

Memory-Augmented Transformers: A Systematic Review from Neuroscience Principles to Enhanced Model Architectures

Parsa Omidi

Huawei Technologies Canada

parsa.omidi@huawei.com

Xingshuai Huang

Huawei Technologies Canada

xingshuai.huang@h-partners.com

Axel Laborieux

Huawei Technologies Switzerland

axel.laborieux@huawei.com

Bahareh Nikpour

Huawei Technologies Canada

bahar.nikpour@h-partners.com

Tianyu Shi

Huawei Technologies Canada

tianyu.shi@h-partners.com

Armaghan Eshaghi

Huawei Technologies Canada

armaghan.eshaghi@huawei.com

Abstract

Memory is fundamental to intelligence, enabling learning, reasoning, and adaptability across biological and artificial systems. While Transformer architectures excel at sequence modeling, they face critical limitations in long-range context retention, continual learning, and knowledge integration. This review presents a unified framework bridging neuroscience principles—dynamic multi-timescale memory, selective attention, and consolidation—with engineering advances in Memory-Augmented Transformers. We organize recent progress through three taxonomic dimensions: functional objectives (context extension, reasoning, knowledge integration, adaptation), memory representations (parameter-encoded, state-based, explicit, hybrid), and integration mechanisms (attention fusion, gated control, associative retrieval). Our analysis of core memory operations—reading, writing, forgetting, and capacity management—reveals a shift from static caches toward adaptive, test-time learning systems. We identify persistent challenges in scalability and interference, alongside emerging solutions including hierarchical buffering and surprise-gated updates. This synthesis provides a roadmap toward cognitively-inspired, lifelong-learning Transformer architectures.

1 Introduction

Memory is fundamental to both biological and artificial intelligence (AI), serving as the foundation for cognition, reasoning, and adaptive learning (Camina & Güell, 2017). In humans, memory enables the retention, retrieval, and manipulation of information across multiple time scales, supporting complex behaviors such as decision-making and problem-solving Wang et al. (2025a). This dynamic process integrates sensory inputs, transient processing, and long-term storage, forming a sophisticated cognitive architecture.

In AI, memory has become increasingly central as models evolve from static pattern recognition to more flexible, human-like cognition. Transformer architectures (Vaswani, 2017) have significantly advanced nat-

ural language processing, vision, and multimodal learning, yet their memory mechanisms remain restricted compared to the flexibility and efficiency of biological systems.

The primary limitation arises from self-attention’s quadratic complexity, constraining context window sizes. To stay within hardware limits, techniques like token pruning, sparse attention, and KV caching extend context but at a fidelity cost: sparse or approximate attention fractures long-range dependencies, and KV caches must evict or compress older entries, discarding vital information and harming coherence (Wang et al., 2024a). Another issue is the static nature of knowledge representation in standard Transformers. Once trained, their parameters are fixed, lacking mechanisms for continual learning or dynamic updates. This rigidity hinders adaptation to new information or user-specific contexts and risks catastrophic forgetting when fine-tuned, unlike the flexible updating seen in biological memory. Transformers also lag far behind biological systems in energy efficiency. The brain uses sparse, distributed, content-addressable memory with localized synaptic dynamics, operating on milliwatts of power (Prince et al., 2016; Gilbert & Brushfield, 2009). In contrast, Transformers require intensive computation: full-context inference scales quadratically with sequence length, while autoregressive decoding must process ever-growing KV caches with linear complexity per token. This computational burden results in orders-of-magnitude higher energy consumption.

To bridge these gaps, memory-augmented Transformers integrate neuroscience-inspired dynamic memory mechanisms. Human memory’s efficiency and adaptability increasingly guide Transformer design, particularly its integration across timescales: sensory memory (brief stimulus retention), prefrontal cortex-maintained working memory (short-term processing), and long-term memory (lifelong learning via neocortical-hippocampal networks). This architecture balances immediate processing with stable knowledge retention, while memory allocation is further regulated by salience and context-focusing attention only on relevant inputs, as described by the global workspace theory (Dehaene et al., 2011; Baars et al., 2021).

These neuroscience-derived principles increasingly shape memory-augmented Transformer architectures. Recent models incorporate multi-timescale memory, dynamic resource allocation, and plasticity-stability trade-offs, drawing explicit inspiration from hippocampal indexing, neuromodulatory gating, and hierarchical organization.

Recent review efforts have explored memory structures in AI models from various angles. For example, Ma et al. (2023) surveys memory augmentation techniques specifically in graph neural networks (GNNs), while Du et al. (2025) offers a broader perspective, covering diverse memory mechanisms across AI models, including long-term memory, long-context memory, parametric memory modification, and multi-source memory. He et al. (2024b) approaches the topic from a human-inspired perspective, focusing on long-term memory in AI models. Other surveys narrow their scope further: Shan et al. (2025) and Wu et al. (2025) focus on memory mechanisms in large language models (LLMs), while Zhang et al. (2024) specifically investigates memory in LLM-based agents. Similarly, Liu et al. (2025) includes a dedicated chapter on memory usage in foundation agents.

However, existing reviews are limited in two key ways. First, most rely on a single taxonomy to categorize memory-augmented methods, failing to provide a multidimensional or interdisciplinary understanding of the field. Second, many focus narrowly on specific model types or memory paradigms, such as long-term memory, LLMs, or agents, without addressing the broader landscape of memory integration across Transformer-based models.

In contrast, our review provides a comprehensive and interdisciplinary examination of memory augmentation techniques across Transformer models of various sizes, types, and applications. Our objectives are as follows:

- **Establish comprehensive taxonomies** linking neuroscience principles to memory mechanisms in Transformers from three different aspects.
- **Analyze core memory operations**, including reading, writing, forgetting, and self-management.
- **Identify current challenges** in memory-augmented Transformer design and highlight emerging paradigms and future directions inspired by biological memory.

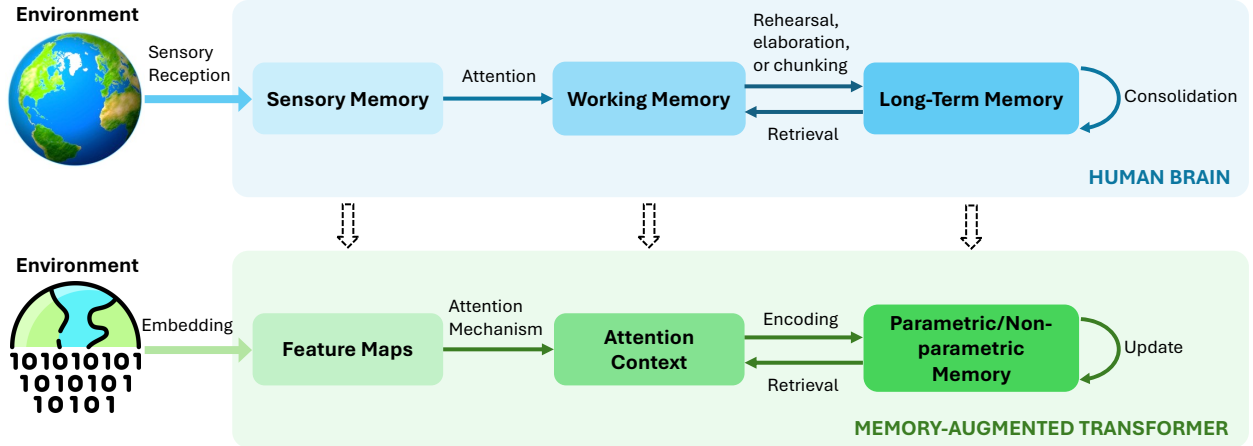


Figure 1: Parallels between the memory systems in the human brain and memory-augmented Transformers. Human memory consists of three interacting subsystems: sensory memory, working memory, and long-term memory. Memory-augmented Transformers mirror this architecture by leveraging embeddings, attention mechanisms, and advanced encoding and retrieval techniques to construct feature maps (analogous to short-term memory), attention contexts (analogous to working memory), and parametric or non-parametric memory (analogous to long-term memory).

By integrating insights from neuroscience and AI, this review aims to provide a conceptual framework and practical guidance for developing more efficient, adaptive, and cognitively inspired memory-augmented Transformers.

In the following sections, we begin by introducing memory architectures in biological cognitive systems, including the structure of human memory (Section 2.1), interactions between different memory systems (Section 2.2), and underlying computational principles (Section 2.3). Section 3 presents our proposed taxonomies from three perspectives: functional objectives (Section 3.1), memory types (Section 3.2), and integration techniques (Section 3.3). We then examine the mechanisms of memory operations adopted in the reviewed methods (Section 4), followed by a discussion of key challenges and future directions (Section 5).

2 Memory Architectures in Biological Cognitive Systems

Human memory operates as an interconnected, multi-layer network that stores, retrieves, and adapts information across several time-scales. Because these operations are hierarchical and widely distributed, stored knowledge is continuously reorganised, supporting rapid perception, flexible reasoning, and lifelong learning. This section reviews the biological architecture of memory and extracts principles that can inform cognitively-inspired AI models.

2.1 Architecture of Human Memory

Rather than a single store, human memory comprises three interacting subsystems, i.e., sensory, working, and long-term memory, as shown in the upper part of Figure 1. Each of them is optimised for a distinct combination of capacity, persistence, and processing depth (Cowan, 2008). Together they enable perception, decision-making, and learning across milliseconds to decades.

Sensory Memory: the Initial Buffer. Sensory memory provides a high-bandwidth, ultra-short buffer for raw perceptual input: visual traces (iconic) persist for ≈ 250 ms and auditory traces (echoic) for up to 2–3s (Reznik et al., 2023). In that brief window the brain analyses many stimuli in parallel; only items flagged by attention transition to working memory, while the rest decay rapidly, preventing overload. Neurally, these transient traces arise from sustained activity in primary sensory cortices and thalamo-cortical loops, organised

into modality-specific registers that filter noise and normalise signals before further processing (Camina & Güell, 2017). Transformers mimic part of this stage via token embeddings and positional encodings, which stabilise raw inputs for downstream layers. Yet, unlike biological circuits that adapt gain and leverage oscillations for temporal binding, current AI pipelines remain static, making robust perception under noise and context-dependent retention an open challenge for memory-augmented models.

Working Memory: the Cognitive Workspace. Working memory provides a transient, capacity-limited workspace that actively maintains and manipulates information required for reasoning, problem-solving, and goal-directed behaviour (Miller, 1956; Baddeley, 2003). Empirical estimates place its span at roughly four to seven “chunks,” a limit mitigated by chunking strategies and sustained by oscillatory activity in the prefrontal-parietal network.

Persistent firing, which is often organised through theta-gamma coupling, keeps multiple representations simultaneously accessible, while dopaminergic signals from the ventral tegmental area gate updates, suppress distractions, and prioritise task-relevant items (Roux & Uhlhaas, 2014). Cross-modal binding is supported by beta-band synchrony that links prefrontal cortex with hippocampal and sensory regions, enabling flexible recombination of auditory, visual, and spatial cues during complex tasks (Quak et al., 2015).

Functionally, the prefrontal cortex operates as a central executive, allocating attention, switching tasks, and coordinating specialised buffers (e.g., phonological loop, visuospatial sketchpad) (Russett et al., 2020). This distributed control balances stability with rapid updating, allowing the system to adapt to changing demands while avoiding interference.

Transformer self-attention partially echoes these operations by selectively weighting tokens within a fixed context window. Yet current models lack biologically inspired features such as neuromodulatory gating, oscillatory binding, and energy-efficient recall; external memories and recurrent variants narrow the gap but have yet to match the flexibility and robustness of human working memory.

Long-Term Memory: the Knowledge Repository. Long-term memory (LTM) is the brain’s durable storehouse, capable of retaining knowledge and experience for years or even a lifetime. Its defining strength is persistence: after consolidation, a trace can remain accessible indefinitely, provided it is periodically reactivated. Information is organised hierarchically into interconnected schemas that accelerate retrieval and support broad generalisation, yet the system stays plastic because each act of recall can render a trace temporarily labile and open to updating before it is re-stored during reconsolidation (Luo et al., 2022; Lee et al., 2017).

Two complementary consolidation processes underpin this durability. **Synaptic consolidation**, completed within hours, strengthens hippocampal circuits through activity-dependent events such as calcium spikes and sharp-wave ripples (Mujawar et al., 2021). **Systems consolidation** unfolds over days to years, as coordinated oscillations during sleep transfer memory indices from the hippocampus to distributed neocortical networks, creating resilient, cortex-based representations that can survive hippocampal damage (Luo et al., 2022).

LTM comprises episodic and semantic subsystems. **Episodic memory** records personally experienced events tied to a specific time and place and relies on hippocampal pattern completion for cue-based recall. **Semantic memory** stores abstract facts and concepts in widely distributed cortical networks, allowing individuals to answer questions like a capital city’s name without re-living the original learning episode (Kumar, 2021). The interplay of these subsystems enables both vivid recollection and flexible inference.

Adult neurogenesis in the dentate gyrus adds further adaptability, inserting new neurons that improve pattern separation and support the incorporation of novel information without erasing older traces (Anacker & Hen, 2017). This continual renewal helps the brain distinguish similar experiences and maintain cognitive flexibility across the lifespan.

Current AI systems approximate LTM with a mix of parameter-encoded knowledge and external memories. Parameter storage offers instant access but is costly to update, whereas external key-value banks such as Memformer’s fixed-size slots (Wu et al., 2020) or EMAT’s compressed QA memories (Wu et al., 2022b) allow on-the-fly writes and reads at inference time. Retrieval-Augmented Generation (RAG) (Lewis et al., 2020)

extends this idea by fetching fresh documents from external indices before every response, giving models a dynamic knowledge base. Despite these advances, artificial LTM still suffers from limited consolidation and vulnerability to catastrophic forgetting when new data overwrites old weights (Ranjith & Baskaran, 2024). Closing this gap will require biologically inspired mechanisms, e.g., dynamic consolidation, adaptive forgetting, and hierarchical memory layouts, that mirror the robustness and context sensitivity of human long-term memory.

2.2 Interactions Between Memory Systems

Human memory functions as a dynamic network, where sensory, working, and long-term stores communicate continuously to maximize learning and behaviour. Instead of isolated modules, these systems exchange activity through converging cortical–subcortical loops that adapt to context, attention, and emotional salience.

Encoding, Consolidation, and Retrieval. Encoding begins when sensory traces reach the prefrontal cortex, which filters and amplifies task-relevant inputs before they flow into working and long-term stores. Subsequent consolidation—particularly during slow-wave sleep—relies on hippocampal replay that drives neocortical reorganisation, stabilising both episodic and semantic traces (Klinzing et al., 2019). Retrieval completes this cycle: a partial cue reactivates hippocampal indices, triggering pattern-completion processes that reconstruct the distributed cortical representation and return it to working memory for use or further updating Teyler & Rudy (2007).

Top-Down and Bottom-Up Modulation. During retrieval, top-down signals from prefrontal regions bias processing toward current goals, suppressing irrelevant information, while bottom-up inputs from sensory and limbic areas flag novelty or emotional significance. Neuromodulators such as dopamine and acetylcholine strengthen synapses that encode behaviourally important events, fine-tuning what is stored or updated (Gazzaley & Nobre, 2012).

Emotional and Multimodal Integration. Emotionally charged or multisensory experiences recruit coordinated activity in the amygdala, hippocampus, and prefrontal cortex, yielding more persistent memories (Dolcos et al., 2004). The thalamus binds inputs from different senses, and hippocampal pattern completion links them into context-rich episodes that can be triggered later by a single cue.

Competitive and Co-operative Dynamics. Memory systems shift between competition and cooperation. Under stress or heavy cognitive load, control can pass from flexible hippocampal networks to faster, habit-based striatal circuits, ensuring rapid action (Schwabe & Wolf, 2013). In calmer conditions, episodic and semantic stores collaborate: detailed recollections supply context while abstract schemas guide generalisation and planning (Moscovitch et al., 2016).

Default Mode Network and Predictive Processing. The Default Mode Network supports offline consolidation, autobiographical recall, and mental simulation (Higgins et al., 2021). By replaying prior experiences, it updates internal models, enabling predictive processing that helps the brain (and, by extension, AI agents) anticipate future events and adapt behaviour accordingly (Liu et al., 2021).

Understanding this balance of stability and plasticity offers a template for AI: memory architectures that coordinate fast buffers with slower, more permanent stores, employ selective gating, and integrate cross-modal information can move beyond static storage toward lifelong, context-aware learning.

2.3 Computational Principles from Biological Memory

The architectural and functional properties of biological memory systems reveal fundamental computational principles that guide memory-augmented transformer design. These principles address universal computational challenges: managing limited resources, balancing stability with plasticity, and coordinating information flow across multiple timescales. Abstracting these neurobiological solutions into engineering heuristics yields a practical design playbook, now guiding the development of the most effective memory-augmented Transformer architectures.

Hierarchical Resource Allocation. Biological memory demonstrates that computational efficiency emerges from hierarchical organization rather than uniform processing (Hasson et al., 2015). Sensory

memory’s high-bandwidth, ultra-short retention enables parallel pre-processing, working memory’s capacity-limited workspace allows flexible manipulation, and long-term memory’s distributed storage supports both rapid recall and gradual consolidation. Multimodal evidence suggests that these hierarchical dynamics emerge as a global organizing principle of mammalian brains, with cortical timescale gradients topographically mirrored in striatum, thalamus, and cerebellum (Raut et al., 2020). *This hierarchical structure suggests that artificial systems benefit from multi-tier memory architectures that match storage characteristics to computational demands.*

Attention-Memory Bidirectional Coupling. The interaction between attention and memory reveals a crucial computational principle: memory systems both shape and are shaped by attentional mechanisms (Chun & Turk-Browne, 2007). Extensive evidence demonstrates that attention and memory cannot operate without each other: memory has limited capacity and attention determines what will be encoded, while memory from past experience guides what should be attended. Brain areas important for memory, such as the hippocampus and medial temporal lobe structures, are recruited in attention tasks, and memory directly affects frontal-parietal networks involved in spatial orienting. This bidirectional coupling enables adaptive resource allocation and context-sensitive processing through attention-dependent coupling between forebrain and brainstem neuromodulatory systems (Cicero et al., 2025). *These principles suggest AI memory systems should incorporate feedback loops between retrieval mechanisms and encoding processes.*

Neuromodulatory Gating and Significance Filtering. Biological memory formation relies on the interplay between Hebbian plasticity and neuromodulatory systems, making memory encoding inherently state-dependent and gated by the behavioral significance of information (Bazzari & Parri, 2019). Neuro-modulators such as dopamine and acetylcholine have distinct and complementary roles: dopamine regulates the induction of synaptic plasticity by modulating glutamatergic signaling, while acetylcholine orchestrates neuronal activity at both synaptic and network-wide levels. These systems establish computational principles of selective attention and adaptive thresholding, which allow the brain to prioritize salient information for encoding. Notably, selective neuromodulatory gating reflects a fundamental asymmetry in biological cognition: less than 5% of brain activity is devoted to conscious processes, while over 95% operates unconsciously, thus maximizing efficiency and resource allocation Raichle et al. (2001) Nail (2021). The strict limitations of working memory—estimated at 4–7 meaningful chunks—essentially define the conscious mind’s computational budget. *This insight suggests that artificial memory systems should carefully reserve costly, conscious-like processing for high-priority tasks such as novelty detection or conflict resolution, while routine memory operations are best delegated to automatic, parallel processing pathways analogous to the brain’s unconscious majority.*

Replay-Based Consolidation and Interference Management. The brain’s solution to the stability-plasticity dilemma through dual-phase consolidation, i.e., rapid hippocampal encoding followed by gradual neocortical integration, reveals essential computational principles for managing memory interference (Squire et al., 2015). Neural replay during sleep drives consolidation by reactivating patterns of network activity that occurred during previous experience, leading to potentiation of relevant synaptic connections in the cortex. This process enables rapid learning without catastrophic forgetting through replay-based consolidation and systems-level reorganization, where hippocampal replay propagates to cortex with reprocessing to extract statistical overlap from different encoding episodes. The stability-plasticity dilemma reflects a fundamental challenge in learning systems: retaining stored memory while learning new information (Mermillod et al., 2013). *The implication for AI is that effective memory systems require complementary fast and slow learning mechanisms with explicit consolidation phases.*

Content-Addressable Associative Retrieval. Biological memory systems excel at content-addressable retrieval through associative networks that enable pattern completion from partial cues (Rolls, 2013). The hippocampal CA3 subfield functions as an autoassociative network that stores experiences as memories, with abundant recurrent connections exhibiting spike-timing-dependent plasticity that allows pattern completion and recovery of stored patterns from noisy cues (Kang & Toyozumi, 2024). Empirical evidence from direct hippocampal recordings reveals pattern completion mechanisms where reinstatement of encoding patterns occurs during successful recollection, linked to gamma power fluctuations that coordinate selection of target-relevant neurons (Staresina et al., 2016). The CA3 region exemplifies this through its ability to retrieve complete episodic memories from fragmentary input, while semantic memory supports flexible infer-

ence through conceptual associations mediated by distributed cortical networks. *These mechanisms suggest that artificial memory architectures should prioritize associative rather than positional indexing and support similarity-based retrieval that mirrors biological pattern completion processes.*

Cross-Modal Integration and Binding. Biological memory systems demonstrate sophisticated cross-modal integration capabilities essential for unified cognitive processing (Nyhus & Curran, 2010; Shi et al., 2023). Theta and gamma oscillations enable interaction between cortical structures and the hippocampus for encoding and retrieval of episodic memories, where cortical gamma oscillations bind relevant stimulus features for perceptual representations and gamma phase synchronization between cortical and hippocampal neurons provides the mechanism for encoding diverse cortical information into hippocampal representations (Nyhus & Curran, 2010). Evidence from spatial decision-making tasks reveals that hippocampal-prefrontal interactions show maximum coherence during cross-modal binding, with theta rhythm dynamically modulating neurons in both regions (Tavares & Tort, 2022). Cross-modal prediction is supported by indirect pathways mediated by higher-order areas that receive convergent sensory inputs (Shi et al., 2023), *suggesting that artificial memory systems should incorporate mechanisms for binding information across modalities through oscillatory coordination.*

These computational principles collectively point toward memory-centered cognitive architectures where memory systems serve as the substrate for all cognitive operations rather than passive storage devices. The memory-centered cognition perspective places an active association substrate at the heart of cognition, making prediction and priming based on prior experience fundamental aspects of processing. Understanding these principles provides the foundation for designing artificial memory systems that move beyond static storage toward dynamic, adaptive, and context-aware memory architectures that can support the flexible, hierarchical, and associative processing characteristic of biological cognition.

3 Taxonomy of Memory-Augmented Transformers

Memory-augmented Transformers aim to overcome the fixed-context and static-knowledge constraints of standard models by drawing inspiration from the dynamic nature of human memory. This section presents a taxonomy of existing architectures along three dimensions: functional objectives, memory types, and integration technique (Figure 2). We relate these categories to biological memory principles to show how they help bridge the gap between current Transformers and human-like cognition.

3.1 Categorization by Functional Objectives

Memory-augmented Transformers address fundamental AI challenges, each mapped to a distinct functional objective:

Temporal Context Extension. Transformers struggle to process sequences beyond a fixed window, unlike the human brain’s ability to integrate experiences over long timescales. The evolution of temporal context extension reveals a clear trajectory from static windowing mechanisms toward sophisticated, biologically-inspired adaptive memory systems.

Sliding Window Attention (SWA) (Beltagy et al., 2020) establishes the foundational approach to linear-complexity context extension, where each token attends to a fixed window of w neighboring tokens, achieving $\mathcal{O}(n.w)$ complexity while maintaining parallelization. However, SWA operates as a static sensory buffer without adaptive selection or contextual awareness, limiting its effectiveness for complex temporal dependencies.

The field has witnessed a clear evolutionary trajectory from this static windowing toward adaptive, memory-augmented mechanisms that progressively incorporate working memory principles. ABC (Attention with Bounded-Memory Control) (Peng et al., 2021) transforms static windowing through learned, contextualized control strategies that dynamically determine token retention within fixed memory budgets. Transformer-FAM (Hwang et al., 2024) introduces feedback attention loops creating sustained activations across unlimited contexts, effectively transforming static windows into dynamic working memory systems. NAMMs (Cetin et al., 2025) employ STFT spectrogram analysis with genetic algorithms to evolve attention patterns for

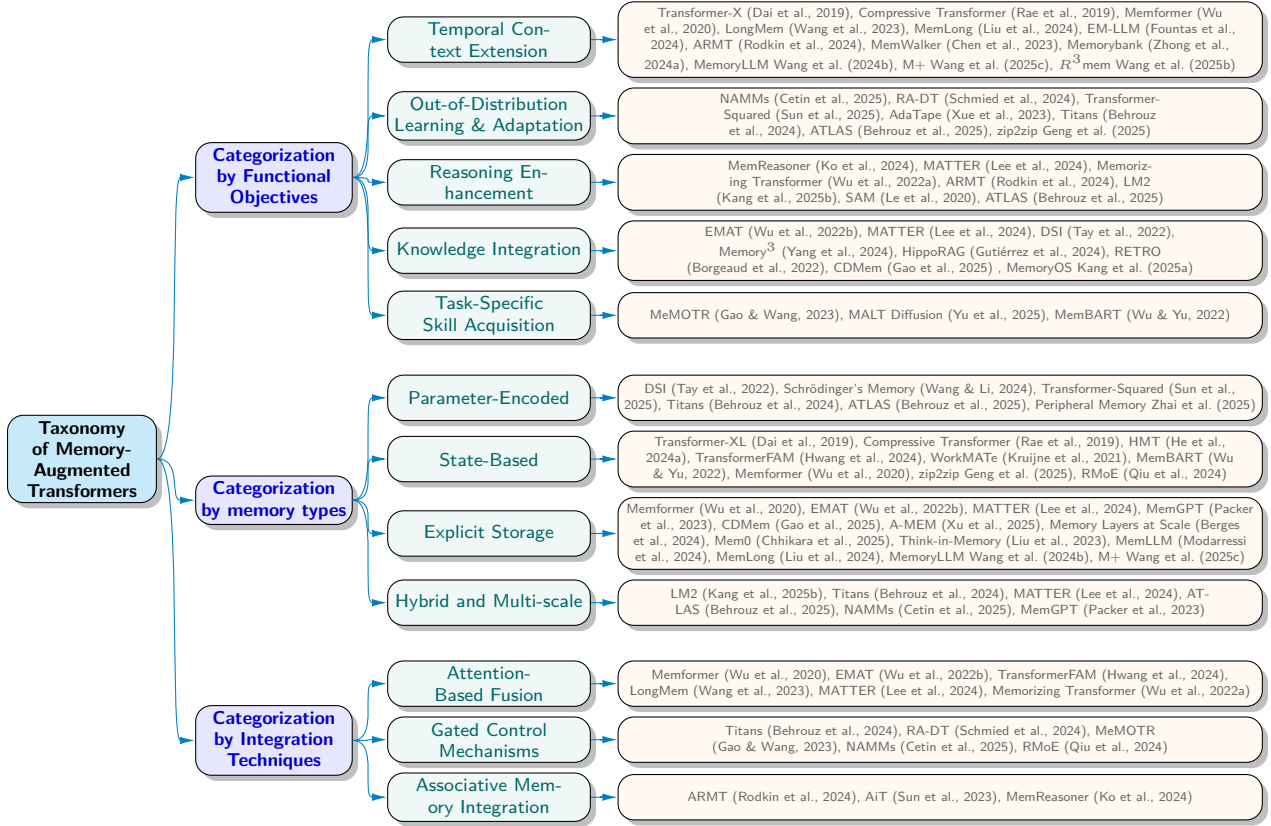


Figure 2: Taxonomy of Memory-Augmented Transformers.

zero-shot cross-modal transfer, while AdaTape (Xue et al., 2023) extends this through adaptive tape tokens that dynamically adjust sequence content and computational allocation. ATLAS (Behrouz et al., 2025) represents the culmination of this trend, adding sophisticated memory mechanisms after each sliding window through the Omega rule and polynomial feature mapping, achieving super-linear memory capacity.

Building upon these adaptive windowing foundations, practical implementations began with simple KV caching mechanisms. Transformer-XL (Dai et al., 2019) introduced the concept of caching key-value pairs from previous segments with relative positional encoding, where the cache itself functions as the memory, storing compressed representations of past context to extend processing beyond fixed windows. This established the fundamental principle that memory in transformers is essentially intelligent caching. Successive approaches focused on making this cache more selective and efficient at storing important information. Compressive Transformer (Rae et al., 2019) enhanced the basic KV cache through vector compression, boosting temporal range by 38% by intelligently compressing older cached states rather than simply discarding them. However, it still operated under fixed storage constraints, requiring memory eviction strategies. MemoryLLM Wang et al. (2024b) elevated the idea by adding a learnable write-gate, compression-on-evict, and a neural router that selects the top-k relevant keys, enabling ≈ 20 k-token context with near-constant compute. M+ Wang et al. (2025c) then removes this 20 k ceiling: it splits the cache into a small on-GPU working store and a large CPU-resident long-term bank, coordinated by a co-trained retriever and read-write scheduler. The hierarchy preserves the compress-on-evict principle yet sustains coherent generation across >160 k tokens while adding $<3\%$ throughput overhead. R^3 mem Wang et al. (2025b) introduces a reversible compression architecture that enables bidirectional transformation between raw context and compressed memory representations through hierarchical chunking across multiple semantic levels—segmenting content from paragraphs to sentences to sub-sentence units—ensuring both efficient compression and faith-

ful reconstruction of long contexts while maintaining semantic coherence across compression-decompression cycles.

Memformer (Wu et al., 2020) represented a breakthrough by decoupling computation from memory through similarity-based cache management—the cache became truly adaptive, updating based on content relevance rather than simple temporal recency. Its MRBP optimization cut training memory costs by 55% by learning which cached representations were most valuable to retain. LongMem (Wang et al., 2023) and MemLong (Liu et al., 2024) further refined cache intelligence: LongMem freezes the backbone LLM and uses a trainable SideNet to selectively retrieve and fuse the most relevant cached key-value pairs from a growing memory bank, while MemLong employs Retrieval Causal Attention to actively prune less important cached entries, demonstrating that the cache can learn what to forget as well as what to remember.

EM-LLM (Fountas et al., 2024) represents the pinnacle of intelligent caching through episodic memory segmentation. It uses Bayesian surprise detection to partition the cache into meaningful episodes, enabling retrieval that combines both semantic similarity and temporal contiguity across sequences up to 10 million tokens. ARMT (Rodkin et al., 2024) scales this concept to 50 million tokens through Hopfield-inspired associative caching with explicit erase operations, while MemWalker (Chen et al., 2023) creates hierarchical cache structures using trees of text summaries, and Memorybank (Zhong et al., 2024a) implements cognitively-inspired cache decay following human memory patterns like the spacing effect.

This progression mirrors the hierarchical integration of biological memory systems, where simple sensory buffering (SWA, basic KV caching) evolves into sophisticated working memory mechanisms (adaptive windowing, intelligent cache management) that balance immediate processing needs with longer-term contextual understanding, moving toward truly cognitive attention systems that integrate multiple timescales of temporal context.

Out-of-Distribution (OOD) Learning and Adaptation. Memory-augmented Transformers address the challenge of adapting to novel data distributions while preserving performance on familiar content through surprise-driven mechanisms that mirror biological memory systems’ ability to encode novel experiences while maintaining stable knowledge representations (Barry & Gerstner, 2024; Sinclair et al., 2021; Frank et al., 2022).

Surprise-driven adaptation forms the core of effective OOD learning. EM-LLM (Fountas et al., 2024) demonstrates this through Bayesian surprise detection and graph-theoretic boundary refinement to segment sequences into episodic events. This training-free approach automatically detects distribution shifts, creating distinct memory episodes for novel patterns while preserving performance on familiar content, embodying how novelty detection enables rapid adaptation to unexpected domains without compromising existing knowledge.

Titans (Behrouz et al., 2024) advance surprise-driven adaptation through prediction error gating, where KL divergence thresholds at the single-token level determine when memory updates occur. This token-by-token surprise detection enables fine-grained, test-time learning without parameter modification, allowing models to selectively memorize novel information while avoiding interference with established knowledge. ATLAS (Behrouz et al., 2025) similarly employs surprise signals through the Omega rule at the local context level, using sliding windows to determine which multi-token contexts warrant long-term memorization based on prediction error magnitudes across sliding windows rather than individual tokens. Dynamic Input Pruning (Federici et al., 2024) achieves zero-parameter adaptation through magnitude-based pruning of MLP activations per token, implementing a predictor-free strategy for real-time efficiency improvements without model retraining.

zip2zip Geng et al. (2025) demonstrates adaptive tokenization as a novel OOD adaptation strategy, dynamically expanding vocabulary at inference time through compression-based token merging that enables models to efficiently process unfamiliar token patterns and domains without retraining, achieving significant latency improvements while maintaining adaptability to new linguistic distributions.

Evolutionary and adaptive mechanisms enable cross-domain generalization without domain-specific training. NAMMs (Cetin et al., 2025) demonstrate evolutionary optimization of attention patterns using genetic algorithms and STFT spectrogram analysis to evolve token retention policies for zero-shot cross-modal trans-

fer. AdaTape (Xue et al., 2023) extends adaptive allocation through elastic input sequences with adaptive tape tokens that dynamically adjust sequence content and computational allocation based on problem complexity. RA-DT (Schmied et al., 2024) combines episodic memory with surprise-based pruning to retain high-error experiences, boosting multitask efficiency by 40% while mimicking dopaminergic learning mechanisms. Transformer-Squared (Sun et al., 2025) encodes procedural expertise directly into parameter space using SVD, dynamically blending expert vectors during inference to reach 90% accuracy on unseen tasks, despite a 15% latency overhead. Dutta & Sra (2024) extends the Memformer architecture to procedural computation by storing and combining past optimization gradients with learned coefficients, achieving 98% convergence accuracy on OOD tasks through trial-and-error learning patterns that mirror biological motor learning systems.

These approaches collectively demonstrate that effective OOD adaptation requires selective memory updating mechanisms that balance novelty detection with stability preservation, enabling memory-augmented Transformers to achieve flexible adaptation characteristics of biological memory systems while maintaining computational tractability and avoiding catastrophic forgetting.

Reasoning Enhancement. Extended context windows fundamentally enhance reasoning capabilities by providing access to larger knowledge bases and longer chains of inference (Yang et al., 2025). However, the relationship between context length and reasoning performance is not simply linear - it requires sophisticated memory mechanisms to maintain coherence across extended sequences, as standard attention mechanisms struggle with long-range dependencies and coherence degradation in extended contexts (Press et al., 2021).

Memory-augmented Transformers address these challenges by integrating scattered information over extended contexts for multi-hop inference and relational reasoning. Several approaches demonstrate significant improvements in reasoning performance through different memory architectures. MemReasoner (Ko et al., 2024) bridges encoders and decoders using a temporally-aware memory module with bidirectional GRUs and iterative updates, improving multi-hop QA by 18%. MATTER (Lee et al., 2024) unifies unstructured text and QA pairs into neural memories, using a cross-encoder to link questions to relevant data, boosting throughput by 100 \times and HotpotQA accuracy by 12%.

For tasks requiring extensive retrieval and pattern completion, associative memory approaches show particular promise. The Memorizing Transformer (Wu et al., 2022a) uses a kNN-retrievable memory to dynamically integrate distant context for tasks like theorem proving and code generation, scaling to 262K tokens while outperforming baselines in long-range reasoning, mirroring hippocampal episodic retrieval for problem-solving. ARMT (Rodkin et al., 2024) scales reasoning across 50 million tokens with associative memory blocks for pattern completion and interference mitigation, echoing hippocampal mechanisms. Self-Attentive Associative Memories (SAM) (Le et al., 2020) uses dual memory units and outer-product attention for updating item and relationship memories, improving performance on graph and geometric reasoning tasks such as the Traveling Salesman Problem and shortest path finding.

Gated memory mechanisms provide another effective approach to reasoning enhancement. LM2 (Kang et al., 2025b) adds memory modules with gated mechanisms to each decoder layer, outperforming standard Transformers on multi-hop reasoning over 128k-token contexts. ATLAS (Behrouz et al., 2025) shows that optimal context memorization enables complex reasoning by learning which historical information remains relevant for current inferences.

Alternative architectures explore hierarchical reasoning approaches that prioritize computational depth over extended context. HRM (Wang et al., 2025a) addresses reasoning depth through hierarchical convergence using coupled recurrent modules that achieve enhanced computational depth for problems requiring extensive search and backtracking. For resource-constrained settings, Memory-R+ (Le et al., 2025) enhances reasoning in tiny LLMs (≤ 1 B parameters) through dual episodic memory modules that provide intrinsic rewards for exploration and exploitation, achieving 2-14% performance improvements.

The key insight is that reasoning benefits from both quantity and quality of accessible context. Longer contexts provide more potential information, but sophisticated memory mechanisms are required to selectively attend to relevant information while avoiding interference from irrelevant details.

Knowledge Integration. Knowledge Integration encompasses the synthesis, storage, and retrieval of diverse information types into unified representations that support reasoning and generation. This process requires sophisticated indexing mechanisms that enable models to dynamically combine information from multiple sources for context-aware generation.

Retrieval-Augmented and Hierarchical Approaches demonstrate efficient knowledge incorporation strategies. RETRO (Borgeaud et al., 2022) combines frozen BERT retrievers with differentiable cross-attention, achieving GPT-3 performance with $25\times$ fewer parameters through access to 2 trillion token databases. CDMem (Gao et al., 2025) implements hierarchical three-stage encoding, i.e., expert, short-term, and long-term, through graph-structured, context-dependent indexing, achieving 85.8% success on ALFWorld and 56.0% on ScienceWorld by enabling multilevel knowledge recall tailored to current contexts.

Heterogeneous Memory Integration unifies diverse knowledge formats within a single architecture. EMAT (Wu et al., 2022b) encodes millions of QA pairs into key-value memory using fast MIPS for sub-millisecond querying, improving Natural Questions performance from 25.8 to 44.3 EM while maintaining 1000 queries/s throughput. MATTER (Lee et al., 2024) integrates both unstructured and semi-structured sources into type-agnostic neural memories, achieving $100\times$ throughput improvement over conventional retrieve-and-read models.

Parameter-Encoded and Brain-Inspired Systems explore direct knowledge embedding and neurobiological architectures. Memory³ (Yang et al., 2024) converts a textual knowledge base into a memory bank which can be seen as a bank of sparse retrievable parameters, enabling smaller language models to match the performance of bigger models, as well as reducing hallucinations and increasing factuality. DSI (Tay et al., 2022) encodes document corpora directly into model weights, enabling direct query-to-document mapping. HippoRAG (Gutiérrez et al., 2024) brings knowledge integration to RAG systems through the construction of a concept graph inspired by the hippocampus, outperforming RAG in multi-hop QA by up to 20% while being $10\text{-}30\times$ cheaper and $6\text{-}13\times$ faster.

MemoryOS Kang et al. (2025a) introduces an operating system-inspired hierarchical memory architecture for AI agents, featuring three-tier storage (short-term for immediate context, mid-term for recent interactions, and long-term for persistent personal memory) managed through four core modules (Storage, Updating, Retrieval, and Generation), which enable evolutionary adaptation via heat-based segment prioritization and dialogue-chain FIFO mechanisms; this results in a 49.11% F1 score improvement on long-term conversational benchmarks like LoCoMo, outperforming baselines by enhancing factual consistency and personalization in extended dialogues.

These approaches demonstrate that effective knowledge integration requires semantic organization, efficient access patterns, and scalable architectures that handle massive knowledge bases while maintaining precision. The convergence of retrieval-augmented methods, hierarchical encoding, and neurobiologically-inspired designs enables memory-augmented Transformers to bridge static parametric models with dynamic knowledge systems capable of large-scale, multi-format integration.

Task-Specific Skill Acquisition. Task-Specific Skill Acquisition enables models to learn and apply procedural knowledge for specialized tasks—such as object tracking, video generation, or dialogue, by encoding operations for robust, context-aware performance. Notable architectures include MeMOTR (Gao & Wang, 2023), which uses object-specific long-term memory with exponential decay and confidence-based updates for multi-object tracking. MALT Diffusion (Yu et al., 2025) employs recurrent attention and memory vectors to generate temporally coherent videos over long durations. In dialogue, MemBART (Wu & Yu, 2022) preserves memory states across turns, enhancing response quality. These models demonstrate that specialized memory mechanisms, whether persistent, episodic, or stateful, are essential for robust skill acquisition and deployment, echoing the compartmentalization of procedural memory in biological systems.

3.2 Categorization by memory types

Memory-augmented Transformers can be systematically differentiated by memory types, each offering distinct computational and cognitive properties: parameter-encoded, state-based, explicit storage, and hybrid/multi-scale systems.

Parameter-Encoded Memory. Parameter-encoded memory systems store knowledge directly within model weights, analogous to synaptic consolidation in biological systems where knowledge becomes distributed across neural connections. This approach offers fundamental advantages including immediate access without external retrieval operations and unified processing where memory and computation share the same parameter space. However, capacity constraints emerge as a critical limitation since memory capacity is bounded by the number of parameters available for knowledge storage.

Training-time parameter encoding provides stable, consolidated knowledge but lacks adaptability. DSI (Differentiable Search Index) (Tay et al., 2022) revolutionizes retrieval by encoding entire document corpora directly into standard Transformer parameters, transforming traditional retrieval into a generative task where models learn direct query-to-document mappings through the existing attention and feedforward mechanisms. Schrödinger’s Memory (Wang & Li, 2024) reveals the latent memory capabilities of large language models, demonstrating that LLMs can reconstruct complete datasets from minimal contextual cues through parameter-encoded associations formed during training. The key insight is that memory exists in a "superposition" state, remaining hidden until specific contextual triggers activate associative recall patterns, much like human memory retrieval from partial cues. Memory³ (Yang et al., 2024) converts textual knowledge bases into explicit memory banks functioning as sparse retrievable parameters, implemented through specialized embedding layers and sparse attention mechanisms. The system uses aggressive sparsification techniques and two-stage pretraining to efficiently store 1.1×10^8 text chunks within modified feedforward networks that enable smaller language models to match larger model performance.

Test-time parameter learning represents a revolutionary advance where parameters adapt dynamically during inference, addressing the fundamental limitation of static knowledge storage. Titans (Behrouz et al., 2024) uses MLP-based memory with KL divergence thresholds for surprise-driven, real-time parameter updates, maintaining stability via gating mechanisms that prevent catastrophic interference. ATLAS (Behrouz et al., 2025) enhances test-time learning by expanding MLP capacity through polynomial feature mapping and employs the Omega rule for sliding window optimization, achieving super-linear memory growth without traditional gradient descent. Transformer-Squared (Sun et al., 2025) enables real-time task adaptation by encoding procedural expertise directly into parameter space using SVD decomposition of feedforward layers, dynamically blending expert vectors during inference through specialized MLP mixing networks. Similarly, Peripheral Memory Zhai et al. (2025) introduces a CPU-RAM analogous architecture where LLMs function as processors interfacing with parameter-encoded memory banks modeled through Kolmogorov-Arnold Networks, enabling dynamic memory operations controlled by internal model states while maintaining direct integration with the model’s parameter space. This approach demonstrates how parameter-encoded systems can capture and recombine procedural knowledge using adaptive feedforward architectures rather than just declarative facts.

The evolution from static parameter encoding (DSI, Schrödinger’s Memory) to dynamic parameter learning (Titans, ATLAS, Transformer-Squared) represents a paradigm shift toward adaptive parameter-encoded systems that combine the efficiency advantages of parameter encoding with flexible adaptation capabilities. While training-time approaches provide stable knowledge consolidation, test-time parameter learning enables real-time adaptation with capacity enhancement techniques like polynomial feature mapping addressing fundamental scalability constraints. This progression points toward future architectures where parameter-encoded memory becomes truly dynamic, supporting both stable knowledge consolidation and adaptive capacity expansion during deployment.

State-Based Memory. State-based memory maintains information through persistent activations or hidden states that carry forward across processing steps, fundamentally differing from parameter-encoded approaches in that memory resides in dynamic activations rather than static weights. This approach mirrors biological working memory systems where information is maintained through sustained neural firing patterns, enabling temporal continuity and context preservation across extended sequences.

Transformer-XL (Dai et al., 2019) pioneered this approach through segment-level recurrence, caching hidden states from previous segments with relative positional encoding to extend context beyond fixed windows. While achieving substantial improvements in perplexity and long-range dependency capture, this method requires significant memory resources as cached states accumulate. Compressive Transformer (Rae et al.,

2019) addressed memory intensity through compressed state buffering, maintaining recent states in full resolution while compressing older memories using learned functions, extending temporal range by 38% and reflecting biological memory’s tendency to retain vivid recent experiences while abstracting older information. Hierarchical Memory Transformer (HMT) (He et al., 2024a) extends this idea by layering three progressively coarser caches—token-, chunk-, and segment-level—on top of the basic recurrent buffer, allowing 100 K-token streams on a single GPU while hierarchically pruning stale activations and keeping recent ones in full detail.

Advanced state-based mechanisms have emerged with sophisticated memory management capabilities. TransformerFAM (Hwang et al., 2024) introduces feedback attention loops where each layer attends to its own latent representations from previous time steps, creating internal working memory that enables indefinite sequence processing with $\mathcal{O}(L)$ complexity. This sustained activation mechanism transforms static attention layers into dynamic working memory systems capable of maintaining coherent representations across arbitrarily long sequences. WorkMATE (Kruijne et al., 2021) implements biologically-inspired gated memory circuits controlled by internal actions through reinforcement learning, successfully handling hierarchical tasks like 12-AX where multiple context levels must be maintained simultaneously, demonstrating how selective gating enables multiple independent representations within shared activation spaces. RMoE Qiu et al. (2024) extends state-based memory concepts to mixture-of-experts routing, where GRU-maintained hidden states capture routing history across consecutive layers, enabling each routing decision to leverage accumulated routing patterns from previous layers for improved expert selection and utilization.

Specialized applications further demonstrate state-based memory’s versatility. MemBART (Wu & Yu, 2022) applies persistent memory states for dialogue modeling, preserving conversation context across multiple turns to enable coherent long-term interactions. Memformer (Wu et al., 2020) employs a unified memory approach combining internal state representations with external memory banks, using Memory Replay Backpropagation (MRBP) to optimize memory usage and reduce training costs by 55%, representing a hybrid between state-based and explicit storage approaches. The Hierarchical Reasoning Model (HRM) Wang et al. (2025a) implements dual recurrent modules operating at different timescales—high-level for abstract planning and low-level for detailed computations—achieving enhanced effective depth of deep computation while maintaining training stability, demonstrating near-perfect performance on complex reasoning tasks with only 27M parameters.

State-based memory systems offer fundamental computational advantages including temporal continuity through persistent activations, efficient information propagation across time steps, and seamless integration with sequential processing architectures. These systems excel at maintaining coherent information flow within processing sequences, enabling models to preserve contextual representations across multiple computation steps without external storage overhead.

However, state-based approaches face inherent architectural constraints: memory capacity is fundamentally bounded by hidden state dimensions, creating potential bottlenecks for complex information storage. Interference between heterogeneous information types stored within shared activation spaces poses additional challenges, as different memory contents must compete for the same representational resources. Furthermore, the computational overhead of maintaining and updating persistent states throughout processing can become substantial for long sequences.

Explicit Storage Memory. Explicit storage memory employs external modules for scalable information storage and retrieval, maintaining persistent memory banks that survive beyond individual inference sessions. Unlike parameter-encoded memory that stores knowledge within model weights or state-based memory that maintains information through activations, explicit storage systems utilize dedicated external storage modules that can be accessed, updated, and scaled independently of the core model architecture, analogous to hippocampal indexing where sparse representations point to distributed memory traces.

Foundational approaches established core principles through dedicated external modules. Memformer (Wu et al., 2020) pioneered fixed-size external key-value stores with similarity-based retrieval, demonstrating efficient integration with Transformer architectures. EMAT (Wu et al., 2022b) implements compact neural memories for structured knowledge storage with fast retrieval capabilities, while MATTER (Lee et al., 2024) integrates heterogeneous data sources into unified external memory frameworks, achieving orders of magnitude throughput improvements while maintaining type-agnostic storage capabilities.

Advanced external storage systems introduce sophisticated organization and management strategies. MemGPT (Packer et al., 2023) implements OS-inspired hierarchical storage with main context and archival stores managed through function calls, enabling unbounded context through intelligent paging policies. CD-Mem (Gao et al., 2025) exemplifies graph-structured external storage through context-dependent indexing that organizes agent experiences into comprehensive external knowledge bases. A-MEM (Xu et al., 2025) advances this through automatically generated and evolving memory notes that form dynamic external knowledge graphs capturing semantic relationships over time.

Specialized applications demonstrate external storage versatility across different domains. Memory Layers at Scale (Berges et al., 2024) embeds external key-value slots within Transformer layers using product-key lookup for web-scale deployment, while Mem0 (Chhikara et al., 2025) targets production environments through external memory systems that blend vector embeddings with graph-structured representations for persistent user-specific memory. Think-in-Memory (Liu et al., 2023) and MemLLM (Modarressi et al., 2024) construct external triplet memory systems storing subject-object-relation structures, enabling models to query relationships through external memory rather than parameter-encoded associations. MemLong (Liu et al., 2024) demonstrates context extension by retrieving past embeddings from external storage systems, handling up to 80K tokens while preserving core model parameters. Memory-R+ Le et al. (2025) demonstrates intrinsic motivation applications where separate success and failure memory modules use kNN-based retrieval to compute rewards that guide reinforcement learning in tiny LLMs.

Key distinguishing characteristics define explicit storage memory’s advantages. Memory persistence enables permanent knowledge banks that accumulate information across sessions, unlike temporary state-based memory. Independent scalability allows external modules to expand knowledge capacity without requiring proportional increases in core model parameters. Structured organization predominantly stores structured and semi-structured data through indexing schemes including vector similarity search, graph traversal, and hierarchical clustering, enabling efficient access to diverse knowledge types.

Unlike simple retrieval-augmented generation approaches, explicit storage memory-augmented Transformers implement tightly integrated, differentiable memory modules that enable end-to-end optimization and sophisticated memory management strategies. These systems provide seamless integration between external memory operations and model computations, supporting dynamic knowledge updates and context-aware memory management. Explicit storage architectures enable Transformers to scale beyond fixed context limitations, support continual knowledge integration, and provide efficient retrieval for complex reasoning applications through persistent, structured external memory systems that evolve independently of core model constraints.

Hybrid and Multi-Scale Memory Systems. Hybrid memory systems combine multiple memory types, including parameter-encoded, state-based, and explicit storage, within unified architectures, creating hierarchical memory organizations that leverage the complementary strengths of different memory mechanisms. This architectural approach mirrors the brain’s integration of multiple memory subsystems, where different temporal scales and storage mechanisms work together to support flexible cognition.

LM2 (Kang et al., 2025b) demonstrates sophisticated parameter-state hybrids by integrating external memory modules with learnable gates into each decoder layer, enabling dynamic coordination between internal representations and external storage. Titans (Behrouz et al., 2024) advance this by combining state-based attention with parameter-encoded long-term memory modules that adapt during test time, while MATTER (Lee et al., 2024) represents parameter-explicit hybrids that encode diverse knowledge into model weights while maintaining external retrieval capabilities. These examples demonstrate how different memory types can coexist within single architectures to handle both immediate processing needs and long-term knowledge access.

These systems establish memory hierarchies based on temporal characteristics, creating multi-tiered architectures where fast state-based memory handles immediate context, medium-speed explicit storage manages session-persistent information, and slow parameter-encoded memory provides consolidated knowledge foundations. MemGPT (Packer et al., 2023) exemplifies this through OS-inspired memory management that coordinates working context (state-based) with archival storage (explicit) through learned paging policies.

Advanced hybrid systems like ATLAS (Behrouz et al., 2025) and NAMMs (Cetin et al., 2025) implement dynamic memory allocation that adaptively distributes memory resources across different types based on task demands and surprise signals. ATLAS uses context-aware optimization to determine when information should transition between memory stores, while NAMMs employ evolutionary algorithms to optimize memory allocation patterns across modalities. This adaptive coordination demonstrates that effective memory systems require intelligent arbitration between memory types rather than static allocation schemes.

The evolution toward hybrid architectures represents a fundamental shift from single-memory-type systems toward cognitively-inspired memory ecosystems that mirror the distributed, hierarchical organization of biological memory. These systems achieve computational flexibility by combining the immediate access of parameter-encoded memory, the temporal continuity of state-based memory, and the scalable capacity of explicit storage within unified frameworks that can dynamically adapt to diverse cognitive demands.

3.3 Categorization by Integration Techniques

The effectiveness of memory-augmented Transformers depends not only on what they remember, but also on how stored knowledge is integrated into ongoing computations. Integration techniques determine how retrieved information influences model behavior and how memory states evolve, reflecting the sophisticated coordination mechanisms found in biological memory systems.

Attention-Based Fusion remains the primary method for integrating memory content, enabling dynamic selection and weighting of stored information. Memformer (Wu et al., 2020) pioneered cross-attention between layer activations and external memory banks, gating semantically salient tokens much like thalamo-cortical loops filter relevant information in the brain. EMAT (Wu et al., 2022b) accelerates this approach by issuing retrieval queries at early layers and propagating key-value pairs through decoder stages, achieving millisecond-scale throughput for real-time applications. TransformerFAM (Hwang et al., 2024) advances fusion through feedback attention loops within each layer, creating internal working memory that supports indefinitely long contexts without external cache management. LongMem (Wang et al., 2023) introduces hybrid fusion via its SideNet module, which decouples memory retrieval from backbone updates while adaptively blending live inputs with cached representations through residual connections. MATTER (Lee et al., 2024) demonstrates heterogeneous fusion by encoding diverse content types into fixed-length neural memories accessed through universal attention heads, while the Memorizing Transformer (Wu et al., 2022a) implements kNN-based attention over rolling buffers to approximate human-like recency bias with logarithmic complexity.

Gated Control Mechanisms implement neuromodulatory-inspired regulation of memory updates and retention, mirroring how biological systems selectively encode and maintain information. Titans (Behrouz et al., 2024) employs surprise-driven writes triggered by KL divergence thresholds, mimicking norepinephrine’s role in novelty detection and memory consolidation. RA-DT (Schmied et al., 2024) combines episodic memory with adaptive forgetting gates based on statistical surprise, reducing catastrophic forgetting by 40% in multi-task reinforcement learning scenarios. MeMOTR (Gao & Wang, 2023) integrates exponential decay with confidence-driven pruning for object tracking, replicating striatal pathway dynamics that balance stability with adaptability. NAMMs (Cetin et al., 2025) takes an evolutionary approach, using genetic algorithms to evolve token retention policies that balance stability and plasticity through GABAergic-like inhibition mechanisms. RMoE (Qiu et al., 2024) demonstrates how GRU-based gated control can enhance routing efficiency in Mixture-of-Experts architectures by leveraging historical routing patterns across layers, establishing dependencies between routing decisions to improve parameter efficiency and expert selection diversity.

Associative Memory Integration enables content-addressable recall and efficient pattern completion across large contexts, shifting from positional to semantic indexing. ARMT (Rodkin et al., 2024) implements Hopfield-inspired associative blocks for $\mathcal{O}(1)$ retrieval over 50 million tokens, directly mirroring hippocampal CA3 circuits’ role in relational memory and pattern completion. Associative Transformer (AiT) (Sun et al., 2023) employs low-rank memory priors as attractors within a global workspace architecture, mimicking cortical column dynamics and distributed representation schemes. MemReasoner (Ko et al., 2024) enhances

associative integration through bidirectional GRUs that support iterative read-update cycles, maintaining coherence across long documents through sustained memory interactions.

These integration strategies collectively represent a paradigm shift from fixed positional indexing toward content-sensitive memory access that bridges artificial attention mechanisms with neural memory systems. By implementing biologically-inspired fusion, gating, and associative mechanisms, memory-augmented Transformers achieve more flexible and context-aware information integration that approaches the adaptive capabilities of human cognition. The convergence of these techniques enables models to dynamically coordinate multiple memory systems while maintaining computational efficiency and biological plausibility.

4 Mechanisms of Memory Operations

Memory-augmented Transformers overcome fixed-context limits of standard architectures by integrating neuroscience and engineering advances for dynamic, scalable memory. This section reviews core mechanisms, i.e., reading, writing, forgetting, capacity optimization, and self-management/adaptation, highlighting key techniques and representative models from the recent literature.

Read Operations. Early neural memories such as the Neural Turing Machine (Graves et al., 2014), DNC (Graves et al., 2016), and Kanerva Machines (Wu et al., 2018) introduced content-based addressing, but modern memory-augmented Transformers refine the read step with specialised retrieval mechanisms tailored to massive stores. Memory Layers at Scale (Berges et al., 2024) replaces dense feed-forward blocks with trainable key-value layers that perform product-key lookup, giving sub-linear top-k search across billions of entries while preserving end-to-end differentiability. EMAT (Wu et al., 2022b) shows that maximum-inner-product search can return millions of QA pairs in sub-millisecond latency, letting the model integrate external knowledge at every decoding step without harming throughput. The Memorizing Transformer (Wu et al., 2022a) augments attention with approximate k-nearest-neighbour queries into a continually growing cache, scaling recall to 262 k tokens and matching the perplexity of much larger dense models.

Associative designs push retrieval to constant time: ARMT (Rodkin et al., 2024) stores tokens in Hopfield-style energy basins for $O(1)$ pattern completion over 50 M-token contexts, and AiT (Sun et al., 2023) adds low-rank priors that reconstruct missing tokens from partial cues, outperforming sparse Transformers on relational reasoning benchmarks. For multi-hop discourse, MemReasoner (Ko et al., 2024) iteratively re-reads a temporal memory with bidirectional GRUs until the readout stabilises, boosting long-document question answering. MemLong (Liu et al., 2024) couples local attention with retrieval-causal attention that selects semantically relevant chunks from an 80 k-token cache, maintaining single-GPU efficiency.

More adaptive schemes appear in CDMem (Gao et al., 2025), which navigates a graph-indexed memory to fetch task-specific subgraphs, and ABC (Peng et al., 2021), which learns neural policies that decide when and how deeply to probe memory rather than relying on fixed heuristics. Finally, NAMMs (Cetin et al., 2025) demonstrate that the attention matrix itself can encode reusable retrieval plans, enabling zero-shot read strategies that transfer across modalities. Together these mechanisms move the field from uniform similarity search toward context-sensitive, learned, and even evolutionary reading policies that approach the flexibility of biological episodic recall.

Write Operations. Memory-augmented Transformers now treat writing as an active, learned decision rather than an unconditional overwrite. Titans (Behrouz et al., 2024) triggers a write only when prediction-error-derived surprise exceeds a KL-based threshold, mirroring dopamine-gated consolidation and allowing the model to memorise rare events without destabilising prior knowledge. LM2 (Kang et al., 2025b) introduces per-layer input/forget/output gates around an external store, so each decoder layer decides in real time how much of its state should be committed, yielding controllable long-context reasoning without extra fine-tuning. Memformer (Wu et al., 2020) ports LSTM-style gates into a key-value memory, giving fine-grained retention and erasure that stabilise sequence modelling, while MeMOTR (Gao & Wang, 2023) adds exponential decay plus confidence gating to keep only high-value object tracks in video streams.

Beyond simple gating, A-MEM (Xu et al., 2025) writes “memory notes” that are later linked and evolved into a graph, creating a self-organising semantic store that grows with the agent’s experience. Memory Layers at Scale (Berges et al., 2024) spreads writes across product-key memory shards on multiple GPUs, enabling

Table 1: Mechanisms of memory operations in memory-augmented Transformers, with key techniques and representative models.

Operation	Key Mechanism	High-fidelity Representative Models
Read	Content-based addressing	Neural Turing Machine (Graves et al., 2014); DNC (Graves et al., 2016); Kanerva Machine (Wu et al., 2018)
	Specialised similarity search	Memory Layers at Scale (Berges et al., 2024); EMAT (Wu et al., 2022b); Memorizing Transformer (Wu et al., 2022a)
	Associative retrieval	ARMT (Rodkin et al., 2024); AiT (Sun et al., 2023); MemReasoner (Ko et al., 2024); MemLong (Liu et al., 2024)
	Adaptive graph / policy-driven reads	CDMem (Gao et al., 2025); ABC (Peng et al., 2021); NAMMs (Cetin et al., 2025)
Write	Surprise / uncertainty-gated writes	Titans (Behrouz et al., 2024); LM2 (Kang et al., 2025b); MeMOTR (Gao & Wang, 2023)
	LSTM-style input-forget gating	Memformer (Wu et al., 2020); WorkMATE (Kruijne et al., 2021); RMoE (Qiu et al., 2024)
	Confidence-filtered updates	A-MEM (Xu et al., 2025); MemBART (Wu & Yu, 2022); MemoryLLM (Wang et al., 2024b); M+ (Wang et al., 2025c); Memory-R+ (Le et al., 2025)
	Reinforcement / optimisation traces	Memformers (Dutta & Sra, 2024); ATLAS (Behrouz et al., 2025)
Forget	Selective pruning	MemLong (Liu et al., 2024); MeMOTR (Gao & Wang, 2023); Titans (Behrouz et al., 2024)
	Exponential decay	MeMOTR (Gao & Wang, 2023); LM2 (Kang et al., 2025b)
	Adaptive (gate-controlled) decay	ARMT (Rodkin et al., 2024); Memformer (Wu et al., 2020); MemoryBank (Zhong et al., 2024a)
	Surprise-triggered erase	Titans (Behrouz et al., 2024); EM-LLM (Fountas et al., 2024)
	Task-aware forgetting	RA-DT (Schmied et al., 2024)
Capacity	Learned compression	Compressive Transformer (Rae et al., 2019); MATTER (Lee et al., 2024); EMAT (Wu et al., 2022b); zip2zip (Geng et al., 2025)
	Hierarchical chunk / tree buffers	MemLong (Liu et al., 2024); LM2 (Kang et al., 2025b); Meaningful Memory (Zhong et al., 2024b); M+ (Wang et al., 2025c); HRM (Wang et al., 2025a)
	Sharded / product-key KV	Memory Layers at Scale (Berges et al., 2024)
Self-Management	Dynamic allocation at test time	Transformer-Squared (Sun et al., 2025); NAMMs (Cetin et al., 2025); Titans (Behrouz et al., 2024); Peripheral Memory (Zhai et al., 2025)
	Sub-system specialisation	MATTER (Lee et al., 2024); MemBART (Wu & Yu, 2022); MemoryOS (Kang et al., 2025a)
	Interference control	ARMT (Rodkin et al., 2024); MemReasoner (Ko et al., 2024); RA-DT (Schmied et al., 2024); Schrödinger’s Memory (Wang & Li, 2024)

continual learning at web scale without bottlenecks. MemBART (Wu & Yu, 2022) mitigates read–write interference in dialogue by running parallel attention streams and merging them through residual gates. In procedural settings, Memformers (Dutta & Sra, 2024) treat past optimisation gradients as first-class memory registers, letting the model cache and reuse computation traces during new tasks. ATLAS (Behrouz et al., 2025) pushes test-time learning further: its Omega rule adjusts memory weights over sliding windows with polynomial feature mapping, achieving super-linear capacity growth without gradient descent. Finally, WorkMATE (Kruijne et al., 2021) shows that reinforcement-learned gating policies can independently open or close multiple working-memory slots, supporting concurrent, interference-free storage of task rules. Collectively, these writing mechanisms shift the focus from passive storage to selective, context-aware, and scalable writing, a prerequisite for lifelong, low-interference memory in Transformer systems.

Forgetting Dynamics. Effective forgetting sustains continual learning by pruning obsolete traces and freeing capacity for salient information. Modern memory-augmented Transformers therefore implement selective, learned erase policies rather than indiscriminate decay. MemLong (Liu et al., 2024) prunes keys

whose retrieval counts fall below a threshold, ensuring its external cache stays focused on behaviourally relevant chunks. MeMOTR (Gao & Wang, 2023) adds confidence-weighted exponential decay so unreliable object tracks vanish naturally as a video unfolds. Titans (Behrouz et al., 2024), LM2 (Kang et al., 2025b), and Atlas (Behrouz et al., 2025) all regulate forgetting with adaptive gates: Titans couples KL-surprise with a trainable decay factor, LM2 ties gate strength to layer-wise uncertainty, and ATLAS uses the Omega-rule’s sliding-window optimisation—effectively down-weighting contributions from tokens that lie outside a polynomially mapped context window, allowing new information to replace stale traces without gradient descent.

Aggressive cleanup appears in ARMT (Rodkin et al., 2024), whose Hopfield memory periodically normalises and hard-deletes outdated vectors, preventing spurious attractors. Memformer (Dutta & Sra, 2024) mixes LSTM-style forget gates with memory-replay back-propagation so rarely used slots fade while important ones are refreshed, and MemoryBank (Zhong et al., 2024a) models retention with an Ebbinghaus-shaped counter that decays unless the entry is reaccessed. EM-LLM (Fountas et al., 2024) reinforces this trend by coupling prediction-error spikes to simultaneous write-and-prune cycles, mirroring neuromodulatory control of consolidation.

Together these mechanisms mark a shift from passive decay toward context-sensitive, learned forgetting that protects critical memories while continuously liberating capacity for new experiences.

Capacity Optimization. Capacity optimization addresses how memory-augmented Transformers expand storage and retrieval ability without linearly inflating computation or parameters. Current work converges on three complementary tactics—compression, hierarchy, and sparsity—to keep memory growth compatible with practical hardware budgets.

Compressive techniques shrink inactive activations or knowledge chunks before eviction. The Compressive Transformer (Rae et al., 2019) auto-encodes aged hidden states into coarse vectors, doubling usable context while holding FLOPs steady and matching baseline perplexity on WikiText-103 . At the knowledge level, EMAT (Wu et al., 2022b) and MATTER (Lee et al., 2024) map millions of QA pairs or mixed documents to short neural codes; maximum-inner-product search then delivers sub-millisecond retrieval without adding trainable weights.

Hierarchical organization spreads capacity across tiers with different granularity. MemLong (Liu et al., 2024) chunks sequences and prunes rarely accessed blocks, maintaining 80 K-token windows on a single GPU. LM2 (Kang et al., 2025b) builds tree-indexed memories that let local detail and global context be fetched at equal cost, sustaining reasoning over 128 K tokens. HMT (He et al., 2024a) stacks sensory, short-term, and long-term buffers, matching large long-context models while using $\approx 2\%$ of their parameters.

Sparse look-ups push size further by reducing per-query work. Memory Layers at Scale shards (Berges et al., 2024) product-key tables across GPUs, supporting billion-entry memories with sub-linear compute and intact end-to-end gradients . Dynamic Memory Compression Nawrot et al. (2024) learns head- and layer-specific KV sharing, cutting inference memory up to $4\times$ with negligible accuracy drop , while MLKV Zuhri et al. (2024) shares KV heads across layers to trim cache by up to $6\times$ at similar quality.

Together, these advances show that intelligent compression, hierarchical buffering, and sparse retrieval make large-capacity memory feasible, allowing even modest-sized Transformers to reason over book-length context or web-scale knowledge without prohibitive cost.

Self-Management and Adaptation. After compression and hierarchical layout tame raw capacity, the next hurdle is deciding how that capacity is used in real time. Recent models treat memory as an autonomous resource that can be allocated, specialised, or pruned during inference, bringing Transformers closer to the selective plasticity of biological cognition.

Transformer-Squared (Sun et al., 2025) routes activations through a pool of expert vectors selected on-the-fly, letting the model enlarge functional capacity without weight updates while preserving high accuracy on unseen procedural tasks. Titans (Behrouz et al., 2024) adds a neuromodulatory gate: only tokens whose KL-surprise clears a learned threshold are written, and low-surprise traces decay, reducing interference while the long-term store grows during deployment. ATLAS (Behrouz et al., 2025) generalises this to sliding

windows; its Omega rule re-weights entire spans, down-scoring stale patterns and allowing super-linear memory growth without gradient descent. NAMMs (Cetin et al., 2025) evolve layer-wise retention masks from attention statistics, trimming key-value caches by up to 80% yet improving long-context benchmarks through zero-shot transfer across modalities.

Interference control is handled by orthogonal or gated rewrites. ARMT (Rodkin et al., 2024) projects new vectors onto an orthogonal subspace before insertion, preventing outdated attractors and keeping $O(1)$ retrieval stable over tens of millions of tokens. MemReasoner (Ko et al., 2024) iteratively re-reads and updates a temporal store with bidirectional GRUs until the representation converges, preventing early facts from being overwritten and boosting multi-hop question answering on 128 k-token documents. RA-DT (Schmied et al., 2024) links episodic memory to a reinforcement-learning critic, retaining only high-error trajectories and lifting multi-task sample efficiency while bounding memory size. MemBART (Wu & Yu, 2022) isolates dialogue context from world knowledge through dual attention streams and residual gates, and Schrödinger’s Memory (Wang & Li, 2024) stores traces in a latent “superposition” that surface only when cued, lowering hallucination rates in factual probing.

Collectively, these systems replace static buffers with self-monitoring stores that learn what to remember, where to place it, and when to forget. By coupling dynamic allocation with interference-aware rewriting, they extend the compression and hierarchy tools of capacity optimisation into a full feedback loop, allowing Transformers to balance stability and plasticity throughout their lifetime.

5 Discussion, Challenges, and Future Directions

Memory-augmented Transformers have progressed from simple context extensions to sophisticated cognitive architectures, narrowing the gap between learning and memory. Our taxonomic analysis in Table 2 shows a rapid shift from static pattern recognition to adaptive, experience-driven intelligence. From the earliest systems in 2019 to today’s production-ready designs, development has converged toward hybrid storage, adaptive dynamics, and intelligent forgetting, while also exposing persistent challenges in scaling, evaluation, and integration. This chapter distills insights from that evolution, highlights constraints that limit current models, and outlines research directions to bridge artificial and biological memory systems—offering a roadmap toward architectures that not only extend computational capacity but also support genuine artificial cognition.

5.1 Overview and Synthesis

Evolutionary trajectory and convergence

- **Foundation (2019–2021):** Early systems established explicit memory management beyond standard attention via state-based recurrence and compression, demonstrating that long-range modeling benefits from persistent activations and hierarchical reduction (e.g., Transformer-XL Dai et al. (2019); Compressive Transformer Rae et al. (2019)).
- **Expansion (2022–2024):** Retrieval-augmented modeling scaled access from thousands to billions of entries using kNN/MIPS indexing and chunked cross-attention (e.g., Memorizing Transformer Wu et al. (2022a); RETRO Borgeaud et al. (2022); EMAT Wu et al. (2022b); Memory Layers at Scale Berges et al. (2024)), while architectures diversified to associative, hierarchical, and graph-based organization (e.g., AiT Sun et al. (2023), MemGPT Packer et al. (2023), MemWalker Chen et al. (2023), HippoRAG Gutiérrez et al. (2024), CDMem Gao et al. (2025)). Surprise-gated updates emerged as a biologically motivated write policy (e.g., Titans Behrouz et al. (2024)), complementing selective reset/decay and LRU strategies for stability under growth.
- **Maturation (2025):** Production-oriented designs emphasized hybrid storage, test-time adaptation, and specialized access, including expert routing and compression-based interfaces (e.g., zip2zip Geng et al. (2025)), operational memory OS abstractions (e.g., MemoryOS Kang et al. (2025a)), and hierarchical controllers for reasoning (e.g., HRM Wang et al. (2025a)), reflecting a shift from static pattern recognition to adaptive, experience-driven intelligence.

Table 2: Comprehensive feature matrix for memory-augmented Transformer models: Evolution from 2019-2025

Year	Model	Architecture			Generality		Memory Dynamics			Management	
		Storage Class	Integration Method	Backbone Compatibility	Input Modality	Memory Span	Write Trigger	Plasticity	Memory Scope	Retrieval Mechanism	Forgetting Mechanism
2019	Transformer-XL (Dai et al., 2019)	S	Wrp	×	T	S	Stc	F	Lyr	Attn	FIFO
	Compressive Transformer Rae et al. (2019)	S	Wrp	×	T	M	Stc	F	Lyr	Attn	Dec+FIFO
2020	SAM (Le et al., 2020)	E	Plg	✓	M	A	Pol	TT	Gbl	Outer	Rst
	Memformer Wu et al. (2020)	H(SE)	Plg	✓	M	L	G	TT	Gbl	Attn	Dec
2021	ABC Peng et al. (2021)	E	Plg	✓	T	S	Pol	F	Gbl	Pol	—
	WorkMATE (Kruijne et al., 2021)	S	Plg	✓	T	S	Pol	TT	Lyr	Attn	—
2022	EMAT (Wu et al., 2022b)	E	Plg	✓	T	L	Stc	F	Gbl	MIPS	—
	RETRO (Borgeaud et al., 2022)	E	Plg	✓	T	L	Stc	F	Gbl	kNN	—
	DSI (Tay et al., 2022)	P	Bsp	×	T	L	Stc	F	Lyr	Attn	—
	Memorizing Transf. Wu et al. (2022a)	E	Plg	✓	T	L	Stc	F	Gbl	kNN	—
	MemBART (Wu & Yu, 2022)	S	Plg	×	T	L	G	TT	Hrch	Dual	Rst
2023	LongMem (Wang et al., 2023)	E	Plg	✓	T	L	G	TT	Gbl	Attn	Prn
	MemGPT Packer et al. (2023)	E	Plg	✓	T	M	G	TT	Hrch	kNN	LRU
	Think-in-Memory (Liu et al., 2023)	E	Plg	✓	T	L	G	TT	Gbl	Trip	Dec
	AdaTape (Xue et al., 2023)	H(PS)	Plg	✓	T	S	Pol	F	Lyr	Tape	—
	MemWalker (Chen et al., 2023)	E	Bsp	✓	T	M	Pol	TT	Hrch	Tree	Rst
	AI ² (Sun et al., 2023)	H	Plg	✓	M	A	Pol	TT	Gbl	Assoc	Dec
	MeMOTR Gao & Wang (2023)	H(SE)	Bsp	×	M	L	Stc	TT	Gbl	Attn	Dec
2024	MemoryBank Zhong et al. (2024a)	E	Plg	✓	T	L	G	TT	Hrch	kNN	Dec
	TransformerFAM Hwang et al. (2024)	S	Plg	✓	T	S	Stc	F	Lyr	Attn	—
	HMT He et al. (2024a)	H(SE)	Plg	✓	T	L	Stc	F	Hrch	Attn	Dec
	MemoryLLM (Wang et al., 2024b)	H(SE)	Plg	✓	T	L	Stc	TT	Gbl	Attn	Dec
	HippoRAG Gutiérrez et al. (2024)	E	Plg	✓	T	L	Stc	TT	Gbl	Graph	Rst
	MATTER Lee et al. (2024)	H(PE)	Wrp	✓	T	M	Stc	F	Gbl	MIPS	—
	Memory3 Yang et al. (2024)	H	Plg	✓	T	M	Stc+Pol	F	Gbl	kNN	—
	ARMT (Rodkin et al., 2024)	E	Plg	✓	T	A	Pol	TT	Hrch	Assoc	Cyc
	MemLong Liu et al. (2024)	E	Plg	✓	T	L	G	TT	Gbl	kNN	Prn
	Schrödinger's Memory Wang & Li (2024)	P	Bsp	×	T	L	Stc	F	Gbl	Attn	—
	MemReasoner Ko et al. (2024)	E	Plg	✓	T	L	Stc	TT	Gbl	Attn	Rst
	EM-LLM (Fountas et al., 2024)	E	Plg	✓	T	L	Sur	TT	Gbl	Seg+kNN	Dec
	RA-DT (Schmied et al., 2024)	E	Plg	✓	M	L	Sur	TT	Gbl	Pol	Sel
	Memory Layers at Scale (Berges et al., 2024)	P	Bsp	×	T	L	Stc	F	Gbl	PK-MIPS	—
	Titans Behrouz et al. (2024)	H(PE)	Wrp	✓	T	M	Sur	TT	Gbl	Attn	Dec
2025	Transformer-Squared (Sun et al., 2025)	P	Wrp	✓	T	L	G	TT	Lyr	Attn	—
	LM2 Kang et al. (2025b)	H(SE)	Wrp	✓	T	M	G	TT	Hrch	Attn	Dec
	NAMMs (Cetin et al., 2025)	H(PSE)	Plg	✓	M	M	Pol	TT	Hrch	Pol	Prn
	R ³ mem (Wang et al., 2025b)	H(PE)	Plg	✓	T	M	G	TT	Gbl	Comp	Sel
	RMoe (Qiu et al., 2024)	S	Plg	✓	T	S	G	TT	Lyr	Hier	Sel
	Memory-R+ (Le et al., 2025)	E	Plg	✓	T	L	Sur	TT	Gbl	kNN	Sel
	Mem0 Chhikara et al. (2025)	E	Plg	✓	T	L	G	TT	Hrch	kNN+Graph	LRU
	CDMem (Gao et al., 2025)	E	Plg	✓	T	M	Pol	TT	Hrch	Graph	Rst
	ATLAS Behrouz et al. (2025)	H(PE)	Wrp	✓	T	M	Sur	TT	Gbl	Attn	Dec
	MemoryOS (Kang et al., 2025a)	E	Plg	✓	T	M	G	TT	Hrch	Seg	LRU
	zip2zip (Geng et al., 2025)	P	Wrp	✓	T	S	Pol	TT	Gbl	Comp	—
	Peripheral Memory (Zhai et al., 2025)	P	Plg	✓	T	S	G	TT	Gbl	Attn	—
	MALT Diffusion (Yu et al., 2025)	S	Plg	✓	M	L	Stc	F	Lyr	RecAtt	—
	A-MEM (Xu et al., 2025)	E	Plg	✓	M	L	G	TT	Gbl	kNN	Evol
	HRM (Wang et al., 2025a)	S	Bsp	×	T	M	Pol	TT	Hrch	Hier	Rst

Legend:

Architecture:

Storage Class: P = Parameter-encoded, S = State-based, E = External store, H = Hybrid, H(PS) = Parameter+State, H(PE) = Parameter+External, H(SE) = State+External, H(PSE) = All three
Integration Method: Plg = Plug-in, Wrp = Wrapper/adaptor, Bsp = Bespoke redesign
Backbone Compatibility: ✓ = Universal, × = Architecture-specific

Generality:

Input Modality: T = Text-only, M = Multi-modal

Memory Span: S = Short-term, L = Long-term, M = Multi-scale, A = Associative

Memory Dynamics:

Write Trigger: Stc = Static, Sur = Surprise-gated, Pol = Policy-learned, G = Gated
Plasticity: F = Fixed after training, TT = Test-time adaptable
Memory Scope: Lyr = Layer-local, Gbl = Global, Hrch = Hierarchical

Management:

Retrieval Mechanism: Attn = Attention-based, kNN = k-Nearest neighbor, Assoc = Associative, Graph = Graph-based, MIPS = Max inner product search, PK-MIPS = Product-key

MIPS, Pol = Policy-driven, Seg = Segmentation, Outer = Outer-product, Trip = Triplet-based, Expert = Expert-routing, Tree = Tree-based, Tape = Tape-based, RecAtt = Recurrent attention, Dual = Dual-stream, Comp = Compression-based, Hier = Hierarchical

Forgetting Mechanism: FIFO = First-in-first-out, Dec = Decay, Prn = Pruning, Cyc = Cycle-based, Rst = Reset, LRU = Least-recently-used, Sel = Selective, Evol = Evolutionary

Timeline: 2019 2020 2021 2022 2023 2024 2025
 Note: — = Not reported

Architecture: hybrid dominance. Parameter-encoded memory offers immediate access but risks catastrophic interference when updated; state-based memory supports rapid adaptation but is capacity-limited; external stores scale but add retrieval/consistency overhead. Hybrid designs increasingly combine these modalities to balance latency, scalability, and plasticity through division of labor and policy-driven coordination.

Memory dynamics: from rules to policies. Write operations progressed from static schedules to surprise-gated consolidation and learned policies, mitigating stability–plasticity trade-offs by updating on prediction errors and adapting allocation/eviction to task demands. Test-time plasticity became the default, enabling personalization and continual adaptation in deployment.

Retrieval and forgetting: specialization matters. Access evolved beyond attention and pure similarity toward structure-aware methods: graph navigation for relational queries, associative retrieval for content-addressable access, and hierarchical/expert routing for specialization and efficiency. Forgetting moved from FIFO/decay to LRU, selective, cycle-based, and evolutionary strategies, showing that intelligent erasure—aligned with utility and hierarchy—is as consequential as storage and retrieval for sustained performance under growth.

5.2 Challenges

Despite remarkable progress toward cognitive memory architectures, our comprehensive analysis reveals fundamental challenges that continue to constrain practical deployment and theoretical understanding of memory-augmented systems.

Scalability and Retrieval Bottlenecks. Despite significant architectural innovations, memory-augmented systems face fundamental scalability constraints that limit practical deployment at scale. Current approaches demonstrate distinct trade-offs between computational efficiency, retrieval accuracy, and resource requirements.

Retrieval mechanisms exhibit characteristic scaling limitations. Approximate similarity search methods, while computationally efficient, suffer from accuracy degradation as memory size increases (Wu et al., 2022b). Product-key decomposition approaches successfully reduce lookup complexity from linear to sub-linear scaling (Berges et al., 2024), yet encounter parameter overhead that constrains expansion to billion-entry systems. Graph-based retrieval methods enable sophisticated multi-hop reasoning but face exponential complexity growth with increasing graph density and traversal depth.

Compression-based solutions present complementary challenges. While techniques like adaptive tokenization can achieve substantial sequence reduction, they introduce inference-time computational overhead that may offset retrieval gains. Hierarchical memory organization similarly requires careful balance between compression efficiency and information fidelity.

Infrastructure considerations increasingly limit deployment viability. Memory-augmented architectures impose substantial storage bandwidth requirements, introduce novel security vulnerabilities through persistent external memory, and exhibit non-linear energy scaling patterns. Distributed implementations face additional consistency and latency challenges that can negate theoretical performance advantages. Cross-modal systems compound these issues by requiring unified similarity metrics that may inadequately represent heterogeneous data types, leading to systematic retrieval degradation across modalities.

Memory Interference and Coordination. Memory systems face fundamental coordination challenges beyond scalability, particularly in multi-task scenarios. The stability-plasticity dilemma manifests distinctly across architectures: while current literature confirms surprise-gated systems excel at novelty detection, evidence for their specific struggles with gradual knowledge drift remains limited. Policy-learned approaches show strong task adaptation capabilities but face inherent overfitting risks during continual learning.

Memory interference emerges as a critical bottleneck when similar contexts trigger conflicting information retrieval. External memory systems suffer from catastrophic collisions during concurrent access, while parameter-encoded approaches experience gradient interference during continual updates. Hybrid archi-

tures attempt to mitigate these issues through functional partitioning, yet optimal allocation strategies remain highly domain-dependent.

Forgetting policies introduce additional coordination complexities. A-MEM’s Xu et al. (2025) evolutionary approach demonstrates promise for utility-based memory curation but requires careful tuning to avoid inadvertent erasure of rare but valuable information. Systems employing LRU and selective forgetting strategies perform effectively in structured environments but struggle under non-stationary conditions where relevance patterns shift unpredictably.

Evaluation and Standardization Gaps. The field suffers from fragmented evaluation methodologies that prevent systematic cross-architecture comparisons. Current benchmarks exhibit dramatic variation in memory requirements, task complexity, and evaluation metrics, making it difficult to assess fundamental trade-offs between different memory strategies. Critical deficiencies include inconsistent context length protocols, divergent benchmark emphases on retrieval versus reasoning capabilities, and absent robustness evaluations for adversarial scenarios, memory corruption, and distribution shift. Most significantly, existing evaluations rarely assess long-term adaptation, memory utilization efficiency, or interference mitigation—precisely the capabilities that distinguish sophisticated memory systems from basic approaches.

5.3 Future Directions

The evolution of memory-augmented Transformers toward truly cognitive architectures requires coordinated advances across multiple research frontiers. The following directions represent the most promising paths for achieving human-like memory capabilities while addressing current technical limitations.

Toward Cognitive Flexibility and Lifelong Learning. Emerging paradigms in memory-augmented transformers focus on building systems that can dynamically store, retrieve, and update knowledge in ways that reflect the adaptability of biological memory. Neuroscientific insights from Dijksterhuis et al. (2024) highlight the value of memory consolidation, revealing how concept cells in the human hippocampus reactivate when pronouns reference specific nouns, seamlessly linking new linguistic input to stored concepts—a mechanism comparable to integrating episodic memories into an LLM’s parametric memory to bypass capacity constraints and achieve lasting retention. Complementing this, the position paper by Pink et al. (2025) argues that episodic memory is a vital missing component for long-term LLM agents, proposing a framework with five essential properties to foster adaptive behavior and outlining a research roadmap to embed these capabilities. A pivotal trend in this direction is the decoupling of computation from storage, enabling models to tap into external or hybrid memory banks for real-time, current information without the need for retraining. Such architectures facilitate personalized, context-responsive outputs in evolving environments. Furthermore, innovations like test-time training and memory-driven optimization empower models to learn and adapt during deployment, bolstered by selective forgetting and zero-shot transfer mechanisms that enhance generalization. The integration of multimodal memory and collaborative networks also holds promise for deeper reasoning and shared learning among agents. To sustain these advancements, progress in hierarchical storage, memory compression, and hardware-aware design is driving scalable, energy-efficient deployment across varied platforms.

Toward Human-Like Cognition: The Role of Memory in Intelligent Agents. As intelligent agents evolve toward more human-like reasoning and autonomy, the integration of sophisticated memory systems becomes a central design challenge. Unlike conventional LLMs that operate in a stateless fashion, truly interactive agents must preserve context, interpret ongoing events, and adapt their behavior across time (Liang et al., 2024; Yi et al., 2025). Drawing from cognitive science, recent agent architectures incorporate short-term memory for maintaining dialogue context, working memory for in-the-moment reasoning, and long-term memory for accumulating knowledge and past experiences (Li et al., 2024). Vector databases have emerged as a popular solution for implementing long-term memory, enabling fast, similarity-based retrieval of episodic and procedural knowledge (Hatalis et al., 2023). However, realizing robust memory-driven behavior introduces significant difficulties. Agents often fail to separate memory types, leading to conflicts between episodic and semantic recall, and may repeatedly attempt failed subtasks without effective use of episodic feedback (Wang et al., 2024c). As memory grows over time, retrieval speed and storage cost become critical concerns, especially when managing large volumes of data. Static or manually defined metadata can limit

retrieval quality, pointing to a need for agents to learn metadata attributes dynamically to support smarter decision-making (Sarch et al., 2023). Moreover, integrating long-term memory with external knowledge bases like ontologies or knowledge graphs could enhance contextual grounding and reasoning (Wang et al., 2024c). Addressing these issues is essential to building agents capable of flexible, adaptive cognition that mirrors the structure and function of human memory systems.

Future Architectures and Ethical Considerations. Test-time training, memory-driven optimization, and zero-shot transfer learning allow models to adapt during deployment, offering the promise of lifelong learning. Multimodal memory systems and collaborative agent networks open new paths for collective intelligence, deeper reasoning, and shared learning across environments. Yet, these capabilities also introduce ethical and societal considerations. As memory-augmented Transformers are adopted in sensitive domains like healthcare, education, and personalized services, ensuring transparency, privacy, and user control over memory becomes vital. Techniques for explainable memory operations, data auditing, and bias mitigation will be critical to build trust and prevent misuse.

In summary, the future of memory-augmented Transformers lies in bridging engineering efficiency with cognitive flexibility. By combining continual learning, dynamic memory adaptation, and biologically inspired design, alongside ethical safeguards, these systems have the potential to transform AI from static pattern recognizers into adaptive, intelligent agents.

References

- Christoph Anacker and René Hen. Adult hippocampal neurogenesis and cognitive flexibility—linking memory and mood. *Nature Reviews Neuroscience*, 18(6):335–346, 2017.
- Bernard J Baars, Natalie Geld, and Robert Kozma. Global workspace theory (gwt) and prefrontal cortex: Recent developments. *Frontiers in psychology*, 12:749868, 2021.
- Alan Baddeley. Working memory: looking back and looking forward. *Nature reviews neuroscience*, 4(10):829–839, 2003.
- Martin LLR Barry and Wulfram Gerstner. Fast adaptation to rule switching using neuronal surprise. *PLoS computational biology*, 20(2):e1011839, 2024.
- Amjad H Bazzari and H Rheinallt Parri. Neuromodulators and long-term synaptic plasticity in learning and memory: A steered-glutamatergic perspective. *Brain sciences*, 9(11):300, 2019.
- Ali Behrouz, Peilin Zhong, and Vahab Mirrokni. Titans: Learning to memorize at test time. *arXiv preprint arXiv:2501.00663*, 2024.
- Ali Behrouz, Zeman Li, Praneeth Kacham, Majid Daliri, Yuan Deng, Peilin Zhong, Meisam Razaviyayn, and Vahab Mirrokni. Atlas: Learning to optimally memorize the context at test time. *arXiv preprint arXiv:2505.23735*, 2025.
- Iz Beltagy, Matthew E Peters, and Arman Cohan. Longformer: The long-document transformer. *arXiv preprint arXiv:2004.05150*, 2020.
- Vincent-Pierre Berges, Barlas Oğuz, Daniel Haziza, Wen-tau Yih, Luke Zettlemoyer, and Gargi Ghosh. Memory layers at scale. *arXiv preprint arXiv:2412.09764*, 2024.
- Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Millican, George Bm Van Den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, et al. Improving language models by retrieving from trillions of tokens. In *International conference on machine learning*, pp. 2206–2240. PMLR, 2022.
- Eduardo Camina and Francisco Güell. The neuroanatomical, neurophysiological and psychological basis of memory: Current models and their origins. *Frontiers in pharmacology*, 8:438, 2017.
- Edoardo Cetin, Qi Sun, Tianyu Zhao, and Yujin Tang. An evolved universal transformer memory. In *The Thirteenth International Conference on Learning Representations*, 2025.

-
- Howard Chen, Ramakanth Pasunuru, Jason Weston, and Asli Celikyilmaz. Walking down the memory maze: Beyond context limit through interactive reading. *arXiv preprint arXiv:2310.05029*, 2023.
- Prateek Chhikara, Dev Khant, Saket Aryan, Taranjeet Singh, and Deshraj Yadav. Mem0: Building production-ready ai agents with scalable long-term memory. *arXiv preprint arXiv:2504.19413*, 2025.
- Marvin M Chun and Nicholas B Turk-Browne. Interactions between attention and memory. *Current opinion in neurobiology*, 17(2):177–184, 2007.
- Nicholas G Cicero, Elizabeth Riley, Khena M Swallow, Eve De Rosa, and Adam Anderson. Attention-dependent coupling with forebrain and brainstem neuromodulatory nuclei differs across the lifespan. *Geroscience*, pp. 1–20, 2025.
- Nelson Cowan. What are the differences between long-term, short-term, and working memory? *Progress in brain research*, 169:323–338, 2008.
- Zihang Dai, Zhilin Yang, Yiming Yang, Jaime Carbonell, Quoc V Le, and Ruslan Salakhutdinov. Transformer-xl: Attentive language models beyond a fixed-length context. *arXiv preprint arXiv:1901.02860*, 2019.
- Stanislas Dehaene, Jean-Pierre Changeux, and Lionel Naccache. The global neuronal workspace model of conscious access: from neuronal architectures to clinical applications. *Characterizing consciousness: From cognition to the clinic?*, pp. 55–84, 2011.
- Doris E Dijksterhuis, Matthew W Self, Jessy K Possel, Judith C Peters, ECW van Straaten, Sander Idema, Johannes C Baaijen, Sandra MA van der Salm, Erik J Aarnoutse, Nicole CE van Klink, et al. Pronouns reactivate conceptual representations in human hippocampal neurons. *Science*, 385(6716):1478–1484, 2024.
- Florin Dolcos, Kevin S LaBar, and Roberto Cabeza. Interaction between the amygdala and the medial temporal lobe memory system predicts better memory for emotional events. *Neuron*, 42(5):855–863, 2004.
- Yiming Du, Wenyu Huang, Danna Zheng, Zhaowei Wang, Sebastien Montella, Mirella Lapata, Kam-Fai Wong, and Jeff Z Pan. Rethinking memory in ai: Taxonomy, operations, topics, and future directions. *arXiv preprint arXiv:2505.00675*, 2025.
- Sanchayan Dutta and Suvrit Sra. Memory-augmented transformers can implement linear first-order optimization methods. *arXiv preprint arXiv:2410.07263*, 2024.
- Marco Federici, Davide Belli, Mart Van Baalen, Amir Jalalirad, Andrii Skliar, Bence Major, Markus Nagel, and Paul Whatmough. Efficient llm inference using dynamic input pruning and cache-aware masking. *arXiv preprint arXiv:2412.01380*, 2024.
- Zafeirios Fountas, Martin A Benfeghoul, Adnan Oomerjee, Fenia Christopoulou, Gerasimos Lampouras, Haitham Bou-Ammar, and Jun Wang. Human-like episodic memory for infinite context llms. *arXiv preprint arXiv:2407.09450*, 2024.
- Darya Frank, Alex Kafkas, and Daniela Montaldi. Experiencing surprise: The temporal dynamics of its impact on memory. *Journal of Neuroscience*, 42(33):6435–6444, 2022.
- Pengyu Gao, Jinming Zhao, Xinyue Chen, and Long Yilin. An efficient context-dependent memory framework for llm-centric agents. In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 3: Industry Track)*, pp. 1055–1069, 2025.
- Ruopeng Gao and Limin Wang. Memotr: Long-term memory-augmented transformer for multi-object tracking. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 9901–9910, 2023.
- Adam Gazzaley and Anna C Nobre. Top-down modulation: bridging selective attention and working memory. *Trends in cognitive sciences*, 16(2):129–135, 2012.

-
- Saibo Geng, Nathan Ranchin, Maxime Peyrard, Chris Wendler, Michael Gastpar, Robert West, et al. zip2zip: Inference-time adaptive vocabularies for language models via token compression. *arXiv preprint arXiv:2506.01084*, 2025.
- Paul E Gilbert and Andrea M Brushfield. The role of the ca3 hippocampal subregion in spatial memory: a process oriented behavioral assessment. *Progress in Neuro-Psychopharmacology and Biological Psychiatry*, 33(5):774–781, 2009.
- Alex Graves, Greg Wayne, and Ivo Danihelka. Neural turing machines. *arXiv preprint arXiv:1410.5401*, 2014.
- Alex Graves, Greg Wayne, Malcolm Reynolds, Tim Harley, Ivo Danihelka, Agnieszka Grabska-Barwińska, Sergio Gómez Colmenarejo, Edward Grefenstette, Tiago Ramalho, John Agapiou, et al. Hybrid computing using a neural network with dynamic external memory. *Nature*, 538(7626):471–476, 2016.
- Bernal Jiménez Gutiérrez, Yiheng Shu, Yu Gu, Michihiro Yasunaga, and Yu Su. Hipporag: Neurobiologically inspired long-term memory for large language models. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- Uri Hasson, Janice Chen, and Christopher J Honey. Hierarchical process memory: memory as an integral component of information processing. *Trends in cognitive sciences*, 19(6):304–313, 2015.
- Kostas Hatalis, Despina Christou, Joshua Myers, Steven Jones, Keith Lambert, Adam Amos-Binks, Zohreh Dannenhauer, and Dustin Dannenhauer. Memory matters: The need to improve long-term memory in llm-agents. In *Proceedings of the AAAI Symposium Series*, volume 2, pp. 277–280, 2023.
- Zifan He, Yingqi Cao, Zongyue Qin, Neha Prakriya, Yizhou Sun, and Jason Cong. Hmt: Hierarchical memory transformer for efficient long context language processing. *arXiv preprint arXiv:2405.06067*, 2024a.
- Zihong He, Weizhe Lin, Hao Zheng, Fan Zhang, Matt W Jones, Laurence Aitchison, Xuhai Xu, Miao Liu, Per Ola Kristensson, and Junxiao Shen. Human-inspired perspectives: A survey on ai long-term memory. *arXiv preprint arXiv:2411.00489*, 2024b.
- Cameron Higgins, Yunzhe Liu, Diego Vidaurre, Zeb Kurth-Nelson, Ray Dolan, Timothy Behrens, and Mark Woolrich. Replay bursts in humans coincide with activation of the default mode and parietal alpha networks. *Neuron*, 109(5):882–893, 2021.
- Dongseong Hwang, Weiran Wang, Zhuoyuan Huo, Khe Chai Sim, and Pedro Moreno Mengibar. Transformerfam: Feedback attention is working memory. *arXiv preprint arXiv:2404.09173*, 2024.
- Jiazheng Kang, Mingming Ji, Zhe Zhao, and Ting Bai. Memory os of ai agent. *arXiv preprint arXiv:2506.06326*, 2025a.
- Jikun Kang, Wenqi Wu, Filippas Christianos, Alex J Chan, Fraser Greenlee, George Thomas, Marvin Purtorab, and Andy Toulis. Lm2: Large memory models. *arXiv preprint arXiv:2502.06049*, 2025b.
- Louis Kang and Taro Toyoizumi. Distinguishing examples while building concepts in hippocampal and artificial networks. *Nature Communications*, 15(1):647, 2024.
- Jens G Klinzing, Niels Niethard, and Jan Born. Mechanisms of systems memory consolidation during sleep. *Nature neuroscience*, 22(10):1598–1610, 2019.
- Ching-Yun Ko, Sihui Dai, Payel Das, Georgios Kollias, Subhajit Chaudhury, and Aurelie Lozano. Memreasoner: A memory-augmented llm architecture for multi-hop reasoning. In *The First Workshop on System-2 Reasoning at Scale, NeurIPS’24*, 2024.
- Wouter Kruijne, Sander M Bohte, Pieter R Roelfsema, and Christian NL Olivers. Flexible working memory through selective gating and attentional tagging. *Neural Computation*, 33(1):1–40, 2021.

-
- Abhilasha A Kumar. Semantic memory: A review of methods, models, and current challenges. *Psychonomic bulletin & review*, 28(1):40–80, 2021.
- Hung Le, Truyen Tran, and Svetha Venkatesh. Self-attentive associative memory. In *International conference on machine learning*, pp. 5682–5691. PMLR, 2020.
- Hung Le, Dai Do, Dung Nguyen, and Svetha Venkatesh. Reasoning under 1 billion: Memory-augmented reinforcement learning for large language models. *arXiv preprint arXiv:2504.02273*, 2025.
- Dongkyu Lee, Chandana Satya Prakash, Jack FitzGerald, and Jens Lehmann. Matter: Memory-augmented transformer using heterogeneous knowledge sources. *arXiv preprint arXiv:2406.04670*, 2024.
- Jonathan LC Lee, Karim Nader, and Daniela Schiller. An update on memory reconsolidation updating. *Trends in cognitive sciences*, 21(7):531–545, 2017.
- Patrick Lewis, Ethan Perez, Aleksandra 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.
- Zeyuan Li, Yangfan He, Lewei He, Jianhui Wang, Tianyu Shi, Bin Lei, Yuchen Li, and Qiuwu Chen. Falcon: Feedback-driven adaptive long/short-term memory reinforced coding optimization system. *arXiv preprint arXiv:2410.21349*, 2024.
- Xuechen Liang, Yangfan He, Yinghui Xia, Xinyuan Song, Jianhui Wang, Meiling Tao, Li Sun, Xinhang Yuan, Jiayi Su, Keqin Li, et al. Self-evolving agents with reflective and memory-augmented abilities. *arXiv preprint arXiv:2409.00872*, 2024.
- Bang Liu, Xinfeng Li, Jiayi Zhang, Jinlin Wang, Tanjin He, Sirui Hong, Hongzhang Liu, Shaokun Zhang, Kaitao Song, Kunlun Zhu, et al. Advances and challenges in foundation agents: From brain-inspired intelligence to evolutionary, collaborative, and safe systems. *arXiv preprint arXiv:2504.01990*, 2025.
- Lei Liu, Xiaoyan Yang, Yue Shen, Binbin Hu, Zhiqiang Zhang, Jinjie Gu, and Guannan Zhang. Think-in-memory: Recalling and post-thinking enable llms with long-term memory. *arXiv preprint arXiv:2311.08719*, 2023.
- Weijie Liu, Zecheng Tang, Juntao Li, Kehai Chen, and Min Zhang. Memlong: Memory-augmented retrieval for long text modeling. *arXiv preprint arXiv:2408.16967*, 2024.
- Yunzhe Liu, Marcelo G Mattar, Timothy EJ Behrens, Nathaniel D Daw, and Raymond J Dolan. Experience replay is associated with efficient nonlocal learning. *Science*, 372(6544):eabf1357, 2021.
- Wenhan Luo, Di Yun, Yi Hu, Miaomiao Tian, Jiajun Yang, Yifan Xu, Yong Tang, Yang Zhan, Hong Xie, and Ji-Song Guan. Acquiring new memories in neocortex of hippocampal-lesioned mice. *Nature communications*, 13(1):1601, 2022.
- Guixiang Ma, Vy A Vo, Theodore L Willke, and Nesreen K Ahmed. Memory-augmented graph neural networks: A brain-inspired review. *IEEE Transactions on Artificial Intelligence*, 5(5):2011–2025, 2023.
- Martial Mermillod, Aurélie Bugaiska, and Patrick Bonin. The stability-plasticity dilemma: Investigating the continuum from catastrophic forgetting to age-limited learning effects, 2013.
- George A Miller. The magical number seven, plus or minus two: Some limits on our capacity for processing information. *Psychological review*, 63(2):81, 1956.
- Ali Modarressi, Abdullatif Köksal, Ayyoob Imani, Mohsen Fayyaz, and Hinrich Schütze. Memllm: Finetuning llms to use an explicit read-write memory. *arXiv preprint arXiv:2404.11672*, 2024.
- Morris Moscovitch, Roberto Cabeza, Gordon Winocur, and Lynn Nadel. Episodic memory and beyond: the hippocampus and neocortex in transformation. *Annual review of psychology*, 67(1):105–134, 2016.

-
- Swaleha Mujawar, Jaideep Patil, Bhushan Chaudhari, and Daniel Saldanha. Memory: Neurobiological mechanisms and assessment. *Industrial psychiatry journal*, 30(Suppl 1):S311–S314, 2021.
- Thomas Nail. Most brain activity is ‘background noise’—and that’s upending our understanding of consciousness, 2021.
- Piotr Nawrot, Adrian Łańcucki, Marcin Chochowski, David Tarjan, and Edoardo M Ponti. Dynamic memory compression: Retrofitting llms for accelerated inference. *arXiv preprint arXiv:2403.09636*, 2024.
- Erika Nyhus and Tim Curran. Functional role of gamma and theta oscillations in episodic memory. *Neuroscience & Biobehavioral Reviews*, 34(7):1023–1035, 2010.
- Charles Packer, Vivian Fang, Shishir_G Patil, Kevin Lin, Sarah Wooders, and Joseph_E Gonzalez. Memgpt: Towards llms as operating systems. 2023.
- Hao Peng, Jungo Kasai, Nikolaos Pappas, Dani Yogatama, Zhaofeng Wu, Lingpeng Kong, Roy Schwartz, and Noah A Smith. Abc: Attention with bounded-memory control. *arXiv preprint arXiv:2110.02488*, 2021.
- Mathis Pink, Qinyuan Wu, Vy Ai Vo, Javier Turek, Jianing Mu, Alexander Huth, and Mariya Toneva. Position: Episodic memory is the missing piece for long-term llm agents. *arXiv preprint arXiv:2502.06975*, 2025.
- Ofir Press, Noah A Smith, and Mike Lewis. Train short, test long: Attention with linear biases enables input length extrapolation. *arXiv preprint arXiv:2108.12409*, 2021.
- Luke Y Prince, Travis J Bacon, Cezar M Tigaret, and Jack R Mellor. Neuromodulation of the feedforward dentate gyrus-ca3 microcircuit. *Frontiers in synaptic neuroscience*, 8:32, 2016.
- Zihan Qiu, Zeyu Huang, Shuang Cheng, Yizhi Zhou, Zili Wang, Ivan Titov, and Jie Fu. Layerwise recurrent router for mixture-of-experts. *arXiv preprint arXiv:2408.06793*, 2024.
- Michel Quak, Raquel Elea London, and Durk Talsma. A multisensory perspective of working memory. *Frontiers in human neuroscience*, 9:197, 2015.
- Jack W Rae, Anna Potapenko, Siddhant M Jayakumar, and Timothy P Lillicrap. Compressive transformers for long-range sequence modelling. *arXiv preprint arXiv:1911.05507*, 2019.
- Marcus E Raichle, Ann Mary MacLeod, Abraham Z Snyder, William J Powers, Debra A Gusnard, and Gordon L Shulman. A default mode of brain function. *Proceedings of the national academy of sciences*, 98(2):676–682, 2001.
- J Ranjith and Santhi Baskaran. Adaptive knowledge consolidation: A dynamic approach to mitigating catastrophic forgetting in text-based neural networks. 2024.
- Ryan V Raut, Abraham Z Snyder, and Marcus E Raichle. Hierarchical dynamics as a macroscopic organizing principle of the human brain. *Proceedings of the National Academy of Sciences*, 117(34):20890–20897, 2020.
- Daniel Reznik, Robert Trampel, Nikolaus Weiskopf, Menno P Witter, and Christian F Doeller. Dissociating distinct cortical networks associated with subregions of the human medial temporal lobe using precision neuroimaging. *Neuron*, 111(17):2756–2772, 2023.
- Ivan Rodkin, Yuri Kuratov, Aydar Bulatov, and Mikhail Burtsev. Associative recurrent memory transformer. *arXiv preprint arXiv:2407.04841*, 2024.
- Edmund T Rolls. The mechanisms for pattern completion and pattern separation in the hippocampus. *Frontiers in systems neuroscience*, 7:74, 2013.
- Frédéric Roux and Peter J Uhlhaas. Working memory and neural oscillations: alpha–gamma versus theta–gamma codes for distinct wm information? *Trends in cognitive sciences*, 18(1):16–25, 2014.

-
- Jacob Russin, Randall C O'Reilly, and Yoshua Bengio. Deep learning needs a prefrontal cortex. *Work Bridging AI Cogn Sci*, 107(603-616):1, 2020.
- Gabriel Sarch, Yue Wu, Michael J Tarr, and Katerina Fragkiadaki. Open-ended instructable embodied agents with memory-augmented large language models. *arXiv preprint arXiv:2310.15127*, 2023.
- Thomas Schmied, Fabian Paischer, Vihang Patil, Markus Hofmarcher, Razvan Pascanu, and Sepp Hochreiter. Retrieval-augmented decision transformer: External memory for in-context rl. *arXiv preprint arXiv:2410.07071*, 2024.
- Lars Schwabe and Oliver T Wolf. Stress and multiple memory systems: from ‘thinking’ to ‘doing’. *Trends in cognitive sciences*, 17(2):60–68, 2013.
- Lianlei Shan, Shixian Luo, Zezhou Zhu, Yu Yuan, and Yong Wu. Cognitive memory in large language models. *arXiv preprint arXiv:2504.02441*, 2025.
- Liang Shi, Chuqi Liu, Xiaojing Peng, Yifei Cao, Daniel A Levy, and Gui Xue. The neural representations underlying asymmetric cross-modal prediction of words. *Human Brain Mapping*, 44(6):2418–2435, 2023.
- Alyssa H Sinclair, Grace M Manalili, Iva K Brunec, R Alison Adcock, and Morgan D Barense. Prediction errors disrupt hippocampal representations and update episodic memories. *Proceedings of the National Academy of Sciences*, 118(51):e2117625118, 2021.
- Larry R Squire, Lisa Genzel, John T Wixted, and Richard G Morris. Memory consolidation. *Cold Spring Harbor perspectives in biology*, 7(8):a021766, 2015.
- Bernhard P Staresina, Sebastian Michelmann, Mathilde Bonnefond, Ole Jensen, Nikolai Axmacher, and Juergen Fell. Hippocampal pattern completion is linked to gamma power increases and alpha power decreases during recollection. *elife*, 5:e17397, 2016.
- Qi Sun, Edoardo Cetin, and Yujin Tang. Transformer-squared: Self-adaptive llms. In *The Thirteenth International Conference on Learning Representations*, 2025.
- Yuwei Sun, Hideya Ochiai, Zhirong Wu, Stephen Lin, and Ryota Kanai. Associative transformer. *arXiv preprint arXiv:2309.12862*, 2023.
- Lucas CS Tavares and Adriano BL Tort. Hippocampal–prefrontal interactions during spatial decision-making. *Hippocampus*, 32(1):38–54, 2022.
- Yi Tay, Vinh Tran, Mostafa Dehghani, Jianmo Ni, Dara Bahri, Harsh Mehta, Zhen Qin, Kai Hui, Zhe Zhao, Jai Gupta, et al. Transformer memory as a differentiable search index. *Advances in Neural Information Processing Systems*, 35:21831–21843, 2022.
- Timothy J Teyler and Jerry W Rudy. The hippocampal indexing theory and episodic memory: updating the index. *Hippocampus*, 17(12):1158–1169, 2007.
- Ashish Vaswani. Attention is all you need. *Advances in neural information processing systems*, 30:I, 2017.
- Guan Wang, Jin Li, Yuhao Sun, Xing Chen, Changling Liu, Yue Wu, Meng Lu, Sen Song, and Yasin Abbasi Yadkori. Hierarchical reasoning model. *arXiv preprint arXiv:2506.21734*, 2025a.
- Wei Wang and Qing Li. Schrodinger’s memory: Large language models. *arXiv preprint arXiv:2409.10482*, 2024.
- Weizhi Wang, Li Dong, Hao Cheng, Xiaodong Liu, Xifeng Yan, Jianfeng Gao, and Furu Wei. Augmenting language models with long-term memory. *Advances in Neural Information Processing Systems*, 36:74530–74543, 2023.
- Xiaoqiang Wang, Suyuchen Wang, Yun Zhu, and Bang Liu. R3mem: Bridging memory retention and retrieval via reversible compression. *arXiv preprint arXiv:2502.15957*, 2025b.

-
- Xindi Wang, Mahsa Salmani, Parsa Omid, Xiangyu Ren, Mehdi Rezagholizadeh, and Armaghan Eshaghi. Beyond the limits: A survey of techniques to extend the context length in large language models. *arXiv preprint arXiv:2402.02244*, 2024a.
- Yu Wang, Yifan Gao, Xiushi Chen, Haoming Jiang, Shiyang Li, Jingfeng Yang, Qingyu Yin, Zheng Li, Xian Li, Bing Yin, et al. Memoryllm: Towards self-updatable large language models. *arXiv preprint arXiv:2402.04624*, 2024b.
- Yu Wang, Dmitry Krotov, Yuanzhe Hu, Yifan Gao, Wangchunshu Zhou, Julian McAuley, Dan Gutfreund, Rogerio Feris, and Zexue He. M+: Extending memoryllm with scalable long-term memory. *arXiv preprint arXiv:2502.00592*, 2025c.
- Zihao Wang, Shaofei Cai, Anji Liu, Yonggang Jin, Jinbing Hou, Bowei Zhang, Haowei Lin, Zhaofeng He, Zilong Zheng, Yaodong Yang, et al. Jarvis-1: Open-world multi-task agents with memory-augmented multimodal language models. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024c.
- Qingyang Wu and Zhou Yu. Stateful memory-augmented transformers for efficient dialogue modeling. *arXiv preprint arXiv:2209.07634*, 2022.
- Qingyang Wu, Zhenzhong Lan, Kun Qian, Jing Gu, Alborz Geramifard, and Zhou Yu. Memformer: A memory-augmented transformer for sequence modeling. *arXiv preprint arXiv:2010.06891*, 2020.
- Yan Wu, Greg Wayne, Alex Graves, and Timothy Lillicrap. The kanerva machine: A generative distributed memory. *arXiv preprint arXiv:1804.01756*, 2018.
- Yaxiong Wu, Sheng Liang, Chen Zhang, Yichao Wang, Yongyue Zhang, Huifeng Guo, Ruiming Tang, and Yong Liu. From human memory to ai memory: A survey on memory mechanisms in the era of llms. *arXiv preprint arXiv:2504.15965*, 2025.
- Yuhuai Wu, Markus N Rabe, DeLesley Hutchins, and Christian Szegedy. Memorizing transformers. *arXiv preprint arXiv:2203.08913*, 2022a.
- Yuxiang Wu, Yu Zhao, Baotian Hu, Pasquale Minervini, Pontus Stenetorp, and Sebastian Riedel. An efficient memory-augmented transformer for knowledge-intensive nlp tasks. *arXiv preprint arXiv:2210.16773*, 2022b.
- 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.
- Fuzhao Xue, Valerii Likhoshesterov, Anurag Arnab, Neil Houlsby, Mostafa Dehghani, and Yang You. Adaptive computation with elastic input sequence. In *International Conference on Machine Learning*, pp. 38971–38988. PMLR, 2023.
- Hongkang Yang, Zehao Lin, Wenjin Wang, Hao Wu, Zhiyu Li, Bo Tang, Wenqiang Wei, Jinbo Wang, Zeyun Tang, Shichao Song, et al. Memory3: Language modeling with explicit memory. *arXiv preprint arXiv:2407.01178*, 2024.
- Wang Yang, Zirui Liu, Hongye Jin, Qingyu Yin, Vipin Chaudhary, and Xiaotian Han. Longer context, deeper thinking: Uncovering the role of long-context ability in reasoning. *arXiv preprint arXiv:2505.17315*, 2025.
- Qiang Yi, Yangfan He, Jianhui Wang, Xinyuan Song, Shiyao Qian, Xinhang Yuan, Miao Zhang, Li Sun, Keqin Li, Kuan Lu, et al. Score: Story coherence and retrieval enhancement for ai narratives. *arXiv preprint arXiv:2503.23512*, 2025.
- Sihyun Yu, Meera Hahn, Dan Kondratyuk, Jinwoo Shin, Agrim Gupta, José Lezama, Irfan Essa, David Ross, and Jonathan Huang. Malt diffusion: Memory-augmented latent transformers for any-length video generation. *arXiv preprint arXiv:2502.12632*, 2025.

-
- Songlin Zhai, Yuan Meng, Yongrui Chen, Yiwei Wang, and Guilin Qi. Peripheral memory for llms: Integration of sequential memory banks with adaptive querying. In *Forty-second International Conference on Machine Learning*, 2025.
- Zeyu Zhang, Xiaohe Bo, Chen Ma, Rui Li, Xu Chen, Quanyu Dai, Jieming Zhu, Zhenhua Dong, and Jirong Wen. A survey on the memory mechanism of large language model based agents. *arXiv preprint arXiv:2404.13501*, 2024.
- Wanjun Zhong, Lianghong Guo, Qiqi Gao, He Ye, and Yanlin Wang. Memorybank: Enhancing large language models with long-term memory. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 19724–19731, 2024a.
- Weishun Zhong, Tankut Can, Antonis Georgiou, Ilya Shnayderman, Mikhail Katkov, and Misha Tsodyks. Random tree model of meaningful memory. *bioRxiv*, pp. 2024–12, 2024b.
- Zayd Muhammad Kawakibi Zuhri, Muhammad Farid Adilazuarda, Ayu Purwarianti, and Alham Fikri Aji. Mlkv: Multi-layer key-value heads for memory efficient transformer decoding. *arXiv preprint arXiv:2406.09297*, 2024.