# TokMem: Tokenized Procedural Memory for Large Language Models

**Anonymous authors**
Paper under double-blind review

## Abstract

Large language models rely heavily on prompts to specify tasks, recall knowledge and guide reasoning. However, this reliance is inefficient as prompts must be re-read at each step, scale poorly across tasks, and lack mechanisms for modular reuse. In this paper, we aim to store and recall seen procedures efficiently. We introduce TokMem, a tokenized procedural memory that stores recurring procedures as compact, trainable embeddings. Each memory token encodes both an address to a procedure and a control signal that steers generation, enabling targeted behavior with constant-size overhead. To support continual adaptation, TokMem keeps the backbone model frozen, allowing new procedures to be added without interfering with existing ones. We evaluate TokMem on 1,000 tasks for atomic recall and multi-step function-calling for compositional recall, where it consistently outperforms retrieval-augmented generation while avoiding repeated context overhead, and fine-tuning with far fewer parameters. These results establish TokMem as a scalable and modular alternative to prompt engineering and fine-tuning, offering an explicit procedural memory for LLMs.[1]

## 1 Introduction

Large language models (LLMs) have become the foundation of modern natural language processing, powering a wide range of applications in text understanding, generation, and coding (Brown et al., 2020; Grattafiori et al., 2024). Prompting is a widely adopted way to steer LLM behavior, where in-context learning enables adaptation to new tasks without parameter updates (Brown et al., 2020). Consequently, prompt and context engineering has emerged as a dominant mechanism for specifying tasks, obtaining relevant information, and guiding multi-step reasoning or tool invocation (Wei et al., 2022b; Yao et al., 2023; Sahoo et al., 2025).

Despite its success, this reliance on long prompts is inherently inefficient. Constructing and maintaining prompts are labor-intensive and difficult to scale across many tasks (Liu et al., 2023). At inference, long prompts increase computational cost because the attention mechanism scales quadratically with sequence length (Vaswani et al., 2017), and they reduce the effective context window available for inputs and outputs, often leading to truncation and loss of details (Liu et al., 2024a). These limitations make it difficult to manage expanding tasks and to execute procedures efficiently.

To address these issues, recent approaches offload prompts into retrieval-based memory. Retrieval-augmented generation (RAG) (Lewis et al., 2020) and memory systems such as MemGPT(Packer et al., 2023) fetch and reinsert documents or conversational state at inference time. While retrieving in-context learning demonstrations (Wei et al., 2022b) can provide procedural cues that guide the model's behavior, the mechanism still largely aligns with declarative memory in cognitive science: knowledge remains as explicit text that must be repeatedly interpreted. This creates two challenges: (1) retrieved content still occupies the context window, reintroducing quadratic compute and truncation pressure, and (2) frequently used procedures are repeatedly re-read as text rather than compiled into compact, reusable procedures, missing the compression opportunity suggested by minimum description length principles (Grünwald, 2007).

We propose **Tok**enized **Mem**ory (**TokMem**), an explicit form of procedural memory that encodes recurring procedures as compact, trainable tokens while keeping the backbone frozen. Here, a *proce-*

---

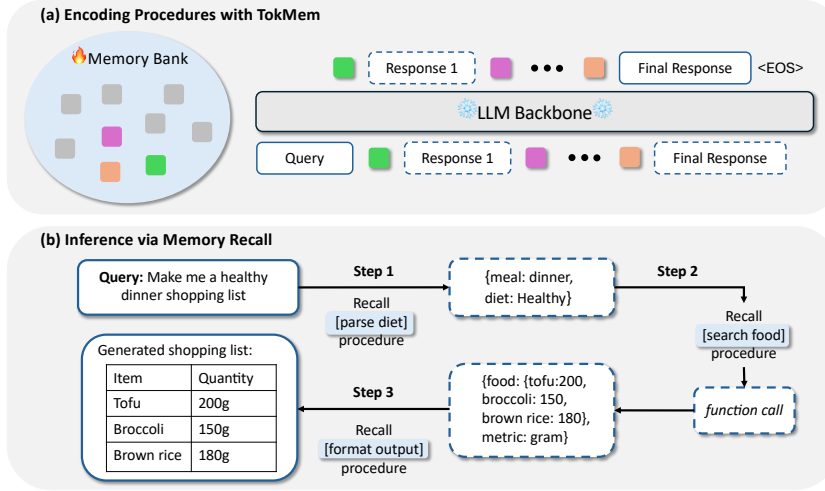[1]Our anonymous code is available here.

Figure 1: Overview of TokMem. (a) New memory (colored) tokens are interleaved with text sequences, learning with next-token-prediction while the LLM backbone remains frozen. (b) An example of inference, a query recalls and chains memory tokens (parse, search, format), enabling multi-step procedural behavior without long prompts.

*dure* means a reusable, context-response mapping that encodes a specific task behavior, inspired by physiological research Anderson & Lebiere (1998). Compared with factual knowledge, procedural knowledge (such as riding a bicycle) is often more nuanced and sophisticated.

In our TokMem, each memory token serves both as an address to a procedure and as a control signal that steer generation, enabling targeted behavior with constant-size overhead. Specifically, rather than front-loading procedures as long prompts, TokMem integrates memory tokens directly into the generation process. As shown in Figure 1b, memory tokens can be invoked and chained across stages: after producing one response segment, the model retrieves the next relevant token, which conditions the next stage, enabling the composition of multi-step behaviors such as parsing, searching, and formatting.

A key advantage of TokMem is that its memory tokens are parameter-isolated from the backbone. This design ensures that the learned procedural knowledge is fully stored in dedicated tokens, allowing new procedures to be added without interfering with existing ones. TokMem thus naturally supports continual learning, where the model can accumulate procedural skills over time while preserving stability. This capability mirrors human procedural memory, where skills are gradually acquired through practice and later invoked by contextual cues (Anderson & Lebiere, 1998). In this way, TokMem enables both efficient learning and continual expansion of procedural knowledge.

We evaluate TokMem in two complementary settings. In the atomic memory recall setting, each task from Super-Natural Instructions (Wang et al., 2022a) is treated as a distinct procedure, TokMem stores and retrieves 1,000 such procedures efficiently, without catastrophic forgetting. In the compositional memory recall setting based on function-calling tasks (Liu et al., 2024b), each tool invocation is modeled as an atomic procedure, and solving a query requires chaining multiple procedures together. TokMem supports this process by composing memory tokens, enabling the model to assemble procedures into coherent multi-step behaviors. Across both settings and multiple LLM backbones, TokMem consistently outperforms retrieval-based baselines and surpasses parametric fine-tuning while using far fewer trainable parameters.

## 2 METHOD

We begin by reviewing how Transformer-based LLMs process input sequences and then describe how TokMem departs from existing approaches to enable procedural memory.

## 2.1 TEXTUALIZED CONTEXT ENGINEERING

A Transformer (Vaswani et al., 2017) processes a sequence of tokens $(a_1, \ldots, a_n) \in \mathbb{N}^n$, where each $a_i$ is an integer representing the index of a token (usually sub-words). The model retrieves the corresponding embedding vector from the embedding layer and produces an input sequence $(\boldsymbol{x}_1, \ldots, \boldsymbol{x}_n) \in \mathbb{R}^{n \times d}$, which is then consumed to predict the next token in sequence.

Recent advances in prompting can be viewed as *textualized context engineering*, where the goal is to carefully choose input tokens that steer the model toward improved behavior. For example, chain-of-thought prompting (Wei et al., 2022a) augments input with intermediate reasoning steps to strengthen logical inference. Retrieval-based methods such as RAG (Lewis et al., 2020) and memory-augmented approaches like MemGPT (Packer et al., 2023) provide relevant information in text form by retrieving external memory. While effective, these approaches are costly: they exhaust limited model context window, and significantly increase compute due to the quadratic complexity of self-attention.

## 2.2 TOKMEM: PROCEDURAL MEMORY AS A TOKEN

Our key idea is that frequently reused procedures can be effectively "compressed" and stored by encoding them into an internalized memory token, bypassing repeated textual specification. Consider $l$ memory tokens. They are added to the vocabulary as special tokens, each represented by an embedding. Thus, they form a memory bank of $l$ special embeddings:

$$M = \begin{bmatrix} \boldsymbol{m}_1^\top \\ \vdots \\ \boldsymbol{m}_l^\top \end{bmatrix} \in \mathbb{R}^{l \times d}, \qquad \boldsymbol{m}_i \in \mathbb{R}^d. \tag{1}$$

Each $\boldsymbol{m}_i$ is a trainable vector with no direct textual translation and represents a unique procedure. For simplicity, we label each memory token $\boldsymbol{m}_i$ to have a special index $a_{m_i} \in \mathbb{N}$.

To connect these tokens with training, we first describe a single training instance. We define a *procedure–response* pair, where a procedure is invoked by a memory token $\boldsymbol{m}_i$, and the corresponding response is a sequence of textual tokens $(r_{i1}, \ldots, r_{in}) \in \mathbb{N}^n$ that encodes the information implied by the procedure. For example, the response could contain arguments for tool calling or certain pattern of text following the procedure.

In training, a procedure-response pair is represented by the concatenation of a procedural memory token $a_{m_i}$ and the embeddings of $(r_{i1}, \ldots, r_{in})$ together. Each training instance may contain multiple turns of procedure-response pairs, modeling tasks that require multi-step reasoning or composition for a query $q$. Formally, the sequentialized training sequence has the layout

$$\boldsymbol{a} = \big(q_1, \ldots, q_k, \ \underbrace{a_{m_i}, a_{r_{i1}}, a_{r_{i2}}, \ldots,}_{\text{procedure-response pair}} \ \underbrace{a_{m_j}, a_{r_{j1}}, a_{r_{j2}}, \ldots,}_{\text{procedure-response pair}} \ \ldots \big). \tag{2}$$

We adopt the standard next-token prediction loss:

$$\mathcal{L}(\boldsymbol{a}; M) = -\sum_{i > k} \log \Pr(a_i \mid \boldsymbol{a}_{<i}; M). \tag{3}$$

During optimization, the memory embeddings $(\boldsymbol{m}_1, \ldots, \boldsymbol{m}_l)$ are trainable and shared in the input embedding layer and LM head, whereas the the pre-trained text token embeddings in the input embedding layer and LM head, as well as the backbone remains frozen. For the training corpus, each memory embedding is exposed to varied queries and responses, allowing it to learn the underlying procedure representation. We visualize our training process in Figure 1a.

## 2.3 INFERENCE WITH MEMORY TOKENS

At inference, TokMem recalls procedures through memory routing and conditional response generation, where *routing* refers to the selection of the appropriate memory token for a given query. Specifically, for a query $q = (q_1, \ldots, q_k)$, the model predicts a distribution over memory tokens from its final hidden state $h_k$:

$$P(a_{m_i}|q) \propto \exp(\text{logit}(m_i|h_k)), \tag{4}$$

3

We choose the memory token with the highest probability and append it to the sequence as $(q_1, \ldots, q_k, a_{m_i})$. Then, the model generates the response autoregressively. For queries that require calling multiple procedures, after generating one response segment, the model may predict another memory token and generate the next response, as seen in Figure 1b. Notably, when a query does not correspond to any learned procedure, the logits of all memory tokens may remain low, and the model naturally defaults to generating regular text tokens.

In summary, the model decides whether to recall memory tokens (and how many) based on its training, as shown in Eqn. (2).

### 2.4 STABILIZING NEW MEMORIES

TokMem allows new procedures to be added incrementally to the memory bank, mimicking how humans continually form new procedural memories without disrupting existing skills (Anderson & Lebiere, 1998; Squire, 2009). This design enables practical deployment scenarios, where an LLM can steadily accumulate routines across domains and tasks rather than retraining from scratch.

However, adding new tokens poses stability challenges. If all procedural memories are introduced at once, the model risks overfitting to spurious patterns. Conversely, when new embeddings are added gradually, they often develop inflated norms that dominate routing logits and suppress older memories. To address this, we introduce renormalization, which is a lightweight post-update calibration to the memory bank $M \in \mathbb{R}^{l \times d}$.

Let $A$ and $I$ denote the indices of the active (new) and inactive (existing) procedural memories, respectively. We estimate the prevailing scale from the inactive set:

$$\bar{n}_I = \text{mean}_{j \in I} \lVert \boldsymbol{m}_j \rVert_2, \tag{5}$$

and rescale each active active embedding as

$$\boldsymbol{m}_i \leftarrow \boldsymbol{m}_i \cdot \frac{\bar{n}_I}{\lVert \boldsymbol{m}_i \rVert_2 + \varepsilon}, \qquad i \in A. \tag{6}$$

This operation preserves the directions of newly added embeddings while aligning their magnitudes to the established scale of the memory bank, ensuring smooth integration without overwhelming the routing dynamics. The computational overhead is negligible, scaling as $O(|A|d)$.

## 3 EXPERIMENTS

We evaluate TokMem in two complementary scenarios. In the *atomic memory recall* setting, each task from the Super-Natural Instructions dataset (Wang et al., 2022a) is framed as a standalone procedure, where a query directly maps to the desired response. In the *compositional memory recall* setting that involves calling multiple procedures, evaluated on the function-calling dataset (Liu et al., 2024b), where invoking a tool is treated as a procedure and solving a query requires composing several function calls. Experiments are conducted on Qwen (Qwen & et al., 2025) and Llama (Grattafiori et al., 2024) model families, ranging from the 0.5B-parameter Qwen to the 8B-parameter Llama model.

### 3.1 EXPERIMENTAL SETUP

**Baselines.** Across both settings, we compare TokMem with textualized context engineering, retrieval-augmented memory, and parameter-efficient fine-tuning.

- **Base**: In the atomic setting, the model answers queries without demonstrations, providing a non-parametric lower bound highlighting the need to recall task knowledge.

- **ICL**: In the compositional setting, we augment input with all tool descriptions and prepend two compositional procedure–response demonstrations, representing a context engineering baseline.

- **RAG**: We retrieve relevant demonstrations or tool usages with Sentence-BERT (Reimers & Gurevych, 2019) and prepend them to the query, following memory-augmented generation (Packer et al., 2023; Chhikara et al., 2025; Xu et al., 2025).

- **Fine-tuning**: We use low-rank adapters (Hu et al., 2022) inserted into the query and value projections of the transformer, updating millions of parameters depending on model sizes.[2] This serves as a parametric form of procedural memory, but is prone to forgetting as new tasks are introduced.

- **Replay Memory**: To mitigate catastrophic forgetting in fine-tuning, we follow the idea of experience replay (Mnih et al., 2015) by maintaining a buffer of previously seen tasks or tools and mixing them with the current training data.

**Training Details.** All methods are implemented in HuggingFace Transformers and trained on a single NVIDIA A6000 GPU with 48GB memory using mixed-precision (bfloat16) training. The backbone models remain frozen; for fine-tuning, only the adapter weights (rank $r = 8$) are updated, while for TokMem, only the embeddings of the newly added procedure IDs are trainable. Tokenizer vocabulary is expanded with these procedure IDs, and their embeddings are initialized by averaging the pretrained embeddings (Hewitt, 2021). For Replay Memory, we mix $20\%$ of replayed samples with the current batch, using a buffer of 500 examples refreshed every 10 tasks in the atomic setting and 1,000 examples updated each round in the compositional setting.

We optimize with AdamW using a learning rate $5 \times 10^{-5}$ for fine-tuning and $5 \times 10^{-3}$ for TokMem; weight decay is $10^{-2}$ for fine-tuning and zero for TokMem. Training runs for one epoch with batch size 4 and maximum sequence length 1024, using teacher forcing and applying the loss only to memory-token and response positions.

**Evaluation methods.** We evaluate TokMem from two perspectives: (1) Memory token routing accuracy, which measures whether the correct memory tokens are selected, and (2) Tasks performance, which measures the generation performance for the task (such as Rouge-L and F1), detailed in the following experiments.

## 3.2 ATOMIC MEMORY RECALL

**Dataset Details.** We evaluate on the Super-Natural Instructions (SNI) dataset (Wang et al., 2022a), which provides diverse QA-style natural language tasks. Here, each task is treated as an individual procedure: a query directly invokes the learned procedure to produce the desired response. We sample 1,000 English tasks, each task contains 500 training and 50 test examples. To reflect how memories are typically acquired over time, we introduce tasks sequentially during training rather than all at once, but training samples of each task are shuffled. We scale the number of tasks from 10 up to 1,000, and record checkpoints after training on $\{10, 50, 200, 500, 1,000\}$ tasks. This resembles incremental domain adaptation (Asghar et al., 2020), where at each checkpoint, the performance is evaluated across all previously seen tasks. Additional tasks details are provided in Appendix D.1.

We follow Wang et al. (2022a) and use Rouge-L (Lin, 2004) to evaluate generation quality. For the methods with memory routing (RAG and TokMem), we additionally report accuracy, reflecting whether the correct procedure was selected and applied.

**Decoupled Memory Embeddings.** We also include an ablation variant where memory tokens are decoupled to an address token and a steering token. This decoupling separates the roles of a memory token and also increase the capacity of TokMem. We refer to this variant as TokMem with decoupled embeddings (TokMem+DC), with further details provided in Appendix A.

**Results and Findings.** TokMem provides the most consistent and scalable performance across models and task scales. As shown in Table 1, non-parametric methods such as Base shows stable but fail to achieve competitive performance. RAG performs reasonably well on when memory is not heavy but quickly degrade as the number of task memory increases, indicating its sensitivity to retriever quality. Parametric methods such as fine-tuning achieve stronger initial accuracy but suffer from forgetting as tasks accumulate; replay memory alleviates this issue but still falls short of TokMem. By contrast, we see that TokMem maintains high accuracy with minimal performance drop when acquiring new task memories, achieving the best average results across all settings.

---

[2]In our preliminary study, we found that adding training parameters by applying low-rank adapters to more projections did not improve performance while greatly increasing computational cost. We therefore follow Hu et al. (2022) and apply LoRA to the query and value projections only.

Table 1: Atomic recall performance on SNI, reported with Rouge-L. TokMem consistently outperforms fine-tuning and RAG across models and scales, maintaining strong results even at 1,000 tasks.

| Model | Method | Number of Tasks | | | | | |
| | | 10 | 50 | 200 | 500 | 1000 | Avg. |
|---|---|---|---|---|---|---|---|
| Qwen 2.5 0.5B | Base | 33.9 | 39.0 | 38.8 | 39.1 | 38.5 | 37.9 |
| | RAG | 50.4 | 43.2 | 38.8 | 36.2 | 34.7 | 40.7 |
| | Fine-Tuing | 52.4 | 48.0 | 40.6 | 41.7 | 43.2 | 45.2 |
| | Replay Memory | 52.4 | 49.5 | 47.2 | 47.7 | 46.7 | 48.7 |
| | TokMem | 52.8 | 51.3 | 49.3 | 50.2 | 50.0 | 50.7 |
| | TokMem+DC | 53.8 | 50.5 | 50.2 | 50.9 | 50.0 | **51.1** |
| Llama 3.2 3B | Base | 16.6 | 19.9 | 20.0 | 18.7 | 18.2 | 18.7 |
| | RAG | 60.0 | 48.7 | 45.8 | 42.3 | 39.9 | 47.3 |
| | Fine-Tuing | 67.1 | 59.1 | 59.5 | 58.4 | 57.9 | 60.4 |
| | Replay Memory | 67.1 | 61.1 | 60.6 | 61.4 | 60.0 | 62.0 |
| | TokMem | 68.0 | 62.3 | 61.2 | 61.5 | 61.5 | **62.9** |
| | TokMem+DC | 68.8 | 62.5 | 58.7 | 61.7 | 61.1 | 62.6 |
| Llama 3.1 8B | Base | 27.2 | 27.8 | 30.4 | 29.6 | 29.5 | 28.9 |
| | RAG | 63.8 | 53.9 | 49.1 | 45.3 | 42.6 | 50.9 |
| | Fine-Tuing | 75.8 | 64.3 | 63.2 | 58.7 | 61.6 | 64.7 |
| | Replay Memory | 75.8 | 65.2 | 64.5 | 63.4 | 63.6 | 66.5 |
| | TokMem | 75.4 | 65.5 | 65.1 | 64.4 | 64.8 | **67.0** |
| | TokMem+DC | 75.6 | 65.8 | 63.7 | 64.2 | 64.4 | 66.7 |

Table 2: Task routing accuracy. TokMem achieves near-perfect routing accuracy at 1,000 tasks, far exceeding RAG retriever whose accuracy falls below 80%.

| Model | Method | Number of Tasks | | | | |
| | | 10 | 50 | 200 | 500 | 1000 |
|---|---|---|---|---|---|---|
| Sentence-Bert | RAG | 99.6 | 92.6 | 88.7 | 83.2 | 79.7 |
| Qwen 2.5 0.5B | TokMem | 99.4 | 98.6 | 97.4 | 96.9 | 94.7 |
| | TokMem+DC | **99.4** | **99.2** | **98.4** | **97.2** | **96.1** |
| Llama 3.2 3B | TokMem | **100.0** | **99.9** | **98.3** | **97.1** | **96.1** |
| | TokMem+DC | 99.8 | 99.3 | 97.2 | 96.2 | 95.4 |
| Llama 3.1 8B | TokMem | **99.8** | **99.6** | **98.9** | **97.7** | **97.5** |
| | TokMem+DC | 99.7 | 99.4 | 97.8 | 97.2 | 97.2 |

The decoupled variant (TokMem+DC) yields modest gains for the smaller Qwen 0.5B model but no improvement for larger Llama models, and in some cases it underperforms TokMem when scaling to many tasks. Overall, although TokMem+DC is a tempting variant, it does not provide additional benefits. We advocate for the simple yet effective TokMem (without DC).

Table 2 further highlights TokMem's robustness in memory routing. Its accuracy remains above 94% even at 1,000 tasks with the smallest 0.5B model, significantly outperforming the Sentence-BERT retriever used in RAG, whose accuracy drops below 80% when have the stress to route for 1,000 tasks. This high-fidelity memory routing enables TokMem to sustain strong performance without relying on external retrieval mechanisms or fine-tuning, demonstrating its advantage in continual and large-scale task acquisition.

**Analysis of Training Efficiency.** We compare the training sample efficiency of LoRA fine-tuning and TokMem on the first 10 tasks from SNI, a mixture setup that removes the impact of forgetting. We set the adapter rank to $r = 1$, which helps prevent overfitting and aligns its parameter scale with TokMem. The result in Figure 2 shows that TokMem consistently achieves higher performance than fine-tuning across all sample budgets, with the greatest advantage appearing in the low-data regime. The decoupled variant (TokMem+DC) offers small but consistent improvements, particularly when more samples are available. Overall, these results highlight TokMem's ability to learn new procedures effectively with limited data, making it a both parameter and data efficient approach to memory acquisition.
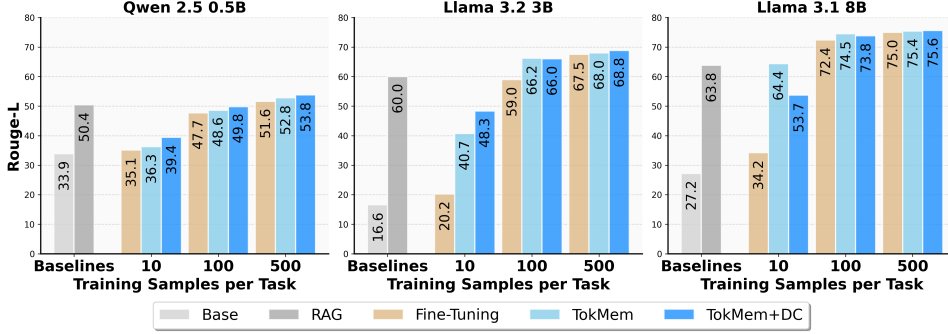
Figure 2: Sample efficiency on a 10-task mixture from SNI. TokMem consistently outperforms fine-tuning in the low-data regime. TokMem can surpass RAG with only 10 training samples, demonstrating strong few-shot learning capability.

Table 3: Compositional tool-use performance on APIGen. TokMem achieves strong tool selection and argument F1 across multiple calls, outperforming ICL and RAG with lower input-augmentation complexity, and surpassing fine-tuning with far fewer trainable parameters.

| Model | Method | #Params | Tool Selection | | | | Argument | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | 2 calls | 3 calls | 4 calls | Avg. | 2 calls | 3 calls | 4 calls | Avg. |
| Llama 3.2 1B | ICL | – | 27.6 | 11.1 | 10.5 | 16.4 | 0.6 | 0.7 | 0.0 | 0.4 |
| | RAG | – | 29.5 | 10.8 | 10.5 | 16.9 | 7.2 | 1.0 | 0.0 | 2.7 |
| | Fine-Tuing | 0.85M | 10.4 | 9.5 | 7.0 | 9.0 | 77.3 | 72.6 | 55.8 | 68.6 |
| | TokMem (w/o adapt) | 0.10M | 86.8 | 80.9 | 90.8 | 86.2 | 68.9 | 61.1 | 73.0 | 67.7 |
| | TokMem (w/ adapt) | 0.10M | 98.4 | 98.0 | 98.9 | **98.4** | 84.3 | 84.3 | 87.8 | **85.5** |
| Llama 3.2 3B | ICL | – | 66.8 | 59.2 | 59.6 | 61.9 | 42.2 | 42.3 | 38.8 | 44.1 |
| | RAG | – | 78.1 | 71.2 | 69.3 | 72.8 | 54.8 | 53.1 | 62.7 | 56.9 |
| | Fine-Tuing | 2.29M | 98.7 | 98.1 | 96.8 | 97.9 | 87.9 | 86.6 | 82.9 | 85.8 |
| | TokMem (w/o adapt) | 0.15M | 82.6 | 79.3 | 67.2 | 76.4 | 65.4 | 57.2 | 50.2 | 57.6 |
| | TokMem (w/ adapt) | 0.15M | 99.2 | 98.2 | 100.0 | **99.2** | 85.9 | 86.7 | 88.3 | **86.3** |
| Llama 3.1 8B | ICL | – | 79.7 | 72.9 | 75.4 | 76.0 | 51.5 | 52.6 | 57.3 | 53.8 |
| | RAG | – | 79.6 | 75.3 | 93.0 | 82.6 | 53.3 | 57.1 | 69.2 | 59.9 |
| | Fine-Tuing | 3.41M | 98.8 | 97.2 | 98.2 | 98.1 | 87.7 | 86.8 | 88.2 | 87.6 |
| | TokMem (w/o adapt) | 0.20M | 84.9 | 82.0 | 81.6 | 82.8 | 65.8 | 56.7 | 65.9 | 62.8 |
| | TokMem (w/ adapt) | 0.20M | 99.4 | 97.9 | 100.0 | **99.1** | 88.1 | 86.5 | 93.4 | **89.3** |

## 3.3 COMPOSITIONAL MEMORY RECALL

**Dataset Details.** We construct a benchmark from the APIGen dataset (Liu et al., 2024b) by sampling the 50 frequently used tools. Here, each tool invocation is treated as an atomic procedure, and solving a query requires composing multiple such procedures. We synthesize 5,000 training queries and 500 test queries, both capped at four calls. Details of the dataset can be found in Appendix D.2

We report performance using two F1 metrics: (i) Tool Prediction F1, which measures whether the correct tools are invoked; and (ii) Argument Generation F1, which evaluates the correctness of function call arguments. For robustness to semantic equivalence, both gold and predicted outputs are normalized into Abstract Syntax Trees before scoring (Patil et al., 2025).

**Adaptation for Compositionality.** We found that TokMem benefits from a brief adaptation phase that exposes the backbone with the compositional structures of memory tokens. Concretely, we fine-tune the backbone on a held-out auxiliary tool set using the same LoRA fine-tuning setup as the baseline. The adapted weights are then merged, after which the backbone remains frozen for memory acquisition and evaluation (see Appendix B.1 for details). Note that the auxiliary tool set is different from the evaluation set, so it does not break our frozen-backbone setup.

The proposed adaptation initialization teaches the model on how to interleave responses with memory tokens. We thus consider it a part of the TokMem approach for compositional tasks and use TokMem with adaptation as the default configuration. We also analyze the effect of this adaptation
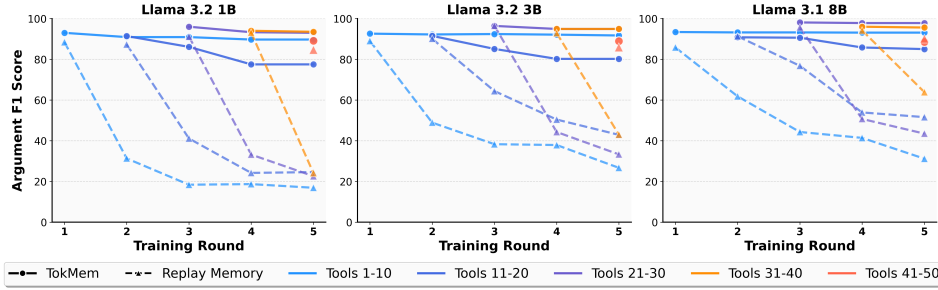
Figure 3: Forgetting analysis in continual adaptation. As new tools are introduced, fine-tuning with replay memory suffers sharp drops on earlier tasks, while TokMem maintains stable performance. Larger models show stronger retention due to greater capacity.
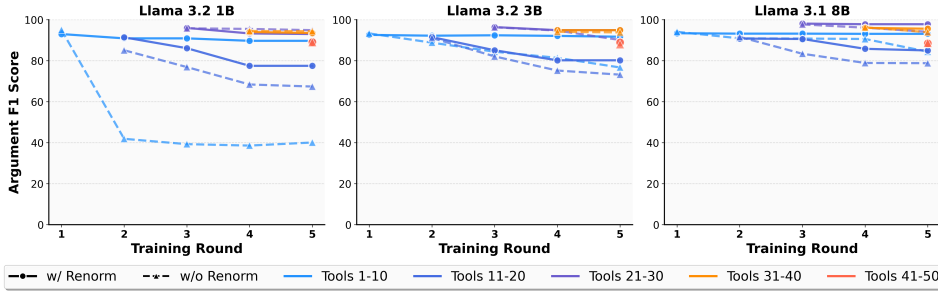


Figure 4: Effect of renormalization on TokMem. Without renormalization, new tokens dominate and older ones are forgotten, particularly in smaller models with limited embedding capacity. Renormalization effectively mitigates this by balancing norms across tokens.

phase on the fine-tuning baseline in Appendix B.2. We find that it degrades baseline performance and we excludes it for a stronger baseline.

**Results and Findings.** Table 3 shows that TokMem achieves the strongest overall performance. Without adaptation for compositionality, it outperforms RAG while avoiding the added complexity from an external retrieval mechanism. Non-parametric baselines such as ICL and RAG perform poorly on both tool prediction and argument generation, particularly with the smaller Llama 1B model, likely due to its weak instruction-following ability.

Compared with parametric baselines, TokMem consistently matches or surpasses LoRA fine-tuning while requiring an order of magnitude fewer trainable parameters. For example, on the Llama 8B model, LoRA requires 3.41M training parameters while TokMem only needs 0.2M to achieves higher performance.

Notably, TokMem exhibits stronger interpretability between tool selection and argument generation, with improvements in the former translating directly into the latter. By contrast, LoRA fine-tuning shows weaker alignment. For example, as seen with the 1B model, it often generates plausible arguments even when tool selection is incorrect, indicating that its argument generation is not properly grounded in the chosen tools.

TokMem also shows significantly better compositional ability than fine-tuning, successfully executing queries with more function calls than it has seen during training. Detailed results are provided in Appendix B.4.

**Analysis of Forgetting.** We compare TokMem with replay memory in a continual adaptation setting, where tools are introduced sequentially over five training rounds (e.g., tools 1–10 in the first round, 11–20 in the second, and so forth). As shown in Figure 3, we see that replay memory struggles to prevent catastrophic forgetting as new tools are introduced. By contrast, TokMem maintains higher performance across tool groups, with only mild declines that primarily reflect the growing

8

Table 4: Comparison of TokMem vs. prefix tuning on memorizing text from the *Fanfics* dataset. TokMem converges faster and achieves lower perplexity than prefix tuning, particularly with few memory tokens.

| Method | 1024 tokens | | 2048 tokens | | 4096 tokens | |
|---|---|---|---|---|---|---|
| | Steps@90%Best ↓ | PPL ↓ | Steps@90%Best ↓ | PPL ↓ | Steps@90%Best ↓ | PPL ↓ |
| *Prefix tuning-1* | 1700 | 3.81 | 2300 | 8.77 | 2200 | 14.32 |
| *TokMem-1* | **1200** | **3.28** | **1400** | **7.07** | **1700** | **12.27** |
| *Prefix tuning-2* | **500** | 1.13 | 1700 | 3.51 | 1800 | 8.38 |
| *TokMem-2* | 600 | **1.09** | **1300** | **2.75** | **1700** | **7.21** |
| *Prefix tuning-5* | 300 | 1.07 | 500 | 1.17 | 1400 | 2.39 |
| *TokMem-5* | **200** | **1.03** | **500** | **1.15** | **1400** | **1.91** |

number of tools. Larger models exhibit better retention for both approaches, likely due to their expanded parameter capacity, which reduces the risk of interference with previously learned tools.

We further investigate the effect of the renormalization step introduced in Section 2.4 on newly added memory tokens, whose norms may otherwise dominate older tokens in the softmax (analyzed in Appendix B.5). As seen in Figure 4, TokMem without renormalization shows noticeable forgetting especially when the size of the model is small. However, larger models are more robust to forgetting even without renormalization, again due to their greater embedding capacity. Overall, renormalization improves TokMem's resistance to forgetting by balancing routing between both new and old memory tokens. Additional analysis on the benefits of keeping the backbone frozen for continual memory acquisition is provided in Appendix B.3.

## 3.4 ANALYSIS ON MEMORY PLACEMENT

An important design choice in TokMem is the placement of memory tokens within the input sequence, which directly influences how the backbone model attends to and integrates procedural knowledge. While TokMem introduces a memory routing mechanism for generating tokens, its effectiveness also depends on this placement strategy. In the absence of routing, TokMem reduces to a prompt-tuning method (Li & Liang, 2021; Lester et al., 2021) with learnable embeddings, but it distinctively adopts an *infix* placement: `query ⊕ MEM ⊕ response`. Our experiments indicate that this infix design allows memory tokens to be activated after the context has been encoded, enabling context-aware conditioning and natural composition of multiple procedures.

However, it is unclear whether this memory placement is strictly better than the more common *prefix* formulation: `MEM ⊕ query ⊕ response` used in prior prompt-tuning work, where prefix tokens influence generation without having observed the query. To study the impact of placement, we compare prefix and infix placements under matched token budgets in the single-task setting.

**Setup.** We compare TokMem with infix memory placement against prefix tuning by stress-testing the capacity of memory tokens (Sastre & Rosá, 2025) using the recent *Fanfics* dataset collected after the pretraining of LLMs (Kuratov et al., 2025). We fix the sequence length to 128 tokens and vary the batch size from 8 to 32, compressing batches of 1024 to 4096 response tokens into 1 to 5 memory tokens. For each sequence, we prepend a randomly generated query that serves only as a marker to distinguish the two placements, while the actual target to be learned remains the response.

We measure learning speed using Steps@90%Best, defined as the number of training steps (evaluated every 100 steps) required to reach 90% of the best perplexity. Results are averaged over five runs. Additional experiments evaluating generalization on a math reasoning dataset are provided in Appendix C.

**Results.** Table 4 shows that TokMem consistently achieves lower perplexity and often converges faster than prefix tuning. With a single token, TokMem reaches 90% of the best perplexity roughly 30% sooner than prefix tuning, indicating that conditioning memory after the query helps the model learn more efficiently. Interestingly, when more tokens are available (e.g., five tokens), the performance gap narrows. This suggests that prior work (Li & Liang, 2021), which typically uses dozens

or even hundreds of tokens, may have underestimated the importance of memory placement in low-token regimes, where each token must compress more procedural information.

## 4 RELATED WORK

Equipping LLMs with memory has been explored through multiple directions. Most existing approaches emphasize declarative memory, where the objective is to store and retrieve explicit information such as facts or conversation history (Packer et al., 2023; Chhikara et al., 2025; Zhong et al., 2024). In contrast, parameter-based approaches internalize task-specific behaviors within model parameters, resembling procedural memory. TokMem builds on this latter view while emphasizing modularity and compositionality.

**Text-based External Memory.** A common approach is to externalize memory as textual content retrieved at inference time. Retrieval-augmented generation (RAG)(Lewis et al., 2020) and its variants(Guu et al., 2020; Karpukhin et al., 2020; Borgeaud et al., 2022; Khandelwal et al., 2020) attach relevant textual chunks during inference, while RET-LLM (Modarressi et al., 2023) encodes knowledge as symbolic triplets. Building on these ideas, more recent systems such as MemGPT (Packer et al., 2023), Mem0 (Chhikara et al., 2025), and A-Mem (Xu et al., 2025) extend these ideas to conversational settings through hierarchical or summarization-based memory states. While effective for factual recall, these approaches are not optimized for procedural control and often incur significant inference-time overhead due to the re-read of textual memory.

**Parameter-based Memory.** Another line of work encodes memory directly into model parameters. Fine-tuning and multitask instruction tuning (Wei et al., 2022a; Sanh et al., 2021), as well as parameter-efficient variants such as LoRA (Hu et al., 2022) allow models to acquire new procedures, but task knowledge is entangled. Mixture-of-LoRAs (Feng et al., 2024) address the entanglement issue but the mixtures are typically invoked independently and are not designed for memory composition. MemoryLLM (Zhong et al., 2024) introduces latent memory pools but remian entangled. Prompt-based methods such as prompt tuning (Lester et al., 2021; Wu et al., 2025) store knowledge implicitly as global embeddings without selective routing, and L2P (Wang et al., 2022b) introduces modular prompt pools but still relies on an external controller to determine which prompts are retrieved. Prompt compression methods (Mu et al., 2023; Chevalier et al., 2023) compress prompts into one-size-fits-all representations, which may distort prompt information. ToolGen (Wang et al., 2025) compresses tools into virtual tokens but focuses on post-training the backbone through multi-stage fine-tuning. By contrast, TokMem keeps the backbone frozen, and introduces discrete memory units that can be added or composed without retraining, supporting continual adaptation.

**Compositional Memory.** A complementary direction explores how models compose skills from simpler building blocks. Chain-of-thought prompting (Wei et al., 2022b) and tool-augmented reasoning frameworks such as Toolformer (Schick et al., 2023) enable multi-step reasoning, but rely on textual instructions that must be re-interpreted at each step. Modular parameter methods (Rosenbaum et al., 2018; Pfeiffer et al., 2021) create specialized adapters that can be recombined, but composition requires parameter merging or heuristic routing. TokMem differs by representing procedures as discrete tokens that can be chained directly in context, enabling lightweight parameter-isolated composition.

## 5 CONCLUSION

We introduced Tokenized Memory (TokMem), a parameter-efficient framework that encodes procedural memory as compact tokens. TokMem enables selective recall and compositional use of procedures without modifying backbone parameters, achieving strong performance across multitask and tool-augmented reasoning benchmarks.

Future directions are discussed in Appendix E, including reinforcement learning for stronger compositional generalization, and personalization through user-specific memory banks. These extensions pave the way for scalable, compact, and user-adaptive memory systems in large language models.

# REFERENCES

John R. Anderson and Christian Lebiere. *The Atomic Components of Thought*. Lawrence Erlbaum Associates, 1998.

Nabiha Asghar, Lili Mou, Kira A. Selby, Kevin D. Pantasdo, Pascal Poupart, and Xin Jiang. Progressive memory banks for incremental domain adaptation. In *International Conference on Learning Representations*, 2020. URL https://openreview.net/forum?id=BkepbpNFwr.

Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Millican, George Van Den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, et al. Improving language models by retrieving from trillions of tokens. *arXiv preprint arXiv:2112.04426*, 2022. URL https://arxiv.org/abs/2112.04426.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in Neural Information Processing Systems*, 33:1877–1901, 2020. URL https://papers.nips.cc/paper/2020/hash/1457c0d6bfcb4967418bfb8ac142f64a-Abstract.html.

Alexis Chevalier, Alexander Wettig, Anirudh Ajith, and Danqi Chen. Adapting language models to compress contexts. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 3829–3846. Association for Computational Linguistics, 2023. URL https://aclanthology.org/2023.emnlp-main.232/.

Prateek Chhikara, Dev Khant, Saket Aryan, Taranjeet Singh, and Deshraj Yadav. Mem0: Building production-ready ai agents with scalable long-term memory, 2025. URL https://arxiv.org/abs/2504.19413.

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. Training verifiers to solve math word problems, 2021. URL https://arxiv.org/abs/2110.14168.

Wenfeng Feng, Chuzhan Hao, Yuewei Zhang, Yu Han, and Hao Wang. Mixture-of-LoRAs: An efficient multitask tuning for large language models. *arXiv preprint arXiv:2403.03432*, 2024. URL https://arxiv.org/abs/2403.03432.

Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, and et al. The llama 3 herd of models, 2024. URL https://arxiv.org/abs/2407.21783.

Peter D Grünwald. *The minimum description length principle*. 2007. URL https://direct.mit.edu/books/monograph/3813/The-Minimum-Description-Length-Principle.

Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Ming-Wei Chang. Retrieval augmented language model pre-training. In *International Conference on Machine Learning*, pp. 3929–3938, 2020. URL http://proceedings.mlr.press/v119/guu20a.html.

John Hewitt. Initializing new word embeddings for pretrained language models, 2021. URL https://nlp.stanford.edu/~johnhew//vocab-expansion.html.

Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2022. URL https://openreview.net/forum?id=nZeVKeeFYf9.

Vladimir Karpukhin, Barlas Oğuz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. Dense passage retrieval for open-domain question answering. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing*, pp. 6769–6781, 2020. URL https://aclanthology.org/2020.emnlp-main.550/.

Urvashi Khandelwal, Omer Levy, Dan Jurafsky, Luke Zettlemoyer, and Mike Lewis. Generalization through memorization: Nearest neighbor language models. In *International Conference on Learning Representations*, 2020. URL `https://openreview.net/forum?id=HklBjCEKvH`.

Yuri Kuratov, Mikhail Arkhipov, Aydar Bulatov, and Mikhail Burtsev. Cramming 1568 tokens into a single vector and back again: Exploring the limits of embedding space capacity. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 19323–19339. Association for Computational Linguistics, 2025. URL `https://aclanthology.org/2025.acl-long.948/`.

Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt tuning. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 3045–3059, 2021. URL `https://aclanthology.org/2021.emnlp-main.243/`.

Patrick Lewis, Ethan Perez, Aleksandar Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented generation for knowledge-intensive NLP tasks. *Advances in Neural Information Processing Systems*, 33: 9459–9474, 2020. URL `https://proceedings.neurips.cc/paper/2020/hash/6b493230205f780e1bc26945df7481e5-Abstract.html`.

Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing*, pp. 4582–4597, 2021. URL `https://aclanthology.org/2021.acl-long.353/`.

Chin-Yew Lin. ROUGE: A package for automatic evaluation of summaries. pp. 74–81. Association for Computational Linguistics, 2004. URL `https://aclanthology.org/W04-1013/`.

Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. Lost in the middle: How language models use long contexts. *Transactions of the Association for Computational Linguistics*, pp. 157–173, 2024a. URL `https://aclanthology.org/2024.tacl-1.9/`.

Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. *ACM Comput. Surv.*, 55(9), 2023. URL `https://doi.org/10.1145/3560815`.

Yujia Liu, Jiacheng Zhang, Ziming Wang, Xiaohan Li, Jing Li, and Hao Wang. APIGen: Automated API code generation for function-calling capabilities in large language models, 2024b. URL `https://arxiv.org/abs/2406.18518`.

Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *nature*, pp. 529–533, 2015. URL `https://www.nature.com/articles/nature14236`.

Ali Modarressi, Ayyoob Imani, Mohsen Fayyaz, and Hinrich Schütze. RET-LLM: Towards a general read-write memory for large language models. In *Advances in Neural Information Processing Systems*, volume 36, pp. 15558–15571, 2023. URL `https://proceedings.neurips.cc/paper_files/paper/2023/hash/6a4cd50db0cad92c4c8d9e6ee01ac8c6-Abstract-Conference.html`.

Jesse Mu, Xiang Lisa Li, and Noah Goodman. Learning to compress prompts with gist tokens. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL `https://openreview.net/forum?id=2DtxPCL3T5`.

Charles Packer, Vivian Fang, Shishir Gururaj Patil, Kevin Lin, Sarah Wooders, and Joseph E Gonzalez. MemGPT: Towards LLMs as operating systems. *arXiv preprint arXiv:2310.08560*, 2023. URL `https://arxiv.org/abs/2310.08560`.

Shishir G Patil, Huanzhi Mao, Fanjia Yan, Charlie Cheng-Jie Ji, Vishnu Suresh, Ion Stoica, and Joseph E. Gonzalez. The berkeley function calling leaderboard (BFCL): From tool use to agentic evaluation of large language models. In *Forty-second International Conference on Machine Learning*, 2025. URL `https://openreview.net/forum?id=2GmDdhBdDk`.

Jonas Pfeiffer, Aishwarya Kamath, Andreas Rücklé, Kyunghyun Cho, and Iryna Gurevych. Adapter-Fusion: Non-destructive task composition for transfer learning. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pp. 487–503. Association for Computational Linguistics, 2021. URL `https://aclanthology.org/2021.eacl-main.39/`.

Qwen and et al. Qwen2.5 technical report, 2025. URL `https://arxiv.org/abs/2412.15115`.

Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, 2019. URL `https://arxiv.org/abs/1908.10084`.

Clemens Rosenbaum, Tim Klinger, and Matthew Riemer. Routing networks: Adaptive selection of non-linear functions for multi-task learning. In *International Conference on Learning Representations*, 2018. URL `https://openreview.net/forum?id=ry8dvM-R-`.

Pranab Sahoo, Ayush Kumar Singh, Sriparna Saha, Vinija Jain, Samrat Mondal, and Aman Chadha. A systematic survey of prompt engineering in large language models: Techniques and applications, 2025. URL `https://arxiv.org/abs/2402.07927`.

Victor Sanh, Albert Webson, Colin Raffel, Stephen H Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Teven Le Scao, Arun Raja, et al. Multitask prompted training enables zero-shot task generalization. *arXiv preprint arXiv:2110.08207*, 2021. URL `https://arxiv.org/abs/2110.08207`.

Ignacio Sastre and Aiala Rosá. Memory tokens: Large language models can generate reversible sentence embeddings. In *Proceedings of the First Workshop on Large Language Model Memorization*, pp. 183–189. Association for Computational Linguistics, 2025. URL `https://aclanthology.org/2025.l2m2-1.14/`.

Timo Schick, Jane Dwivedi-Yu, Roberto Dessi, Roberta Raileanu, Maria Lomeli, Eric Hambro, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. Toolformer: Language models can teach themselves to use tools. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL `https://openreview.net/forum?id=Yacmpz84TH`.

Larry R Squire. Memory and brain systems: 1969–2009. *Journal of Neuroscience*, 29:12711–12716, 2009. URL `https://www.jneurosci.org/content/29/41/12711`.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in Neural Information Processing Systems*, 30:5998–6008, 2017. URL `https://papers.nips.cc/paper_files/paper/2017/hash/3f5ee243547dee91fbd053c1c4a845aa-Abstract.html`.

Renxi Wang, Xudong Han, Lei Ji, Shu Wang, Timothy Baldwin, and Haonan Li. Toolgen: Unified tool retrieval and calling via generation. In *The Thirteenth International Conference on Learning Representations*, 2025. URL `https://openreview.net/forum?id=XLMAMmowdY`.

Yizhong Wang, Swaroop Mishra, Pegah Alipoormolabashi, Yeganeh Kordi, Amirreza Mirzaei, Atharva Naik, Arjun Ashok, Arut Selvan Dhanasekaran, Anjana Arunkumar, David Stap, et al. Super-NaturalInstructions: Generalization via declarative instructions on 1600+ NLP tasks. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp. 5085–5109, 2022a. URL `https://aclanthology.org/2022.emnlp-main.340/`.

Zifeng Wang, Zizhao Zhang, Chen-Yu Lee, Han Zhang, Ruoxi Sun, Xiaoqi Ren, Guolong Su, Vincent Perot, Jennifer Dy, and Tomas Pfister. Learning to prompt for continual learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 139–149, 2022b. URL `https://openaccess.thecvf.com/content/CVPR2022/papers/Wang_Learning_To_Prompt_for_Continual_Learning_CVPR_2022_paper.pdf`.

Jason Wei, Maarten Bosma, Vincent Y Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M Dai, and Quoc V Le. Finetuned language models are zero-shot learners. *International Conference on Learning Representations*, 2022a. URL `https://openreview.net/forum?id=gEZrGCozdqR`.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed H. Chi, Quoc V Le, and Denny Zhou. Chain of thought prompting elicits reasoning in large language models. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), *Advances in Neural Information Processing Systems*, 2022b. URL `https://openreview.net/forum?id=_VjQlMeSB_J`.

Zijun Wu, Yongchang Hao, and Lili Mou. Ulpt: Prompt tuning with ultra-low-dimensional optimization, 2025. URL `https://arxiv.org/abs/2502.04501`.

Wujiang Xu, Zujie Liang, Kai Mei, Hang Gao, Juntao Tan, and Yongfeng Zhang. A-mem: Agentic memory for llm agents. *arXiv preprint arXiv:2502.12110*, 2025. URL `https://arxiv.org/pdf/2502.12110`.

Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. In *The Eleventh International Conference on Learning Representations*, 2023. URL `https://openreview.net/forum?id=WE_vluYUL-X`.

Yu Zhong, Longyue Wang, Jiajun Liu, Guangdong Chen, Minjun Wu, Qifan Zhou, Zerui Wang, Xianzhi Wang, et al. MemoryLLM: Towards self-updatable large language models. *arXiv preprint arXiv:2402.04624*, 2024. URL `https://arxiv.org/abs/2402.04624`.
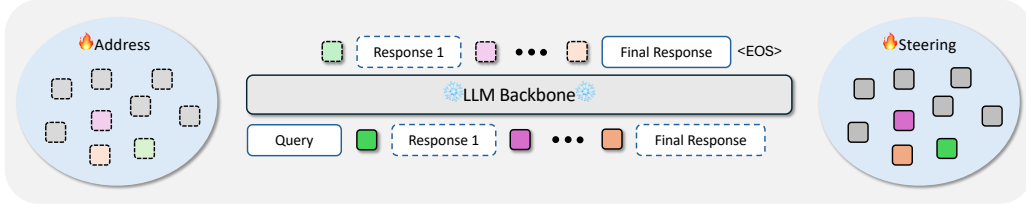
Figure 5: Overview of Decoupled TokMem embeddings, which learns separate memory matrices for address of memories and generation steering.

# A  DECOUPLED EMBEDDING FOR TOKMEM

In the standard TokMem formulation, each memory token embedding $m_i \in \mathbb{R}^d$ is shared shared across two roles: (1) addressing for memory routing and (2) steering for generation. We consider a decoupled (DC) variant that separates these functions into two embedding matrices:

$$M^{\text{addr}} = (u_1; \ldots; u_l) \in \mathbb{R}^{l \times d}, \qquad M^{\text{steer}} = (s_1; \ldots; s_l) \in \mathbb{R}^{l \times d}. \tag{7}$$

Here, $M^{\text{addr}}$ provides *address embeddings* at the output layer. When a memory token is predicted, the model produces a distribution over indices $i$ according to $M^{\text{addr}}$. The chosen index $i$ is then used to retrieve the corresponding steering embedding $s_i$ from $M^{\text{steer}}$, which is injected into the input sequence and influences subsequent generation.

Training follows the standard next-token prediction objective, analogous to Equation 3:

$$\mathcal{L}(a; M^{\text{addr}}, M^{\text{steer}}) = -\sum_{i>k} \log \Pr(a_i \mid a_{<i}; M^{\text{addr}}, M^{\text{steer}}). \tag{8}$$

where $k$ denotes the query length. During optimization, only $M^{\text{addr}}$ and $M^{\text{steer}}$ are updated; the backbone remains frozen. In addition, the renormalization treatment introduced in Section 2.2 is only applied to the address embeddings in $M^{\text{addr}}$.

This decoupled formulation provides a clean separation of functionality: routing is handled via $M^{\text{addr}}$, while steering is controlled by $M^{\text{steer}}$. While conceptually clean, our experiments do not show consistent improvements over the coupled formulation, particularly on larger models, where embedding capacity is sufficient to jointly support both roles.

# B  DETAILS FOR COMPOSITIONAL MEMORY RECALL

## B.1  DETAILS OF ADAPTATION PHASE

In compositional scenarios, the model should not only recall individual procedures but also compose them to solve multi-step queries. To prepare TokMem for such use, we construct a held-out auxiliary training set of 50 tools (5,000 samples) following Section 3.3. The backbone is then fine-tuned for one epoch on this set using LoRA, jointly with the temporary memory embeddings, before the adapted weights are merged and frozen.

The intuition for this adaptation phase to let the LLM learn to align its routing and generation behavior with compositional memory recall. After adaptation, the temporary embeddings are discarded, while the adapted backbone is retained for inference with new tasks. This procedure provides a general inductive bias for modular composition, enabling it to generalize to new tools and procedures without further retraining.

Algorithm 1 summarizes this lightweight procedure. Temporary memory embeddings are inserted into the input sequence, the loss is optimized jointly over memory and response tokens, and once the backbone has adapted, the temporary memory bank is discarded.

15

---

**Algorithm 1** Adaptation Phase for Compositional Memory Recall

---

**Require:** Pretrained backbone $f_{\theta_0}$, adaptation traces $\mathcal{D}_{\text{adapt}}$ from held-out procedures
 1: Initialize backbone $\theta \leftarrow \theta_0$ and temporary memory embeddings $\mathcal{M}$
 2: Set learning rates $\eta_\theta$ and $\eta_{\mathcal{M}}$
 3: **for** each minibatch in $\mathcal{D}_{\text{adapt}}$ **do**
 4:     Insert $\mathcal{M}$ into sequence; forward pass with $f_\theta$
 5:     Compute loss $\mathcal{L}$ on memory and response tokens
 6:     $\theta \leftarrow \theta - \eta_\theta \nabla_\theta \mathcal{L}$
 7:     $\mathcal{M} \leftarrow \mathcal{M} - \eta_{\mathcal{M}} \nabla_{\mathcal{M}} \mathcal{L}$
 8: **end for**
 9: Discard temporary memory $\mathcal{M}$ and freeze backbone $\theta$
10: **return** adapted backbone $f_\theta$

---

Table 5: Comparison of standard fine-tuning vs. fine-tuning with an adaptation phase.

| Model | Configuration | Argument F1 | | | |
| | | 2 calls | 3 calls | 4 calls | Avg. |
|---|---|---|---|---|---|
| | *Fine-tuning* | 77.3 | 72.6 | 55.8 | 68.6 |
| Llama 1B | *+ adapt* | 72.8 | 54.6 | 42.1 | 56.5 |
| | *+ adapt & all linear* | 74.1 | 66.7 | 68.4 | **69.7** |
| | *Fine-tuning* | 87.9 | 86.6 | 82.9 | **85.8** |
| Llama 3B | *+ adapt* | 78.9 | 69.5 | 63.2 | 70.5 |
| | *+ adapt & all linear* | 79.7 | 72.3 | 89.5 | 80.5 |
| | *Fine-tuning* | 87.7 | 86.8 | 88.2 | **87.6** |
| Llama 8B | *+ adapt* | 77.6 | 66.3 | 84.2 | 76.0 |
| | *+ adapt & all linear* | 84.5 | 78.3 | 89.5 | 84.1 |

## B.2    ANALYSIS OF FINE-TUNING WITH ADAPTATION PHASE

TokMem employs an adaptation phase where it is fine-tuned on a held-out auxiliary tools before training for new tools. In Section 3.3, our fine-tuning baseline dose not include this phase. To validate the fairness of our fine-tuning baseline, we investigate whether it could also benefit from this phase by training sequentially on the held-out auxiliary tools and the target tools.

As shown in Table 5, introducing the adaptation phase generally degrades performance compared with the standard fine-tuning baseline. This might be because of the interference between the disjoint held-out auxiliary and target tool sets with completely different functions and parameters. By contrast, TokMem avoids this interference issue by "abandoning" the auxiliary memory tokens after the adaptation phase.

Although adding the training capacity by finetuning with all linear layers alleviate this issue, we see that it still cannot outperform our original setup. This confirms that fine-tuning without adaptation is a stronger baseline for comparison.

## B.3    ANALYSIS OF UNFREEZING LLM BACKBONE FOR TOKMEM

We further examine the importance of freezing the backbone when adding new tool memories, reflecting real-world usage where procedural knowledge grows incrementally over time. This setting contrasts with recent approaches (Wang et al., 2025) that compress tool usage into virtual tokens by post-training the backbone. While such methods improve retrieval efficiency at scale, they rely on modifying backbone parameters, which hinders continual adaptability and risks overwriting prior knowledge.

As shown in Figure 6, unfreezing the backbone during TokMem adaptation leads to severe forgetting of previously learned tools, consistent with catastrophic interference in continual learning. By contrast, freezing the backbone preserves prior capabilities while allowing new tool memories to be incorporated without loss, highlighting TokMem's advantage for incremental adaptation. Notably,
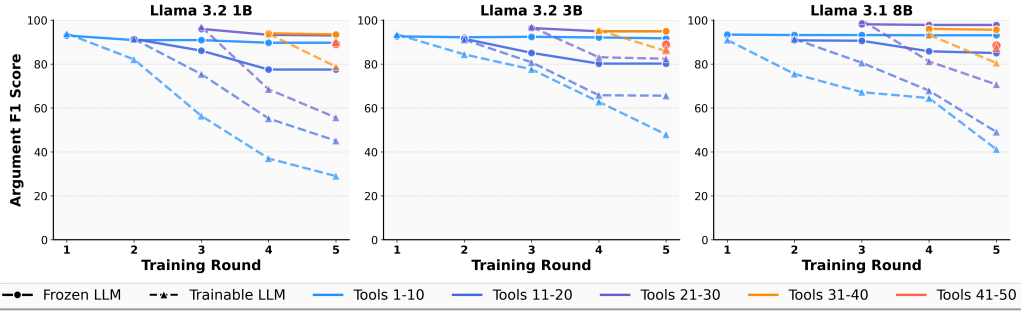
Figure 6: Comparison between Freezing and unfreezing the backbone. Allowing the backbone to update when adding new tool memories causes severe forgetting. Freezing preserves prior tools while enabling new ones.

Table 6: TokMem generalizes to longer tool-call chains at test time, significantly outperforming fine-tuning in zero-shot multi-step settings (training with 1 call and test with 2-4 calls).

| | | Test Time | | | |
|---|---|---|---|---|---|
| **Train Maximum Calls** | **Method** | 2 calls | 3 calls | 4 calls | **Avg.** |
| 1-call | *Fine-tuning* | 34.9 | 21.3 | 14.1 | 23.4 |
| | *TokMem* | 60.3 | 54.3 | 48.9 | **54.5** (+31.1) |
| 2-call | *Fine-tuning* | 86.2 | 78.8 | 64.8 | 76.6 |
| | *TokMem* | 82.0 | 81.8 | 82.3 | **82.0** (+5.4) |
| 3-call | *Fine-tuning* | 86.9 | 85.5 | 80.3 | 84.2 |
| | *TokMem* | 86.8 | 84.0 | 84.7 | **85.2** (+1.0) |
| 4-call | *Fine-tuning* | 87.9 | 86.6 | 82.9 | 85.8 |
| | *TokMem* | 85.9 | 86.7 | 88.3 | **86.3** (+0.5) |

unfreezing offers no meaningful performance gains after the initial training round, suggesting that TokMem strikes an effective balance between performance and continual adaptation.

## B.4 COMPOSITIONAL GENERALIZATION

We observe that TokMem provides clear advantages in compositional generalization over fine-tuning. Table 6 reports Argument F1 when the Llama 3B model is evaluated on queries requiring more function calls than those observed during training.

Notably, when trained solely on single-call data, TokMem achieves much stronger performance than fine-tuning when evaluating on 2 to 4 calls test data. This demonstrates that memory tokens trained for atomic procedures can be effectively composed at test time, enabling strong zero-shot generalization to multi-step behavior.

As the training regime is expanded to include more calls (e.g., up to 3 or 4), the performance gap narrows, but TokMem remains competitive or slightly ahead across all configurations. These results suggest that TokMem naturally supports compositionality, enabling flexible chaining of learned procedures without requiring task-specific fine-tuning.

## B.5 ANALYSIS OF NORM INFLATION FOR NEWER MEMORIES

We analyze the L2 norm of the learned memory embeddings when tools are introduced sequentially using Llama 3.2 3B. We follow the setup in Figure 4 by training TokMem with 5 rounds, adding 10 tools per round without our renormalization treatment. Figure 7 shows that newly added memory tokens gradually develop their L2 norms, which leads to competition with existing frozen tokens for the softmax operation for memory routing.
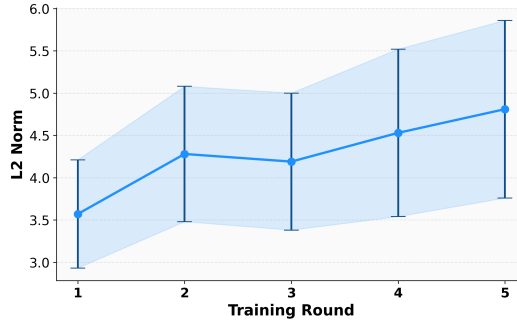
17

Figure 7: L2 norm of newly added memory tokens. Each round introduce 10 new tool memories. Error bar indicate standard deviation across the 10 tokens added in each round.

Table 7: Comparison of prefix tuning vs. TokMem condition embedding on GSM8K with two different size of Llama models. TokMem achieves higher compliance with required output formats and stronger exact-match accuracy than prefix tuning, especially in low-data regimes.

| | | Llama 3.2 1B | | Llama 3.2 3B | |
|---|---|---|---|---|---|
| **Data%** | **Method** | Compliance↑ | EM ↑ | Compliance↑ | EM ↑ |
| 20% | *Prefix tuning* | 0.0 | 0.0 | 45.9 | 33.1 |
| | *TokMem* | **98.0** | **37.7** | **94.6** | **65.6** |
| 100% | *Prefix tuning* | 82.8 | 30.0 | 97.2 | 64.1 |
| | *TokMem* | **97.4** | **39.1** | **98.2** | **66.9** |

## C  ADDITIONAL ANALYSIS ON MEMORY PLACEMENT

We have stress-tested the effect of memory token placement (prefix vs. infix) with randomly generated queries and with varying length of memory tokens in Section 3.4. We now turn to the GSM8K math reasoning dataset (Cobbe et al., 2021), to evaluate generalization and training efficiency.

Our experiments run on Llama 3.2 1B and 3B as backbone models and compare prefix tuning against TokMem under two training setups: using only 20% of the training set that represents a low-data regime, or the full dataset. We report two evaluation metrics.

- **Compliance** measures whether the model follows the required answer format, i.e., producing the final answer after the delimiter "####". This metric isolates the recall of procedural memory from the reasoning abilities already present in the backbone models.
- **Exact Match (EM)** measures the correctness of the final answer after standard normalization (e.g., removing commas or extraneous symbols).

As shown in Table 7, TokMem significantly outperforms prefix tuning, particularly in the low-data setting. With only 20% of the data, prefix-tuning fails to provide meaningful results, yielding zero compliance and EM on the 1B model and underperforming on the 3B model. By contrast, TokMem achieves near-perfect compliance and substantially higher EM scores across both backbones. When trained on the full dataset, prefix tuning improves considerably, yet TokMem continues to deliver stronger compliance and higher EM, underscoring its superior data efficiency and more reliable procedural control.

## D  DETAILS OF DATASETS

### D.1  DETAILS OF SUPER-NATURAL INSTRUCTION

We sample 1,000 English tasks from the SNI dataset, where each task is labeled with a task ID and a short descriptive name. The full list of sampled tasks is provided in Table 8. Tasks are introduced

to the model sequentially in ascending order of their IDs (e.g., the model first sees task 1, then task 2, and so on).

After training on the first $k$ tasks, we save a checkpoint and evaluate performance on the test sets of all $k$ tasks encountered so far. This simulates a continual learning setup where the model is expected to acquire new procedures while retaining previously learned ones. Once the model has been trained on all 1,000 tasks, it should be able to perform all of them without forgetting earlier tasks.

### D.2 Details of Function Calling dataset

For evaluating compositional memory recall, we sample 50 tools from the APIGen dataset (Liu et al., 2024b). The list of tools and their corresponding descriptions is provided in Table 8.

For each tool, we collect 50 query–call pairs, some of which may involve multiple calls to the same tool. This yields a total of $50 \times 50 = 2{,}500$ samples representing the non-compositional use of tools. To avoid data leakage, we split these samples into training and test sets with a 9:1 ratio.

On top of this, we synthesize complex queries by combining calls across different tools. These multi-step queries require the model to invoke multiple tools in sequence. We cap the number of synthesized samples at 5,000 for training and 500 for testing.

## E  Future Work

Our experiments are conducted on the research-oriented SNI and APIGen datasets, which allow controlled analysis of atomic and compositional recall. While these settings demonstrate the feasibility and effectiveness of tokenized procedural memory without backbone training, they do not fully capture the diversity of real-world procedures.

In particular, richer forms of composition, such as interleaving function calls from APIGen with NLP tasks from SNI, as illustrated in Figure 1b, and multi-turn interactions remain unexplored, but could be supported with curated datasets. Overall, advancing TokMem toward practical deployment will require more realistic benchmarks or user-driven data collection pipelines that better reflect open-domain procedural knowledge.

Additional future directions include incorporating reinforcement learning to improve generalization for complex compositional structures, and enabling personalization by allowing users to attach their own memory banks while keeping the backbone frozen. Together, these extensions pave the way for scalable, compact, and user-adaptive memory systems in large language models.

## F  Use of Large Language Models

We used ChatGPT as a general-purpose assistant to improve the writing of our paper, including grammar, readability, and clarity. Additionally, we used it to search for related work, which we then manually verified. All ideas, analyses, and conclusions presented in this paper are our own, and we take full responsibility for the content.

Table 8: Details of the sampled tools from the APIGen dataset, including their names and descriptions.

| ID | Tool | Description |
|---|---|---|
| 1 | auto_complete | Fetch auto-complete suggestions for a given query using the Wayfair API. |
| 2 | binary_addition | Adds two binary numbers and returns the result as a binary string. |
| 3 | binary_search | Performs binary search on a sorted list to find the index of a target value. |
| 4 | cagr | Calculates the Compound Annual Growth Rate (CAGR) of an investment. |
| 5 | calculate_factorial | Calculates the factorial of a non-negative integer. |
| 6 | calculate_grade | Calculates the weighted average grade based on scores and their corresponding weights. |
| 7 | calculate_median | Calculates the median of a list of numbers. |
| 8 | can_attend_all_meetings | Determines if a person can attend all meetings given a list of meeting time intervals. |
| 9 | cosine_similarity | Calculates the cosine similarity between two vectors. |
| 10 | count_bits | Counts the number of set bits (1's) in the binary representation of a number. |
| 11 | create_histogram | Create a histogram based on provided data. |
| 12 | directions_between_2_locations | Fetches the route information between two geographical locations including distance, duration, and steps. |
| 13 | fibonacci | Calculates the nth Fibonacci number. |
| 14 | final_velocity | Calculates the final velocity of an object given its initial velocity, acceleration, and time. |
| 15 | find_equilibrium_index | Finds the equilibrium index of a list, where the sum of elements on the left is equal to the sum of elements on the right. |
| 16 | find_first_non_repeating_char | Finds the first non-repeating character in a string. |
| 17 | find_longest_word | Finds the longest word in a list of words. |
| 18 | find_max_subarray_sum | Finds the maximum sum of a contiguous subarray within a list of integers. |
| 19 | find_minimum_rotated_sorted_array | Finds the minimum element in a rotated sorted array. |
| 20 | flatten_list | Flattens a nested list into a single-level list. |
| 21 | format_date | Converts a date string from one format to another. |
| 22 | generate_password | Generates a random password of specified length and character types. |
| 23 | generate_random_string | Generates a random string of specified length and character types. |
| 24 | get_city_from_zipcode | Retrieves the city name for a given ZIP code using the Ziptastic API. |
| 25 | get_pokemon_move_info | Retrieves information about a Pokémon's move using the PokéAPI. |
| 26 | get_product | Fetches product details from an API using the given product ID. |
| 27 | get_products_in_category | Fetches products in a specified category from the demo project's catalog. |
| 28 | greatest_common_divisor | Computes the greatest common divisor (GCD) of two non-negative integers. |
| 29 | integrate | Calculate the area under a curve for a specified function between two x values. |
| 30 | investment_profit | Calculates the profit from an investment based on the initial amount, annual return rate, and time. |
| 31 | is_anagram_phrase | Checks if two phrases are anagrams of each other, ignoring whitespace and punctuation. |
| 32 | is_leap_year | Checks if a year is a leap year. |
| 33 | is_palindrome | Checks if a string is a palindrome. |
| 34 | is_power | Checks if a number is a power of a given base. |
| 35 | is_rotation | Checks if one string is a rotation of another string. |
| 36 | is_valid_ip_address | Checks if a string is a valid IP address (IPv4). |
| 37 | is_valid_palindrome | Checks if a string is a valid palindrome, considering only alphanumeric characters and ignoring case. |
| 38 | is_valid_sudoku | Checks if a 9x9 Sudoku board is valid. |
| 39 | monthly_mortgage_payment | Calculates the monthly mortgage payment based on the loan amount, annual interest rate, and loan term. |
| 40 | note_duration | Calculates the duration between two musical notes based on their frequencies and the tempo. |
| 41 | place_safeway_order | Order specified items from a Safeway location. |
| 42 | polygon_area_shoelace | Calculates the area of a polygon using the shoelace formula. |
| 43 | potential_energy | Calculates the electrostatic potential energy given the charge and voltage. |
| 44 | project_population | Projects the population size after a specified number of years. |
| 45 | reverse_string | Reverses the characters in a string. |
| 46 | solve_quadratic | Computes the roots of a quadratic equation given its coefficients. |
| 47 | trapezoidal_integration | Calculates the definite integral of a function using the trapezoidal rule. |
| 48 | whois | Fetch the WhoIS lookup data for a given domain using the specified Toolbench RapidAPI key. |
| 49 | whole_foods_order | Places an order at Whole Foods. |
| 50 | wire_resistance | Calculates the resistance of a wire based on its length, cross-sectional area, and material resistivity. |

Figure 8: Overview of the 1,000 English tasks from the SNI dataset used in the atomic recall setting.