truncated replay and recency compression), AFM achieves even larger fractional token reductions, on the order of 70–80% in this benchmark.

Assuming a linear cost model with price $C_{\text{perM}}$ USD per million input tokens—as is standard for commercial LLM APIs—the marginal cost reduction per million tokens is

$$\Delta C = \text{CSR} \times C_{\text{perM}}.$$

For example, with $C_{\text{perM}} = \$3$ and CSR $\approx 0.66$, AFM would save roughly \$2 per million tokens relative to the default baseline. In high-throughput deployments, this cost difference compounds rapidly.

4.11. **Summary.** On the evaluated allergy benchmark, AFM:
– preserves safety and allergy recall in both short and medium-length conversations;
– substantially reduces token usage compared to both a stateless default setup and replay-based baselines;
– avoids the extreme prompt bloat of naive replay while retaining key user-specific constraints more reliably than a simple recency-only compressor.

These results support the central claim of AFM: selectively allocating fidelity across past messages can yield strong safety and factual continuity and accuracy at significantly lower token cost than naive history management strategies.

## 5. Discussion

5.1. **When Context Is Not Comprehension.** One of the core findings in our experiments is that merely including prior messages in an LLM's prompt does not guarantee effective memory. The naive replay strategy—common in many production chat systems—preserves all text but still fails to consistently apply user-specific facts when they matter. In our allergy benchmark, stateless and recency-compressed configurations ignored life-critical details (e.g., a peanut allergy) just one or two turns after they appeared.

This reinforces a central problem in long-context inference: information salience degrades with token distance. LLMs, even with extended context windows, exhibit a recency bias and a tendency to overweight immediately proximal text. This creates an illusion of memory where none exists and highlights the need for systems that explicitly curate and scaffold model attention over time.

5.2. **AFM's Role in the Context Stack.** AFM intervenes at the level of prompt construction—neither modifying the model weights nor introducing complex retrieval infrastructure. This makes it suitable for developers seeking efficient memory management without sacrificing simplicity, interpretability, or control.

Whereas RAG pipelines introduce database latency and risk discourse in coherence due to chunk fragmentation, AFM operates inline, restructuring the prompt without external calls. Compared to static summarization buffers (e.g., fixed-length history windows or periodic compress-and-truncate), AFM dynamically rescores every message per query, allowing real-time adaptation to user intent.

In this sense, AFM functions as a soft attention mechanism over dialogue history—prioritizing based on both semantic relevance and temporal context, and deciding how "heavy" each memory needs to be: full, compressed, or placeholder. AFM clamps cosine similarity at zero, discarding negative similarity information.

5.3. **Efficiency Without Blind Compression.** One of AFM's key contributions is demonstrating that memory compression need not be all-or-nothing. Instead of summarizing the entire past into a single global buffer (which risks flattening important details), AFM preserves structure. Every message is individually scored and passed through a fidelity filter, enabling nuanced context packing.

Empirically, this approach achieves strong token efficiency—reducing prompt size by over 60% on average—without compromising on safety-critical memory. AFM outperforms both naive replay and recency-only baselines on safety recall and token efficiency, and achieves competitive latency in our evaluated scenario.

5.4. **Failure Modes and Challenges.** AFM, while promising, is not immune to edge cases. Several challenges merit deeper scrutiny:

- **Semantic drift from LLM compression.** Abstractive summarization is inherently lossy. While effective for reducing tokens, poor LLM summaries may misrepresent prior user intent or omit subtle qualifiers—especially for long or multi-part messages. Hallucinated details, although rare in our benchmark, remain a risk under low-temperature decoding.
- **Overrepresentation of low-information turns.** Messages with high embedding similarity but low informational value (e.g., "sure," "got it" ) may receive elevated scores unless filtered via length or entropy thresholds. Future versions could incorporate dialogue-act classification to suppress such filler turns.
- **Interaction effects across messages.** AFM scores each message independently, ignoring interactions like pronoun chains or deferred clarifications. This may result in dropped context for referential utterances unless earlier turns are retained. A graph-based salience model might help capture such dependencies.
- **Static hyperparameters.** Current AFM settings use fixed half-life and importance multipliers. Adaptive weighting (e.g., task-conditioned recency schedules) could better balance short-term versus long-term salience.

5.5. **Toward Learned or Reinforced Memory Policies.** AFM's current decision logic is greedy and rule-based. However, as context management becomes more central to multi-agent systems, we envision extensions that learn optimal memory policies:

- **Reinforcement-learned packing.** Agents could be rewarded for efficient context usage that preserves factual recall or reduces hallucination rate.
- **Memory trees.** Instead of a flat log, context could be stored in a hierarchical structure with collapsible nodes and rewritable summaries.
- **Personalization-aware retention.** Messages may be scored differently based on user identity or long-term intent profiles. This would require coupling AFM with identity tracking or goal-prediction modules.

5.6. **Broader Implications.** As LLMs are increasingly embedded in tool-using agents, IDE copilots, or multi-turn reasoning chains, the need for cost-aware, interpretable, and controllable memory systems is growing. AFM's design makes it especially attractive in the following real-world scenarios:

- **Conversational agents with safety constraints.** Applications like healthcare chatbots or legal assistants must retain critical facts across turns without compromising on latency or privacy.
- **Agentic workflows with recurrent state.** Planning agents or goal-oriented bots benefit from a memory that preserves task-relevant state across tool invocations.
- **Model distillation and compression pipelines.** AFM can serve as a teacher signal in training smaller models to simulate long-context behavior without full replay.

In all these contexts, the need is the same: not just to remember more, but to remember better. AFM's selective fidelity framework takes a concrete step toward that goal.

## 6. Conclusion

Large language models are only as effective as the context they are given, yet the standard practice of naïvely replaying the full dialogue history is increasingly untenable in terms of cost,

latency, and memory reliability. Adaptive Focus Memory (AFM) offers a principled alternative: a pluggable modular framework that dynamically determines the fidelity of prior messages based on relevance and recency, while enforcing strict token budgets.

Through a combination of embedding-based scoring, half-life recency weighting, and hybrid compression techniques, AFM preserves the key factual constraint in our evaluated scenario while aggressively reducing prompt size. Our experiments show that AFM outperforms baseline approaches on key metrics of factual retention, safety alignment (AFM matches naive replay on safety alignment in our evaluated scenario at a much lower token cost.), and token efficiency in a safety-critical benchmark. Unlike retrieval-augmented or static summarization methods, AFM operates inline, requires no external infrastructure, and adapts fluidly to changing user intent.

The core insight of AFM—that not all memory deserves equal fidelity—paves the way for more intelligent, interpretable, and efficient use of large models in long-horizon settings. As the field moves toward multi-agent systems, personalized assistants, and stateful copilots, we believe AFM offers a practical foundation for real-world memory management at scale.

Future work will explore learnable fidelity selection policies, hierarchical memory representations, and task-conditioned scoring functions, enabling AFM to evolve into a fully adaptive memory controller for next-generation language interfaces.

## Acknowledgments

## References

[1] C. Cruz. Adaptive Focus Memory (AFM): Dynamic Multi-Fidelity Context Packing for Token-Efficient LLM Conversations. GitHub Repository, 2025. https://github.com/cruz209/AFMforLLM.

[2] A. Vaswani, N. Shazeer, N. Parmar, and others. Attention is All You Need. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2017.

[3] T. Brown, B. Mann, N. Ryder, and others. Language Models are Few-Shot Learners. In *NeurIPS*, 2020.

[4] OpenAI. GPT-4 Technical Report. *arXiv:2303.08774*, 2023.

[5] H. Touvron, T. Lavril, and others. LLaMA: Open and Efficient Foundation Language Models. *arXiv:2302.13971*, 2023.

[6] Anthropic. The Claude 3 Model Family. Technical Report, 2024.

[7] C. Meng and others. Locally Trainable LLMs: Efficient Adaptation via Memory and Context Optimization. *arXiv:2211.XXXX*, 2022.

[8] P. Lewis, E. Perez, A. Kandpal, and others. Retrieval-Augmented Generation for Knowledge-Intensive NLP. In *NeurIPS*, 2020.

[9] H. Ji and others. A Survey on Multi-Document Summarization. *Transactions of the ACL*, 2023.

[10] Y. Wu and others. Efficient Long-Context Memorization in Language Models. In *ICML*, 2022.

[11] A. Roberts, C. Raffel, and others. How Much Knowledge Can You Pack Into the Parameters of a Language Model? *arXiv:2002.08910*, 2020.

[12] M. Geva, R. Schuster, and others. Transformer Feed-Forward Layers Are Key-Value Memories. In *EMNLP*, 2021.

[13] W. Zhou and others. Structured Memory Architectures for Long-Context Reasoning. *arXiv:2403.XXXX*, 2024.

[14] OpenAI. tiktoken: Fast Tokenizer for LLMs. 2023. https://github.com/openai/tiktoken.

[15] OpenAI. OpenAI API Documentation. 2024. https://platform.openai.com/docs.

Purdue University
*Email address*: cchris2004@gmail.com, cruz209@purdue.edu