

# LLM Agent Memory: A Survey from a Unified Representation–Management Perspective

Anonymous ACL submission

## Abstract

Large language models (LLMs) face significant challenges in sustaining long-term memory for agentic applications due to limited context windows. To address this limitation, many work has proposed diverse memory mechanisms to support long-term, multi-turn interactions, leveraging different approaches tailored to distinct memory storage objects, such as KV caches. In this survey, we present a unified taxonomy that organizes memory systems for long-context scenarios by decoupling memory abstractions from model-specific inference and training methods. We categorize LLM memory into three primary paradigms: natural language tokens, intermediate representations and parameters. For each paradigm, we organize existing methods by three management stages, including memory construction, update, and query, so that long-context memory mechanisms can be described in a consistent way across system designs, with their implementation choices and constraints made explicit. Finally, we outline key research directions for long-context memory system design.

## 1 Introduction

Large Language Models (LLMs) are deployed in long-context and interactive settings, such as multi-turn dialogue, task-oriented assistants, and agent-based systems, where task completion requires information to be retained and reused across extended temporal spans rather than within a single prompt (Qian et al., 2023; Gao et al., 2023; Wang et al., 2023e). As input sequences grow to thousands or even millions of tokens, performance often degrades on tasks that require entity tracking, logical consistency, or recall of task-relevant facts across long interaction histories (Gao et al., 2024; Zhang et al., 2024c; Wu et al., 2025b). These observations indicate that the core difficulty in long-context and multi-turn scenarios lies not only in *how much information the model can observe*, but in *how in-*

*formation is selectively retained, retrieved, and integrated during inference* (Shinwari and Usama, 2025; Maharana et al., 2024; Wan and Ma, 2025). This has motivated the introduction of *memory* as a key abstraction in LLM-based systems, enabling task-relevant information to persist beyond the immediate context window and to be reused for future reasoning and decision-making.

Many existing studies have explored diverse LLM memory mechanisms to support long-context reasoning and multi-turn interaction (Shinn et al., 2023; Zhong et al., 2023; Modarressi et al., 2023; Zhong et al., 2024b; Qian et al., 2024). Broadly, memory enables task-relevant information to persist beyond the immediate context window and be reused across interactions, often drawing loose analogies to human cognition, where short-term working memory supports immediate reasoning and long-term memory enables recall over time (Baddeley, 2007; Budson and Kensinger, 2023). The analogy serves as a useful intuition that effective memory depends not only on capacity, but also on selective access and integration (Gao et al., 2024; Wu et al., 2025b). Corresponding surveys have reviewed this literature from multiple perspectives, including LLM-based agents and long-term interaction (Zhang et al., 2025b), long-context modeling and long-term memory (Huang et al., 2023b; Jiang et al., 2024a; Wu et al., 2025b), personalization (Liu et al., 2025a), and system-level efficiency such as inference-time memory management (Pan and Li, 2025; LI et al., 2025; Luohe et al.). These works underscore the central role of memory across application, modeling, and system dimensions.

In this survey, we take a *system-level, operation-centric* view of LLM memory. Instead of cataloging methods by tasks or modules, we organize the literature around a unified abstraction that connects **task knowledge requirements** to **memory representations** through shared **management interfaces**. Figure 1 sketches this logic from ap-

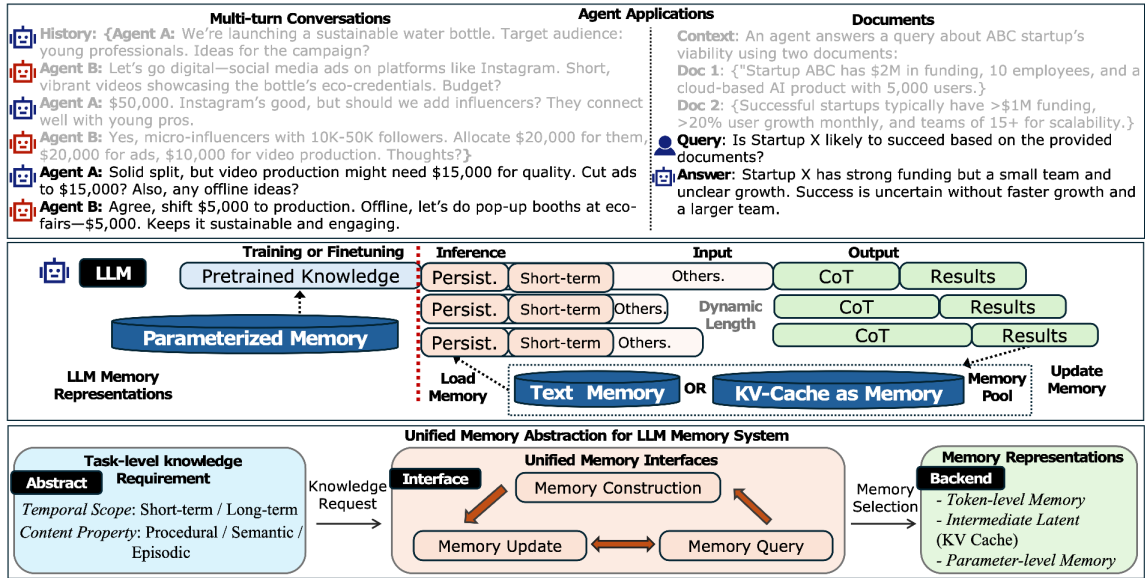


Figure 1: Overview of LLM memory applications, representation, and unified management abstraction.

lications (top) to recurring representation choices (middle) and a common management interface over long interactions (bottom). Under this view, diverse methods can be compared by three questions: *what is stored, where it is stored, and how it is operated*. We structure the survey along two recurring dimensions:

- **Memory representation** describes *where and in what form* information resides: *token-level memory* in the input context, *intermediate latent memory* as inference-time states (e.g., Key–Value caches), and *parameter-level memory* in model weights via adaptation or editing.
- **Memory management** describes *how memory is operated over time* to satisfy task requirements under practical constraints. Across representations, we observe a shared interface of three core operations: *memory construction* (what to store and how to structure it), *memory update* (how to maintain, consolidate, or remove stored content), and *memory query* (how to select and integrate relevant information during inference).

Following this organization, Sections 2–4 review memory mechanisms by aligning each representation with the same construction–update–query interface and its key trade-offs. Section 5 synthesizes insights across memory types and outlines open challenges for reliable long-context and multi-turn LLM memory.

## 2 Natural Language Tokens as Memory

Natural language tokens are the most *explicit* memory format that leverages the context window for non-invasive information reuse. However, it is limited by  $\mathcal{O}(n^2)$  computational costs and the tendency to overlook information buried in the middle of long sequences (Liu et al., 2024c). To address these bottlenecks, **Retrieval-Augmented Generation** and **Agentic Memory** transform fragmented histories into structured, retrievable knowledge.

*Goal:* Keep the *right* evidence in context for reliable reasoning.  
*Challenge:* *Relevance under a context budget.* Token-level memory should be treated as an *auditable working set*. Construction and update must compress and organize information so query returns a small, high-signal set rather than long history.

### 2.1 Retrieval-Augmented Generation (RAG)

RAG here refers to retrieval-based memory systems centered around an external a vector database, rather than full agent frameworks with explicit environment interaction. RAG offers a cost-effective way to extend long-context behavior by retrieving evidence and inserting it into the input.

**Memory Construction.** Construction in RAG specifies *what evidence is stored and how it is indexed for retrieval*. For *what to store*, early RAG systems mainly ingest *unstructured* corpora and

domain datasets (Li et al., 2023c; Yan et al., 2024). More recent settings extend to *semi-structured* documents such as PDFs, where table-aware processing is often required (e.g., table-to-text normalization or Text-to-SQL-style querying (Zha et al., 2023; Luo et al., 2023)). RAG can also build on *structured* sources such as knowledge graphs, where KnowledGPT (Wang et al., 2023j) and G-Retriever (He et al., 2024a) improve graph grounding and evidence selection via soft prompting and PCST-style subgraph optimization. Some approaches leverage the *model’s internal knowledge* to reduce retrieval overhead or bootstrap contexts, for example by selectively invoking retrieval (Wang et al., 2023k), generating aligned contexts (Yu et al., 2022b), or building unbounded memory pools (Cheng et al., 2023).

For *how to index*, most RAG systems first define the index unit, which determines the *retrieval granularity*: coarser units provide broader context but introduce redundancy, while finer units (e.g., tokens or phrases) improve precision at the risk of losing supporting context (Shi et al., 2023; Yu et al., 2023; Chen et al., 2023b; Jin et al., 2023; Wang et al., 2023g). A common approach splits documents into fixed-length spans (typically 100-512 tokens) and builds an index over these units (Teja, 2023). Variants such as recursive splitting, sliding windows, and Small2Big expansion preserve local coherence by attaching surrounding context (Langchain, 2023; Yang, 2023). Beyond unit definition, many systems enrich indexed entries with auxiliary fields (e.g., page numbers or generated cues) to improve matching (Gao et al., 2022). Indexes can further adopt *hierarchical* or *graph-based* structures to support multi-step evidence aggregation across documents (Wang et al., 2023i). Finally, *embeddings* govern retrieval, encoding units as sparse or dense representations and ranking candidates via lexical matching (e.g., BM25) or embedding similarity, often using dense retrievers such as AngIE, Voyage, or BGE (Li and Li, 2023; VoyageAI, 2023; BAAI, 2023). Hybrid sparse-dense retrieval is commonly used to improve robustness, especially for zero-shot queries or rare entities (Zhang et al., 2025c).

**Memory Update.** RAG updates are primarily realized by modifying the external store and its index, including adding new documents, re-chunking and re-embedding content, refreshing metadata, and restructuring the index (e.g., moving from

flat chunking to hierarchical or KG-based organization) to better match evolving domains and query patterns (Wang et al., 2023i; Gao et al., 2022; Langchain, 2023; Yang, 2023). Some methods also adjust *retrieval invocation*, for example, deciding when retrieval is necessary or generating retrieval-aligned contexts to stabilize generation (Wang et al., 2023k; Yu et al., 2022b; Cheng et al., 2023).

**Memory Query.** To improve robustness under underspecified or noisy inputs, existing methods enhance query quality from several complementary angles. One line focuses on *query reformulation*, including expansion (multi-query), decomposition (e.g., least-to-most prompting), validation (CoVe), and transformations (Zhou et al., 2023a; Dhuliawala et al., 2023; Ma et al., 2023; Peng et al., 2023; Gao et al., 2022), where step-back prompting (Zheng et al., 2024a) further abstracts the query to retrieve complementary evidence. Another line emphasizes *query routing*, selecting different retrieval pipelines based on metadata or semantic routers to enable hybrid strategies (Wang et al., 2025a). A third line develops *multi-step and controllable querying*, where retrieval and generation are composed into modular pipelines (e.g., RRR, GenRead, RECITE (Ma et al., 2023; Yu et al., 2022b; Sun et al., 2022)) and governed by dynamic controllers such as DSP, FLARE, and Self-RAG, sometimes combined with fine-tuning or reinforcement learning (Khatab et al., 2022; Jiang et al., 2023c; Asai et al., 2023; Ke et al., 2024; Lin et al., 2023b).

## 2.2 Agentic Memory

In contrast to RAG, which retrieves evidence from an external corpus, agentic memory targets *stateful, multi-turn interaction* and accumulates textual records of observations, actions, and outcomes across steps to support long-horizon reasoning.

**Memory Construction.** Memory construction converts raw interaction traces into compact, retrievable textual units. A common baseline is to summarize dialogue histories, key events, and stable facts (e.g., user preferences or task states), as in MemoryBank and RET-LLM (Zhong et al., 2024b; Modarressi et al., 2023). To enable efficient access, constructed memories are organized using explicit structures, including key-value slots (Modarressi et al., 2023; Salama et al., 2025; Xi et al., 2024), semantic vector representations (Zhong et al., 2024b; Pan et al., 2025; mem0ai, 2024), and relation-

236 aware graphs capturing dependencies among mem- 287  
237 ory fragments, e.g., CGSN, GraphReader, Hip- 288  
238 poRAG (Nie et al., 2022; Li et al., 2024c; Guti- 289  
239 rez et al., 2024). Construction is strengthened 290  
240 by auxiliary signals such as timestamps, sum- 291  
241 maries, or factual tags to improve retrievability, 292  
242 e.g., LongMemEval (Wu et al., 2024a), or by orga- 293  
243 nizing memories along temporal and causal axes 294  
244 for context-sensitive navigation (Theanine (iunn 295  
245 Ong et al., 2025)). To control storage and infer- 296  
246 ence cost, some systems apply token-level pruning, 297  
247 summarization, or soft compression at construc- 298  
248 tion time, and reuse frequent contexts via prompt 299  
249 caching (Jiang et al., 2023a; Chevalier et al., 2023; 300  
250 Liu et al., 2023a; Gim et al., 2024).

251 **Memory Update.** Memory update refines, con- 302  
252 solidates, and revises stored content as interactions 303  
253 proceed. Early work controls growth through peri- 304  
254 odic summarization or restructuring, using explicit 305  
255 summarizers such as MemoryBank and ChatGPT- 306  
256 RSum (Zhong et al., 2024b; Wang et al., 2025d) 307  
257 or prompt-based extraction of salient topics as in 308  
258 MemoChat (Lu et al., 2023a). Beyond compres- 309  
259 sion, many methods treat updating as a *reasoning-* 310  
260 *driven* step: agents reflect on past actions and out- 311  
261 comes and write back reusable artifacts, includ- 312  
262 ing action-thought traces (Yao et al., 2022), self- 313  
263 critique and revision notes (Shinn et al., 2024), 314  
264 distilled reasoning templates (Yang et al., 2024e), 315  
265 and workflow-level records (Wang et al., 2024m). 316  
266 Experience-based agents refine memory through 317  
267 trial-and-error interaction and feedback, revising 318  
268 what to store and how to use it (Liu et al., 2023b; 319  
269 Zhu et al., 2023b; Wang et al., 2023b; Yao et al., 320  
270 2023; Zhao et al., 2024a; Li et al., 2024a). More 321  
271 recent systems emphasize *memory evolution*, al- 322  
272 lowing memories to be edited, linked, or reorga- 323  
273 nized over time. Examples include A-MEM’s inter- 324  
274 connected note-style growth (Xu et al., 2025; Ka- 325  
275 davy, 2021), temporally adaptive structures such as 326  
276 Synapse, R2I, and SCM (Zheng et al., 2024c; Sam- 327  
277 sami et al., 2024; Wang et al., 2024a), as well as 328  
278 selective editing, recursive summarization, memory 329  
279 blending, and self-reflective verification to main- 330  
280 tain relevance and consistency (Bae et al., 2022; 331  
281 Wang et al., 2025c; Kim et al., 2024b; Sun et al., 332  
282 2024).

283 **Memory Query.** Memory query determines how 333  
284 relevant entries are selected and integrated to sup- 334  
285 port ongoing reasoning. Existing methods improve 335  
286 query effectiveness from complementary perspec-

287 tives. *Query-centered* approaches reformulate or 288  
289 refine the query itself, for example via forward- 290  
291 looking rewriting or iterative refinement (Jiang 291  
292 et al., 2023b; Jang et al., 2024). *Memory-centered* 292  
293 approaches enhance ranking and selection through 293  
294 richer indexing signals and reranking strategies, as 294  
295 explored in LongMemEval and personalized long- 295  
296 term memory retrieval (Wu et al., 2024a; Du et al., 296  
297 2024). Finally, *event- and structure-aware* retrieval 297  
298 leverages temporal, causal, or relational structure to 298  
299 traverse memory graphs or timelines, enabling co- 299  
300 herent recall across long interaction histories (e.g., 300  
301 LoCoMo, CC, MSC, and graph-based multi-hop 301  
302 retrieval (Maharana et al., 2024; Jang et al., 2023; 302  
303 Xu et al., 2021; Gutiérrez et al., 2024; Qian et al., 303  
304 2024)). Together, these strategies highlight that 304  
305 effective agentic memory query relies not only on 305  
306 semantic similarity, but also on adaptive, context- 306  
307 aware access over evolving memory states.

### 3 Intermediate Latent as Memory 306

307 Intermediate latent (IL) refer to inference-time in- 307  
308 ternal representations in LLMs, such as attention 308  
309 activations or other continuous vectors, that can 309  
310 be cached and reused as memory. This section fo- 310  
311 cuses on two forms of intermediate latent memory: 311  
312 the Key-Value (KV) cache, and other vector-based 312  
313 memory mechanisms. 313

*Goal:* Enable low-latency continuity within a session.

*Challenge:* State management under fixed capacity. IL memory is best viewed as *runtime state*. Construction, update (merge / compress / evict), and query must be explicitly scheduled to avoid churn, thrashing, and behavior drift.

#### 3.1 KV Cache as Memory 314

315 **Memory Construction.** The KV cache stores in- 315  
316 termediate key and value vectors produced by the 316  
317 attention mechanism to accelerate autoregressive 317  
318 decoding. Its construction is implicit and deter- 318  
319 ministic: during prefilling, KV pairs for all prompt 319  
320 tokens are computed and cached; during decod- 320  
321 ing, only KV pairs for newly generated tokens 321  
322 are appended, while attention reuses the cached 322  
323 states. This reduces per-token complexity from 323  
324  $\mathcal{O}(n^2)$  to  $\mathcal{O}(n)$  and is a core component of modern 324  
325 LLMs (Touvron et al., 2023a,b; Grattafiori et al., 325  
326 2024; DeepSeek-AI et al., 2025a,b). Unlike token- 326  
327 based memory, KV cache construction involves no 327  
328

explicit selection of *what* to store: all tokens initially contribute full KV representations, and memory management is deferred to post-construction. From a memory perspective, the KV cache thus provides a transient, high-fidelity record of recent context tightly coupled to attention computation.

**Memory Update.** Memory update for KV cache focuses on controlling storage cost while preserving attention quality. Existing methods fall into several recurring strategies.

*Eviction and dropping* selectively discard KV entries using static patterns or dynamic importance signals. Representative approaches include fixed sparsity schemes (Xiao et al., 2024a; Han et al., 2024a), layer-wise retention (Wu and Tu, 2024), and attention- or query-aware dropping such as H<sub>2</sub>O, FastGen, Radar, and NACL (Zhang et al., 2023a; Ge et al., 2024; Hao et al., 2025; Chen et al., 2024d). Variants further exploit attention statistics across heads, layers, or tasks (Liu et al., 2023e; Devoto et al., 2024; Yao et al., 2024; Jiang et al., 2024b; Zhong et al., 2024a; Zhou et al., 2025).

*Merging and semantic compression* reduce redundancy by consolidating similar KV entries rather than removing them outright. This includes similarity-based merging (Liu et al., 2024a; Kim et al., 2024a; Agarwal et al., 2024) and semantic-preserving compression at token, chunk, or sentence granularity (Zhang et al., 2025a; Liu et al., 2024d, 2025d; Zhu et al., 2025c).

*Quantization and low-rank approximation* lower per-entry storage cost by reducing numerical precision or exploiting low-rank structure. Representative methods apply low-bit or asymmetric quantization (Liu et al., 2024e; Duanmu et al., 2024; Dong et al., 2024b; Hooper et al., 2024; Zhang et al., 2024b; Li et al., 2025b), attention- or layer-aware scaling (Lin et al., 2025; Yang et al., 2024d), and low-rank compression with residual preservation (Dong et al., 2024a; Saxena et al., 2024; Kang et al., 2024). Dynamic precision schemes further adapt quantization to runtime conditions (Sheng et al., 2023; Zhao et al., 2024b; He et al., 2024b).

*System- and task-aware allocation* adapts KV storage to deployment constraints and task characteristics. Examples include disaggregated and multi-GPU KV storage (Chen et al., 2024c; Li et al., 2025a), layer- or chunk-level budget assignment (Yang et al., 2024a; Liu et al., 2025d), and preference- or workload-aware allocation strategies (Zhu et al., 2025b).

**Memory Query.** Memory query determines how cached KV states are accessed during attention when attending to all entries is inefficient or unnecessary. 1) *KV selection* restricts attention to a subset of relevant entries using query-dependent signals: QUEST and TokenSelect estimate token importance via attention statistics or learned predictors (Tang et al., 2024b; Wu et al., 2025a), Selective Attention prunes targets across heads and layers (Leviathan et al., 2025), and RetrievalAttention treats KV states as retrievable items using approximate nearest neighbor search (Liu et al., 2024b). 2) *KV reuse* avoids redundant computation by sharing cached states across overlapping contexts or requests. Prefix-based reuse organizes caches with tree structures, as in RadixAttention and ChunkAttention (Zheng et al., 2024b; Ye et al., 2024), while cross-request reuse shares KV states based on semantic similarity (Yang et al., 2025a; Agarwal et al., 2025; Tan et al., 2025). More recent systems extend reuse across retrieval and reranking stages in RAG pipelines (Yang et al., 2025b; Yao et al., 2025; Hu et al., 2025; Zhu et al., 2025a; An et al., 2025; Jiang et al., 2025), or externalize KV cache to support long-context inference beyond a single device (Wu et al., 2022; Tworkowski et al., 2023; Di et al., 2025).

### 3.2 Other Vectors as Memory

**External Vectors.** External vector memory augments LLMs with a separate vector store that retains intermediate latents for retrieval and reuse, alleviating the quadratic cost of long-context attention (Al Adel and Burtsev, 2021). Early work explored sentence-level memory slots for sequence modeling, while later systems such as kNN-LM and the Memorizing Transformer leveraged pre-trained embedding spaces and internal representations to enable scalable retrieval over large memory banks (Khandelwal et al., 2019; Wu et al., 2022). Subsequent designs, including MemGPT and Neurocache, maintain vector caches supporting dynamic retrieval and update for long-context or multi-session tasks (Packer et al., 2023; Safaya and Yuret, 2024). More recent architectures introduce structured or associative memory modules, such as CAMELoT, consolidating token representations while balancing novelty and recency (He et al., 2024c), and MemOS or Memory3, which externalize knowledge via vector memory with metadata or sparsification, without modifying core model parameters (Li et al., 2025c; Yang et al., 2024c).

**Steering Vectors.** Steering vectors function as intermediate memory for behavioral control rather than factual storage. Unlike interaction-heavy KV caches or external memories, they encode persistent biases as directions in activation space, originating from PPLM (Dathathri et al.). These vectors modulate hidden states to achieve alignment and interpretability via either contrastive or optimization-based approaches. *Contrastive methods* derive steering vectors from activation differences between datasets exhibiting desired versus undesired behaviors, encoding behavioral preferences as stable directions in activation space. Representative studies demonstrate control over sentiment, toxicity, refusal, or factuality using single or multiple contrastive prompts (Turner et al., 2023; Liu et al., 2023c; Zou et al., 2023; Arditì et al., 2024). While effective, these approaches often require carefully constructed contrastive data and may capture spurious correlations that limit robustness and generalization (Chughtai and Bushnaq, 2025). *Optimization-based methods* instead learn steering vectors by optimizing simple objectives, such as maximizing target sequence likelihood or applying lightweight affine transformations to hidden states, sometimes with only a single example (Subramani et al., 2022; Hernandez et al., 2023; Dunefsky and Cohan, 2025; Mack and Turner, 2024). More recent work explores probe-based or low-shot steering to induce truthfulness, suppress refusal, or support personalization, though empirical effectiveness varies across models and tasks (Li et al., 2024b; Turner et al., 2025; Cao et al., 2024).

#### 4 Parameter as Memory

*Goal:* Persist knowledge across sessions and deployments via weight updates.

*Challenge:* Safe and localized writing.

Parameter-level memory represents *consolidated knowledge*. Updates must add new information while limiting interference with existing capabilities.

**Memory Construction.** Unlike token- or cache-based memory that relies on input context and runtime storage, parametric memory encodes knowledge directly into model weights through pretraining or fine-tuning, providing long-term and context-independent storage. The capacity and structure of this memory are strongly influenced by training-time factors, including *training data composition*

*and augmentation, sequence or context length, and model scale.* Data augmentation strategies such as rephrasing, reordering, or stylistic transformation can significantly increase memorization strength (Allen-Zhu and Li, 2024), while data duplication leads to superlinear memorization effects and raises privacy concerns (Carlini et al., 2021; Lee et al., 2022; Kandpal et al., 2022). Longer training sequences expose models to richer contextual dependencies, increasing the likelihood of verbatim recall (Carlini et al., 2023; Wang et al., 2024). Model size further amplifies parametric memory capacity, with memorization scaling approximately log-linearly with parameter count (Tirumala et al., 2022; Carlini et al., 2023; Freeman et al., 2024). Mechanistic analyses suggest that this constructed memory is not uniformly distributed: MLP layers can be interpreted as key-value memories (Geva et al., 2021), and factual knowledge may localize to specific neurons (Dai et al., 2021), providing a structural basis for later updates.

**Memory Update.** Memory update at the parameter level concerns how knowledge embedded in model weights can be modified, extended, or reorganized without reconstructing parametric memory from scratch. Unlike construction, which writes memory through large-scale training, update mechanisms aim to incorporate new information, personalize behavior, or combine task knowledge while mitigating interference with existing parameters.

*Continual learning* provides a foundational view of parametric memory update by explicitly addressing catastrophic forgetting (Wang et al., 2024e). Regularization-based methods constrain updates to parameters deemed critical for previously learned knowledge, as in EWC, TaSL, SELF-PARAM, and POCL (Kirkpatrick et al., 2017; Feng et al., 2024; Wang et al.; Wu et al., 2024b). Replay-based strategies instead reinforce memory by reintroducing past samples or synthetic pseudo-data, enabling retention without full retraining; for example, DSI++ maintains retrieval performance via generative replay with pseudo-queries (Mehta et al., 2022). Agent-centric extensions such as LSCS further adapt continual learning to interactive settings, incrementally encoding external experiences into parameters over time (Wang et al., 2024k).

*Parameter-efficient fine-tuning (PEFT)* updates parametric memory by introducing lightweight, task- or user-specific adaptations while freezing the backbone (Han et al., 2024b). This paradigm

525 supports personalization and long-term adapta- 577  
526 tion with reduced computational cost. Representa- 578  
527 tive systems encode character traits or personal 579  
528 histories into parameters (e.g., Character- 580  
529 LLM (Shao et al.)), compress evolving user mem- 581  
530 ory through lifelong personal models (AI-Native 582  
531 Memory (Shang et al., 2024)), or enhance dialogue 583  
532 coherence via episodic parametric memory (Mem- 584  
533 oRAG (Qian et al., 2024), Echo (Liu et al., 2025b)). 585

534 *Model merging* updates parametric memory by 586  
535 combining multiple pretrained or fine-tuned mod- 587  
536 els without access to original data. Basic param- 588  
537 eter averaging, widely used in federated learning 589  
538 (e.g., FedAvg), offers simplicity and efficiency but 590  
539 often degrades performance due to parameter con- 591  
540 flicts (McMahan et al., 2017; Marczak et al., 2024). 592  
541 To address this, weighted merging prioritizes im- 593  
542 portant parameters using Taylor approximations, 594  
543 task vectors, or Fisher information (Lee et al., 2019; 595  
544 Qu et al., 2022; Matena and Raffel, 2022; Jhunjhun- 596  
545 wala et al., 2024; Daheim et al., 2024). Subspace- 597  
546 based approaches further mitigate conflicts by prun- 598  
547 ing or masking parameters before merging, leverag- 599  
548 ing over-parameterization to preserve task-relevant 600  
549 memory (e.g., TIES, DARE, Model Breadcrumbs, 601  
550 TALL-masks) (Yadav et al., 2023; Yu et al., 2024a; 602  
551 Davari and Belilovsky, 2023; Wang et al., 2024c). 603  
552 Routing-based methods generalize merging to infer- 604  
553 ence time, dynamically selecting or weighting 605  
554 parameters or experts based on inputs, as in MoE- 606  
555 style or soft-routing frameworks (Shazeer et al., 607  
556 2017; Muqeeth et al., 2024; Lu et al., 2024a; Tang 608  
557 et al., 2024a). LoRA-based routing further en- 609  
558 ables dynamic composition of low-rank updates, 610  
559 although decomposition-induced degradation re- 611  
560 mains a concern (Huang et al., 2023a; Wei et al., 612  
561 2025; Lai et al., 2025).

562 *Task arithmetic* treats parameter updates as vec- 613  
563 tor operations in weight space, enabling addition 614  
564 or subtraction of task-specific knowledge (Ilharco 615  
565 et al., 2022). Methods like TIES and AdaMerg- 616  
566 ing explicitly resolve conflicts by trimming low- 617  
567 magnitude parameters or learning adaptive merging 618  
568 coefficients (Yadav et al., 2023; Yang et al., 2024b), 619  
569 while TwinMerge refines task fusion through super- 620  
570 vised low-rank adaptation (Lu et al., 2024b). These 621  
571 approaches frame parameter update as algebraic 622  
572 composition, balancing flexibility and sensitivity 623  
573 to hyperparameters and task compatibility.

574 *Model editing* focuses on precise, targeted up- 624  
575 dates to parametric memory, enabling the insertion, 625  
576 modification, or deletion of specific facts (Wang

et al., 2024h). Recent systems such as Memo-  
ryLLM support self-updating and long-term re-  
tention, while WISE introduces a dual-memory  
design that separates stable pretrained knowledge  
from edited content and routes queries accord-  
ingly (Wang et al., 2024j,g). Compared to broader  
adaptation methods, model editing offers fine-  
grained control over parametric memory, but often  
relies on assumptions about knowledge localization  
and may accumulate errors over repeated edits.

**Memory Query.** Querying parametric memory  
differs fundamentally from querying token-level or  
KV-cache memory. Rather than retrieving stored  
entries, parametric memory is accessed implic-  
itly through forward computation and analyzed  
via *memorization phenomena*. **Exact memoriza-**  
**tion** refers to verbatim reproduction of training  
sequences under suitable prompts (Carlini et al.,  
2021, 2023; Nasr et al., 2023), while **approximate**  
**memorization** captures semantic or structural sim-  
ilarity without exact copying (Ippolito et al., 2023).  
**Prompt-based memorization** further shows that  
carefully designed prompts can elicit stored con-  
tent from partial prefixes, revealing the conditional  
nature of parametric recall (Biderman et al., 2023).

## 5 Discussion

**From Knowledge Requests to Memory Back-**  
**ends.** Table 1 links task knowledge requirements  
to representation choice by showing what domi-  
nant management challenge are. We summarize  
requirements using two human memory-inspired  
axes (Tulving and Donaldson, 1972; Begg, 1984;  
Squire, 2009): *retention* (short-term vs. long-term)  
and *functional form* (episodic, semantic, procedu-  
ral); details are in Appendix B. Token-level mem-  
ory primarily serves short-term explicit content,  
but is *query-limited* under long contexts, where se-  
lecting and integrating the right evidence becomes  
the bottleneck (Liu et al., 2024c). Intermediate la-  
tent memory supports short-horizon continuity, but  
is *update-limited* because gains depend on cache  
budgeting under fixed capacity (Xiao et al., 2023).  
Parametric memory supports long-term consolida-  
tion, but is *write-limited*: construction and update  
are costly and must control interference and forget-  
ting (Kirkpatrick et al., 2017; Meng et al., 2023).

**Unified Interfaces as System Glue.** The table  
suggests a simple system lesson: representation  
choice shifts the bottleneck to a different operation,

Memory Representation	Preferred Knowledge Requirements		Management Challenge			Strategic Gain (Return)
	Retention	Functional	Construction	Update	Query	
Token-level	Short-term	Episodic + Semantic	Medium	Medium	High	Editability
Intermediate latent	Short-term	Episodic	Low	High	Low	Efficiency
Parameter-level	Long-term	Procedural + Semantic	High	High	Low	Persistence

Table 1: Decision-support matrix for LLM memory management. Knowledge requirements summarize the task regimes each representation most naturally supports. Management challenge indicates where engineering effort is typically concentrated under construction-update-query.

so memory should be engineered through a unified construction-update-query interface rather than as representation-specific pipelines. With shared interfaces, systems can mix backends (tokens for editable evidence, caches for session continuity, parameters for durable consolidation) while keeping a consistent control plane for when to store, how to maintain, and what to use at inference time. This also improves portability across tasks and deployments by decoupling task requirements from backend choice, and makes failures easier to diagnose by localizing them to query selectivity (Liu et al., 2024c), cache policy (Xiao et al., 2023; Zhang et al., 2023c), or unsafe writes (Meng et al., 2023).

**Future Directions** *Specialized Memory Structures.* The growing scale and diversity of LLM memory workloads expose limitations of general-purpose memory hierarchies for long-context and agentic applications (Yao et al., 2022). Memory mechanisms such as KV caches (Pope et al., 2022), vector databases (Lewis et al., 2020), and graph-structured memories (Park et al., 2023) exhibit heterogeneous access patterns and update behaviors, motivating specialized memory structures for different abstractions. Meanwhile, the high cost of data movement highlights software-hardware co-design, where memory-aware algorithms and hardware support jointly reduce latency and energy (Wulf and McKee, 1995; Jouppi et al., 2017; Dao et al., 2022). Future work should explore dedicated memory layouts and tighter software-hardware integration to support memory construction, update, and query at scale (Yu et al., 2022a; Zhong et al., 2024d).

*Unified Training-Inference Systems.* Current LLM systems largely separate training from inference, limiting adaptation to evolving users and environments (Brown et al., 2020a; Bommasani et al., 2022). While this survey focuses on inference memory, long-term agentic settings increasingly blur this boundary through continual learning and personalization during deployment (Park et al.,

2023; Shinn et al., 2023). This trend motivates unified training-inference designs that integrate memory management with lightweight updates and inference-time adaptation while mitigating catastrophic forgetting (Hu et al., 2021). Future work should explore integrated architectures to enable continuous learning and personalized behavior.

*Cross-Domain Methodology Transfer.* LLM memory challenges often parallel classic problems in operating systems (OS) and databases. For example, KV cache management adopts OS-level paradigms like paging, eviction, and tiered hierarchies to handle limited memory (Kwon et al., 2023; Xiao et al., 2024b; Zhang et al., 2023a). RAG systems leverage database techniques for indexing, query optimization, and execution planning (Karpukhin et al., 2020; Khattab and Zaharia, 2020; Izacard et al., 2023; Asai et al., 2023). Furthermore, distributed systems and cloud computing inspire solutions for scaling long-context workloads through partitioning and remote memory (Zhong et al., 2024d; Fu et al., 2024; Jin and Wu, 2025). These established paradigms provide a foundational blueprint for scalable and efficient LLM memory orchestration.

## 6 Conclusion

This survey presents a unified management view of LLM memory that links task knowledge requirements to concrete memory representations through shared interfaces, including construction, update, and query. By treating memory as a system-level capability rather than a collection of task-specific tricks, the survey provides a coherent way to compare diverse designs, reason about effectiveness-efficiency trade-offs, and guide the composition of hybrid memory backends for long-horizon agents. We hope this perspective helps standardize how future work specifies requirements, evaluates memory behavior over extended interaction, and designs reusable memory components that remain reliable under realistic deployment constraints.

## 709 Limitations

710 This survey proposes a unified representation–  
711 management abstraction to organize LLM memory,  
712 but the fast pace of the field means our coverage  
713 may lag behind the newest systems and some in-  
714 dustrial practices are discussed only at a high level.  
715 Our taxonomy is also a simplification: many meth-  
716 ods span multiple representations and the bound-  
717 aries between construction, update, and query can  
718 blur in long-horizon agents. Finally, quantitative  
719 comparisons across papers remain limited due to  
720 inconsistent tasks, models, evaluation protocols,  
721 and deployment settings, so our synthesis empha-  
722 sizes recurring design trade-offs and failure modes  
723 rather than a unified benchmark ranking.

## 724 References

725 Saurabh Agarwal, Bilge Acun, Basil Hosmer, Mostafa  
726 Elhoushi, Yejin Lee, Shivaram Venkataraman, Dim-  
727 itris Papailiopoulos, and Carole-Jean Wu. 2024.  
728 **CHAI: Clustered head attention for efficient LLM**  
729 **inference**. In *Proceedings of the 41st International*  
730 *Conference on Machine Learning*, volume 235 of  
731 *Proceedings of Machine Learning Research*, pages  
732 291–312. PMLR.

733 Shubham Agarwal, Sai Sundaresan, Subrata Mitra, De-  
734 babrata Mahapatra, Archit Gupta, Rounak Sharma,  
735 Nirmal Joshua Kapu, Tong Yu, and Shiv Saini.  
736 2025. **Cache-craft: Managing chunk-caches for**  
737 **efficient retrieval-augmented generation**. *Preprint*,  
738 arXiv:2502.15734.

739 Talkie Ai. 2024. **Talkie | ai-native character community**.

740 Arij Al Adel and Mikhail S. Burtsev. 2021. **Memory**  
741 **transformer with hierarchical attention for long docu-**  
742 **ment processing**. In *2021 International Conference*  
743 *Engineering and Telecommunication*, pages 1–7.

744 Zeyuan Allen-Zhu and Yuanzhi Li. 2024. Physics of  
745 language models: part 3.1, knowledge storage and  
746 extraction. In *Proceedings of the 41st International*  
747 *Conference on Machine Learning*, pages 1067–1077.

748 Yuwei An, Yihua Cheng, Seo Jin Park, and Junchen  
749 Jiang. 2025. **Hyperrag: Enhancing quality-**  
750 **efficiency tradeoffs in retrieval-augmented gener-**  
751 **ation with reranker kv-cache reuse**. *Preprint*,  
752 arXiv:2504.02921.

753 Anysphere. 2025. **Cursor - the ai code editor**. <https://www.cursor.com/en>.

755 A. Ardit, O. Obeso, A. Syed, D. Paleka, N. Panickssery,  
756 W. Gurnee, and N. Nanda. 2024. **Refusal in language**  
757 **models is mediated by a single direction**. *arXiv*.

Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and  
Hannaneh Hajishirzi. 2023. **Self-rag: Learning to**  
retrieve, generate, and critique through self-reflection.  
*arXiv preprint arXiv:2310.11511*.

BAAI. 2023. **Flagembedding**. <https://github.com/FlagOpen/FlagEmbedding>.

Alan Baddeley. 2007. *Working memory, thought, and*  
*action*, volume 45. OuP Oxford.

Alan D. Baddeley and Graham Hitch. 1974. **Working**  
**memory**. volume 8 of *Psychology of Learning and*  
*Motivation*, pages 47–89. Academic Press.

Sanghwan Bae, Donghyun Kwak, Soyoung Kang,  
Min Young Lee, Sungdong Kim, Yui Jeong, Hyeri  
Kim, Sang-Woo Lee, Woomyoung Park, and Nako  
Sung. 2022. **Keep me updated! memory manage-**  
**ment in long-term conversations**. In *Findings of the*  
*Association for Computational Linguistics: EMNLP*  
*2022*, pages 3769–3787, Abu Dhabi, United Arab  
Emirates. Association for Computational Linguistics.

Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang,  
Xiaodong Deng, and 1 others. 2023. **Qwen technical**  
**report**. *Preprint*, arXiv:2309.16609.

I. Begg. 1984. Tulving’s memory [review of the book  
elements of episodic memory, by e. tulving]. *Can-*  
*adian Journal of Psychology / Revue canadienne de*  
*psychologie*, 38(1):144–147.

Iz Beltagy, Matthew E. Peters, and Arman Cohan. 2020.  
**Longformer: The long-document transformer**. *CoRR*,  
abs/2004.05150.

Stella Biderman, USVSN Sai Prashanth, Lintang  
Sutawika, Hailey Schoelkopf, Quentin Anthony,  
Shivanshu Purohit, and Edward Raff. 2023. **Emer-**  
**gent and predictable memorization in large language**  
**models**. *Preprint*, arXiv:2304.11158.

Rishi Bommasani, Drew A. Hudson, Ehsan Adeli, Russ  
Altman, Simran Arora, Sydney von Arx, Michael S.  
Bernstein, Jeannette Bohg, Antoine Bosselut, Emma  
Brunskill, Erik Brynjolfsson, Shyamal Buch, Dallas  
Card, Rodrigo Castellon, Niladri Chatterji, Annie  
Chen, Kathleen Creel, Jared Quincy Davis, Dora  
Demszky, and 95 others. 2022. **On the opportu-**  
**nities and risks of foundation models**. *Preprint*,  
arXiv:2108.07258.

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie  
Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind  
Neelakantan, Pranav Shyam, Girish Sastry, Amanda  
Askell, Sandhini Agarwal, Ariel Herbert-Voss,  
Gretchen Krueger, Tom Henighan, Rewon Child,  
Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu,  
Clemens Winter, and 12 others. 2020a. Language  
models are few-shot learners. In *Proceedings of the*  
*34th International Conference on Neural Information*  
*Processing Systems, NIPS ’20*, Red Hook, NY, USA.  
Curran Associates Inc.

812	Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie	Yilong Chen, Guoxia Wang, Junyuan Shang, Shiyao	865
813	Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind	Cui, Zhenyu Zhang, Tingwen Liu, Shuohuan Wang,	866
814	Neelakantan, Pranav Shyam, Girish Sastry, Amanda	Yu Sun, Dianhai Yu, and Hua Wu. 2024d. <a href="#">NACL:</a>	867
815	Askill, Sandhini Agarwal, Ariel Herbert-Voss,	<a href="#">A general and effective KV cache eviction frame-</a>	868
816	Gretchen Krueger, Tom Henighan, Rewon Child,	<a href="#">work for LLM at inference time.</a> In <i>Proceedings</i>	869
817	Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu,	<i>of the 62nd Annual Meeting of the Association for</i>	870
818	Clemens Winter, and 12 others. 2020b. <a href="#">Lan-</a>	<i>Computational Linguistics (Volume 1: Long Papers),</i>	871
819	<a href="#">guage models are few-shot learners.</a> <i>Preprint,</i>	pages 7913–7926, Bangkok, Thailand. Association	872
820	arXiv:2005.14165.	for Computational Linguistics.	873
821	Andrew E Budson and Elizabeth A Kensinger. 2023.	Zi-Yi Chen, Fan-Kai Xie, Meng Wan, Yang Yuan, Miao	874
822	<a href="#">Why we forget and how to remember better: the sci-</a>	Liu, Zong-Guo Wang, Sheng Meng, and Yan-Gang	875
823	<a href="#">ence behind memory.</a> Oxford University Press.	Wang. 2023c. Matchat: A large language model and	876
824	Y. Cao, T. Zhang, B. Cao, Z. Yin, L. Lin, F. Ma, and	application service platform for materials science.	877
825	J. Chen. 2024. <a href="#">Personalized steering of large lan-</a>	<i>Chinese Physics B</i> , 32(11):118104.	878
826	<a href="#">guage models: Versatile steering vectors through bi-</a>	Xin Cheng, Di Luo, Xiuying Chen, Lemao Liu,	879
827	<a href="#">directional preference optimization.</a> <i>arXiv.</i>	Dongyan Zhao, and Rui Yan. 2023. Lift yourself	880
828	Nicholas Carlini, Daphne Ippolito, Matthew Jagielski,	up: Retrieval-augmented text generation with self	881
829	Katherine Lee, Florian Tramèr, and Chiyuan Zhang.	memory. <i>arXiv preprint arXiv:2305.02437.</i>	882
830	2023. <a href="#">Quantifying memorization across neural lan-</a>	Yuheng Cheng, Ceyao Zhang, Zhengwen Zhang, Xi-	883
831	<a href="#">guage models.</a> <i>Preprint</i> , arXiv:2202.07646.	angrui Meng, Sirui Hong, Wenhao Li, Zihao Wang,	884
832	Nicholas Carlini, Florian Tramèr, Eric Wallace,	Zekai Wang, Feng Yin, Junhua Zhao, and et al. 2024.	885
833	Matthew Jagielski, Ariel Herbert-Voss, Katherine	<a href="#">Exploring large language model based intelligent</a>	886
834	Lee, Adam Roberts, Tom Brown, Dawn Song, Ul-	<a href="#">agents: Definitions, methods, and prospects.</a> <i>arXiv.</i>	887
835	far Erlingsson, Alina Oprea, and Colin Raffel. 2021.	Alexis Chevalier, Alexander Wettig, Anirudh Ajith, and	888
836	<a href="#">Extracting training data from large language models.</a>	Danqi Chen. 2023. <a href="#">Adapting language models to</a>	889
837	<i>Preprint</i> , arXiv:2012.07805.	<a href="#">compress contexts.</a> In <i>Proceedings of the 2023 Con-</i>	890
838	Character AI. 2023. <a href="#">Character ai.</a> Retrieved September	<i>ference on Empirical Methods in Natural Language</i>	891
839	14, 2023 from <a href="https://character.ai/">https://character.ai/</a> .	<i>Processing</i> , pages 3829–3846, Singapore. Associa-	892
840	Dake Chen, Hanbin Wang, Yunhao Huo, Yuzhao Li, and	tion for Computational Linguistics.	893
841	Haoyang Zhang. 2023a. Gamept: Multi-agent col-	B. Chughtai and L. Bushnaq. 2025. <a href="#">Activation space</a>	894
842	laborative framework for game development. <i>arXiv</i>	<a href="#">interpretability may be doomed.</a>	895
843	<i>preprint arXiv:2310.08067.</i>	Fergus IM Craik and Robert S Lockhart. 1972. Lev-	896
844	Jianlv Chen, Shitao Xiao, Peitian Zhang, Kun Luo, Defu	els of processing: A framework for memory re-	897
845	Lian, and Zheng Liu. 2024a. <a href="#">BGE m3-embedding:</a>	search. <i>Journal of verbal learning and verbal be-</i>	898
846	<a href="#">Multi-lingual, multi-functionality, multi-granularity</a>	<i>havior</i> , 11(6):671–684.	899
847	<a href="#">text embeddings through self-knowledge distillation.</a>	Can Cui, Yunsheng Ma, Xu Cao, Wenqian Ye, Yang	900
848	<i>CoRR</i> , abs/2402.03216.	Zhou, Kaizhao Liang, Jintai Chen, Juanwu Lu, Zi-	901
849	Kexin Chen, Junyou Li, Kunyi Wang, Yuyang Du, Jiahui	chong Yang, Kuei-Da Liao, and et al. 2024. A sur-	902
850	Yu, Jiamin Lu, Lanqing Li, Jiezhong Qiu, Jianzhang	vey on multimodal large language models for au-	903
851	Pan, Yi Huang, Qun Fang, Pheng Ann Heng, and	tonomous driving. In <i>Proceedings of the IEEE/CVF</i>	904
852	Guangyong Chen. 2024b. <a href="#">Chemist-x: Large lan-</a>	<i>Winter Conference on Applications of Computer Vi-</i>	905
853	<a href="#">guage model-empowered agent for reaction condition</a>	<i>sion</i> , pages 958–979.	906
854	<a href="#">recommendation in chemical synthesis.</a>	Nico Daheim, Thomas Möllenhoff, Edoardo Ponti,	907
855	Shiyang Chen, Rain Jiang, Dezhi Yu, Jinlai Xu,	Iryna Gurevych, and Mohammad Emtiyaz Khan.	908
856	Mengyuan Chao, Fanlong Meng, Chenyu Jiang,	2024. Model merging by uncertainty-based gradi-	909
857	Wei Xu, and Hang Liu. 2024c. <a href="#">Kvdirect: Dis-</a>	ent matching. In <i>ICLR</i> .	910
858	<a href="#">tributed disaggregated llm inference.</a> <i>Preprint,</i>	Damai Dai, Li Dong, Yaru Hao, Zhifang Sui, Baobao	911
859	arXiv:2501.14743.	Chang, and Furu Wei. 2021. <a href="#">Knowledge neurons in</a>	912
860	Tong Chen, Hongwei Wang, Sihao Chen, Wenhao	<a href="#">pretrained transformers.</a> <i>arXiv.</i>	913
861	Yu, Kaixin Ma, Xinran Zhao, Dong Yu, and Hong-	Damai Dai, Yutao Sun, Li Dong, Yaru Hao, Shuming	914
862	ming Zhang. 2023b. Dense x retrieval: What re-	Ma, Zhifang Sui, and Furu Wei. 2023. <a href="#">Why can GPT</a>	915
863	trieval granularity should we use? <i>arXiv preprint</i>	<a href="#">learn in-context? language models secretly perform</a>	916
864	<i>arXiv:2312.06648.</i>	<a href="#">gradient descent as meta-optimizers.</a> In <i>Findings of</i>	917
		<i>the Association for Computational Linguistics: ACL</i>	918
		2023, pages 4005–4019, Toronto, Canada. Associa-	919
		tion for Computational Linguistics.	920



1032	Muhan Gao, TaiMing Lu, Kuai Yu, Adam Byerly, and Daniel Khashabi. 2024. Insights into LLM long-context failures: When transformers know but don't tell. In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , Miami, Florida, USA. Association for Computational Linguistics.	1089
1033		1090
1034		1091
1035		1092
1036		1093
1037		
1038	Shen Gao, Xiuying Chen, Piji Li, Zhaochun Ren, Lidong Bing, Dongyan Zhao, and Rui Yan. 2019. <a href="#">Abstractive text summarization by incorporating reader comments</a> . In <i>The Thirty-Third AAAI Conference on Artificial Intelligence, AAAI 2019, The Thirty-First Innovative Applications of Artificial Intelligence Conference, IAAI 2019, The Ninth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2019, Honolulu, Hawaii, USA, January 27 - February 1, 2019</i> , pages 6399–6406. AAAI Press.	1094
1039		1095
1040		1096
1041		1097
1042		1098
1043		
1044		1099
1045		1100
1046		1101
1047		
1048	Shivam Garg, Dimitris Tsipras, Percy S Liang, and Gregory Valiant. 2022. What can transformers learn in-context? A case study of simple function classes. In <i>Advances in Neural Information Processing Systems</i> , volume 35, pages 30583–30598.	1102
1049		1103
1050		1104
1051		1105
1052		1106
1053	Suyu Ge, Yunan Zhang, Liyuan Liu, Minjia Zhang, Jiawei Han, and Jianfeng Gao. 2024. <a href="#">Model tells you what to discard: Adaptive KV cache compression for LLMs</a> . In <i>The Twelfth International Conference on Learning Representations</i> .	1107
1054		1108
1055		1109
1056		1110
1057		
1058	Yingqiang Ge, Yujie Ren, Wenyue Hua, Shuyuan Xu, Juntao Tan, and Yongfeng Zhang. 2023. <a href="#">Llm as os (llmao), agents as apps: Envisioning aios, agents and the aios-agent ecosystem</a> . <i>arXiv</i> .	1111
1059		1112
1060		1113
1061		1114
1062	Mor Geva, Roei Schuster, Jonathan Berant, and Omer Levy. 2021. <a href="#">Transformer feed-forward layers are key-value memories</a> . <i>arXiv</i> .	1115
1063		1116
1064		1117
1065		1118
1066		1119
1067	In Gim, Guojun Chen, Seung-seob Lee, Nikhil Sarda, Anurag Khandelwal, and Lin Zhong. 2024. Prompt cache: Modular attention reuse for low-latency inference. <i>Proceedings of Machine Learning and Systems</i> , 6:325–338.	1120
1068		1121
1069		1122
1070	GitHub. 2022. <a href="#">Github copilot</a> .	1123
1071		1124
1072		1125
1073		1126
1074		1127
1075	Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, and 1 others. 2024. <a href="#">The llama 3 herd of models</a> . <i>Preprint</i> , arXiv:2407.21783.	1128
1076		1129
1077		1130
1078	Michael Günther, Jackmin Ong, Isabelle Mohr, Alaeddine Abdessalem, Tanguy Abel, Mohammad Kalim Akram, Susana Guzman, Georgios Mastrapas, Saba Sturua, Bo Wang, Maximilian Werk, Nan Wang, and Han Xiao. 2023. <a href="#">Jina embeddings 2: 8192-token general-purpose text embeddings for long documents</a> . <i>CoRR</i> , abs/2310.19923.	1131
1079		1132
1080		1133
1081		1134
1082		1135
1083		1136
1084		1137
1085	Mandy Guo, Joshua Ainslie, David C. Uthus, Santiago Ontañón, Jianmo Ni, Yun-Hsuan Sung, and Yinfei Yang. 2022. <a href="#">Longt5: Efficient text-to-text transformer for long sequences</a> . In <i>Findings of the Association for Computational Linguistics: NAACL 2022, Seattle, WA, United States, July 10-15, 2022</i> , pages 724–736. Association for Computational Linguistics.	1138
1086		1139
1087		1140
1088		1141
		1142
	Taicheng Guo, Xiuying Chen, Yaqi Wang, Ruidi Chang, Shichao Pei, Nitesh V Chawla, Olaf Wiest, and Xi-angliang Zhang. 2024. <a href="#">Large language model based multi-agents: A survey of progress and challenges</a> . <i>arXiv</i> .	1089
		1090
		1091
		1092
		1093
	Bernal Jiménez Gutiérrez, Yiheng Shu, Yu Gu, Michihiro Yasunaga, and Yu Su. 2024. Hipporag: Neurobiologically inspired long-term memory for large language models. In <i>The Thirty-eighth Annual Conference on Neural Information Processing Systems</i> .	1094
		1095
		1096
		1097
		1098
	Michael Hahn and Navin Goyal. 2023. <a href="#">A theory of emergent in-context learning as implicit structure induction</a> . <i>arxiv</i> , arXiv:2303.07971.	1099
		1100
		1101
	Chi Han, Qifan Wang, Hao Peng, Wenhan Xiong, Yu Chen, Heng Ji, and Sinong Wang. 2024a. <a href="#">LM-infinite: Zero-shot extreme length generalization for large language models</a> . In <i>Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 3991–4008, Mexico City, Mexico. Association for Computational Linguistics.	1102
		1103
		1104
		1105
		1106
		1107
		1108
		1109
		1110
	Zeyu Han, Chao Gao, Jinyang Liu, Jeff Zhang, and Sai Qian Zhang. 2024b. <a href="#">Parameter-efficient fine-tuning for large models: A comprehensive survey</a> . <i>arXiv preprint arXiv:2403.14608</i> .	1111
		1112
		1113
		1114
	Yongchang Hao, Mengyao Zhai, Hossein Hajimirsadeghi, Sepidehsadat Hosseini, and Frederick Tung. 2025. <a href="#">Radar: Fast long-context decoding for any transformer</a> . In <i>The Thirteenth International Conference on Learning Representations</i> .	1115
		1116
		1117
		1118
		1119
	Kostas Hatalis, Despina Christou, Joshua Myers, Steven Jones, Keith Lambert, Adam Amos-Binks, Zohreh Dannenhauer, and Dustin Dannenhauer. 2024. Memory matters: The need to improve long-term memory in llm-agents. <i>Proceedings of the AAAI Symposium Series</i> , 2.	1120
		1121
		1122
		1123
		1124
		1125
	Kai He, Rui Mao, Qika Lin, Yucheng Ruan, Xiang Lan, Mengling Feng, and Erik Cambria. 2023a. <a href="#">A survey of large language models for healthcare: from data, technology, and applications to accountability and ethics</a> . <i>arXiv</i> .	1126
		1127
		1128
		1129
		1130
	Tianyu He, Guanghui Fu, Yijing Yu, Fan Wang, Jianqiang Li, Qing Zhao, Changwei Song, Hongzhi Qi, Dan Luo, Huijing Zou, and et al. 2023b. <a href="#">Towards a psychological generalist ai: A survey of current applications of large language models and future prospects</a> . <i>arXiv</i> .	1131
		1132
		1133
		1134
		1135
		1136
	Xiaoxin He, Yijun Tian, Yifei Sun, Nitesh V Chawla, Thomas Laurent, Yann LeCun, Xavier Bresson, and Bryan Hooi. 2024a. <a href="#">G-retriever: Retrieval-augmented generation for textual graph understanding and question answering</a> . <i>arXiv preprint arXiv:2402.07630</i> .	1137
		1138
		1139
		1140
		1141
		1142

1143	Yefei He, Luoming Zhang, Weijia Wu, Jing Liu, Hong Zhou, and Bohan Zhuang. 2024b. <a href="#">Zipcache: Accurate and efficient kv cache quantization with salient token identification</a> . In <i>Advances in Neural Information Processing Systems</i> , volume 37, pages 68287–68307. Curran Associates, Inc.	1198
1144		1199
1145		1200
1146		1201
1147		
1148		
1149	Zexue He, Leonid Karlinsky, Donghyun Kim, Julian McAuley, Dmitry Krotov, and Rogerio Feris. 2024c. Camelot: Towards large language models with training-free consolidated associative memory. <i>arXiv preprint arXiv:2402.13449</i> .	1202
1150		1203
1151		1204
1152		1205
1153		1206
1154	Zihong He, Weizhe Lin, Hao Zheng, Fan Zhang, Matt W. Jones, Laurence Aitchison, Xuhai Xu, Miao Liu, Per Ola Kristensson, and Junxiao Shen. 2024d. <a href="#">Human-inspired perspectives: A survey on ai long-term memory</a> . <i>arXiv preprint arXiv:2411.00489</i> .	1207
1155		
1156		
1157		
1158		
1159	E. Hernandez, B. Z. Li, and J. Andreas. 2023. <a href="#">Inspecting and editing knowledge representations in language models</a> . <i>arXiv</i> .	1208
1160		1209
1161		1210
1162	Christian Herold and Hermann Ney. 2023. <a href="#">Improving long context document-level machine translation</a> . <i>CoRR</i> , abs/2306.05183.	1211
1163		1212
1164		1213
1165	Lukas Hilgert, Danni Liu, and Jan Niehues. 2024. <a href="#">Evaluating and training long-context large language models for question answering on scientific papers</a> . In <i>Proceedings of the 1st Workshop on Customizable NLP: Progress and Challenges in Customizing NLP for a Domain, Application, Group, or Individual (CustomNLP4U)</i> , pages 220–236, Miami, Florida, USA. Association for Computational Linguistics.	1214
1166		1215
1167		1216
1168		1217
1169		
1170		
1171		
1172		
1173	Coleman Hooper, Sehoon Kim, Hiva Mohammadzadeh, Michael W Mahoney, Sophia Shao, Kurt Keutzer, and Amir Gholami. 2024. KVQuant: Towards 10 million context length LLM inference with KV cache quantization. <i>Advances in Neural Information Processing Systems, NeurIPS 2024</i> , 37:1270–1303.	1218
1174		1219
1175		
1176		
1177		
1178		
1179	Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. <a href="#">Lora: Low-rank adaptation of large language models</a> . <i>Preprint</i> , arXiv:2106.09685.	1220
1180		1221
1181		1222
1182		1223
1183	Junhao Hu, Wenrui Huang, Weidong Wang, Haoyi Wang, Tiancheng Hu, Qin Zhang, Hao Feng, Xusheng Chen, Yizhou Shan, and Tao Xie. 2025. <a href="#">Epic: Efficient position-independent caching for serving large language models</a> . <i>Preprint</i> , arXiv:2410.15332.	1224
1184		1225
1185		1226
1186		1227
1187		
1188		
1189	Wenyue Hua, Lizhou Fan, Lingyao Li, Kai Mei, Jianchao Ji, Yingqiang Ge, Libby Hemphill, and Yongfeng Zhang. 2023. War and peace (waragent): Large language model-based multi-agent simulation of world wars. <i>arXiv preprint arXiv:2311.17227</i> .	1228
1190		1229
1191		1230
1192		1231
1193		1232
1194	Chengsong Huang, Qian Liu, Bill Yuchen Lin, Tianyu Pang, Chao Du, and Min Lin. 2023a. <a href="#">Lorahub: Efficient cross-task generalization via dynamic lora composition</a> .	1233
1195		1234
1196		1235
1197		1236
		1237
		1238
		1239
		1240
		1241
		1242
		1243
		1244
		1245
		1246
		1247
		1248
		1249
		1250
		1251
		1252
		1253
		1254
		1255

1256	Divyansh Jhunjunwala, Shiqiang Wang, and Gauri Joshi. 2024. Fedfisher: Leveraging fisher information for one-shot federated learning. In <i>AISTATS</i> , pages 1612–1620. PMLR.	1313
1257		1314
1258		1315
1259		1316
1260	Huiqiang Jiang, Qianhui Wu, Chin-Yew Lin, Yuqing Yang, and Lili Qiu. 2023a. <a href="#">LLMLingua: Compressing prompts for accelerated inference of large language models</a> . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 13358–13376, Singapore. Association for Computational Linguistics.	1317
1261		1318
1262		1319
1263		1320
1264		1321
1265		1322
1266		1323
1267	Wenqi Jiang, Suvinay Subramanian, Cat Graves, Gustavo Alonso, Amir Yazdanbakhsh, and Vidushi Dadu. 2025. <a href="#">Rago: Systematic performance optimization for retrieval-augmented generation serving</a> . <i>Preprint</i> , arXiv:2503.14649.	1324
1268		1325
1269		1326
1270		1327
1271		1328
1272	Xun Jiang, Feng Li, Han Zhao, Jiaying Wang, Jun Shao, Shihao Xu, Shu Zhang, Weiling Chen, Xavier Tang, Yize Chen, and et al. 2024a. <a href="#">Long term memory: The foundation of ai self-evolution</a> . <i>arXiv</i> .	1329
1273		1330
1274		1331
1275		1332
1276	Yikun Jiang, Huanyu Wang, Lei Xie, Hanbin Zhao, Chao Zhang, Hui Qian, and John C.S. Lui. 2024b. <a href="#">D-llm: A token adaptive computing resource allocation strategy for large language models</a> . In <i>Advances in Neural Information Processing Systems</i> , volume 37, pages 1725–1749. Curran Associates, Inc.	1333
1277		1334
1278		1335
1279		1336
1280		1337
1281		1338
1282	Zhengbao Jiang, Frank Xu, Luyu Gao, Zhiqing Sun, Qian Liu, Jane Dwivedi-Yu, Yiming Yang, Jamie Callan, and Graham Neubig. 2023b. <a href="#">Active retrieval augmented generation</a> . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 7969–7992, Singapore. Association for Computational Linguistics.	1339
1283		1340
1284		1341
1285		1342
1286		1343
1287		1344
1288		1345
1289	Zhengbao Jiang, Frank F Xu, Luyu Gao, Zhiqing Sun, Qian Liu, Jane Dwivedi-Yu, Yiming Yang, Jamie Callan, and Graham Neubig. 2023c. <a href="#">Active retrieval augmented generation</a> . <i>arXiv preprint arXiv:2305.06983</i> .	1346
1290		1347
1291		1348
1292		1349
1293		1350
1294	Bowen Jin, Jinsung Yoon, Jiawei Han, and Sercan Ö. Arik. 2024a. <a href="#">Long-context llms meet RAG: overcoming challenges for long inputs in RAG</a> . <i>CoRR</i> , abs/2410.05983.	1351
1295		1352
1296		1353
1297		1354
1298	Bowen Jin, Hansi Zeng, Guoyin Wang, Xiusi Chen, Tianxin Wei, Ruirui Li, Zhengyang Wang, Zheng Li, Yang Li, Hanqing Lu, and 1 others. 2023. <a href="#">Language models as semantic indexers</a> . <i>arXiv preprint arXiv:2310.07815</i> .	1355
1299		1356
1300		1357
1301		1358
1302		1359
1303	Hanlei Jin, Yang Zhang, Dan Meng, Jun Wang, and Jinghua Tan. 2024b. <a href="#">A comprehensive survey on process-oriented automatic text summarization with exploration of llm-based methods</a> . <i>CoRR</i> , abs/2403.02901.	1360
1304		1361
1305		1362
1306		1363
1307		1364
1308	Hongpeng Jin and Yanzhao Wu. 2025. <a href="#">Ce-collm: Efficient and adaptive large language models through cloud-edge collaboration</a> . In <i>2025 IEEE International Conference on Web Services (ICWS)</i> , pages 316–323.	1365
1309		1366
1310		1367
1311		1368
1312		1369
	Philip Nicholas Johnson-Laird. 1983. <i>Mental models: Towards a cognitive science of language, inference, and consciousness</i> , volume 6. Harvard University Press.	1370
	Norman P. Jouppi, Cliff Young, Nishant Patil, David Patterson, Gaurav Agrawal, Raminder Bajwa, Sarah Bates, Suresh Bhatia, Nan Boden, Al Borchers, Rick Boyle, Pierre-luc Cantin, Clifford Chao, Chris Clark, Jeremy Coriell, Mike Daley, Matt Dau, Jeffrey Dean, Ben Gelb, and 57 others. 2017. <a href="#">In-datacenter performance analysis of a tensor processing unit</a> . <i>SIGARCH Comput. Archit. News</i> , 45(2):1–12.	1371
	David Kadavy. 2021. <i>Digital Zettelkasten: Principles, Methods, and Examples</i> . Kadavy, Inc.	1372
	Zhao Kaiya, Michelangelo Naim, Jovana Kondic, Manuel Cortes, Jiaxin Ge, Shuying Luo, Guangyu Robert Yang, and Andrew Ahn. 2023. <a href="#">Lyfe agents: Generative agents for low-cost real-time social interactions</a> . <i>arXiv</i> .	1373
	Nikhil Kandpal, Eric Wallace, and Colin Raffel. 2022. <a href="#">Deduplicating training data mitigates privacy risks in language models</a> . <i>Preprint</i> , arXiv:2202.06539.	1374
	Hao Kang, Qingru Zhang, Souvik Kundu, Geonhwa Jeong, Zaoxing Liu, Tushar Krishna, and Tuo Zhao. 2024. <a href="#">Gear: An efficient kv cache compression recipe for near-lossless generative inference of llm</a> . <i>Preprint</i> , arXiv:2403.05527.	1375
	Sayash Kapoor, Peter Henderson, and Arvind Narayanan. 2024. <a href="#">Promises and pitfalls of artificial intelligence for legal applications</a> . <i>CoRR</i> , abs/2402.01656.	1376
	Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick SH Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. <a href="#">Dense passage retrieval for open-domain question answering</a> . In <i>EMNLP (1)</i> , pages 6769–6781.	1377
	Zixuan Ke, Weize Kong, Cheng Li, Mingyang Zhang, Qiaozhu Mei, and Michael Bendersky. 2024. <a href="#">Bridging the preference gap between retrievers and llms</a> . <i>arXiv preprint arXiv:2401.06954</i> .	1378
	Urvashi Khandelwal, Omer Levy, Dan Jurafsky, Luke Zettlemoyer, and Mike Lewis. 2019. <a href="#">Generalization through memorization: Nearest neighbor language models</a> . <i>arXiv preprint arXiv:1911.00172</i> .	1379
	Omar Khattab, Keshav Santhanam, Xiang Lisa Li, David Hall, Percy Liang, Christopher Potts, and Matei Zaharia. 2022. <a href="#">Demonstrate-search-predict: Composing retrieval and language models for knowledge-intensive nlp</a> . <i>arXiv preprint arXiv:2212.14024</i> .	1380
	Omar Khattab and Matei Zaharia. 2020. <a href="#">Colbert: Efficient and effective passage search via contextualized late interaction over bert</a> . In <i>Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR</i>	1381

1368	'20, page 39–48, New York, NY, USA. Association	<i>Computational Linguistics</i> . Association for Compu-	1424
1369	for Computing Machinery.	tational Linguistics.	1425
1370	Minsoo Kim, Kyuhong Shim, Jungwook Choi, and	Namhoon Lee, Thalaiyasingam Ajanthan, and Philip	1426
1371	Simyung Chang. 2024a. <a href="#">InfiniPot: Infinite context</a>	Torr. 2019. <a href="#">SNIP: SINGLE-SHOT NETWORK</a>	1427
1372	<a href="#">processing on memory-constrained LLMs</a> . In <i>Pro-</i>	<a href="#">PRUNING BASED ON CONNECTION SENSITIV-</a>	1428
1373	<i>ceedings of the 2024 Conference on Empirical Meth-</i>	<a href="#">ITY</a> . In <i>International Conference on Learning Rep-</i>	1429
1374	<i>ods in Natural Language Processing</i> , pages 16046–	<i>resentations</i> .	1430
1375	16060, Miami, Florida, USA. Association for Com-	Yaniv Leviathan, Matan Kalman, and Yossi Matias.	1431
1376	putational Linguistics.	2025. <a href="#">Selective attention improves transformer</a> . In	1432
1377	Seo Hyun Kim, Keummin Ka, Yohan Jo, Seung-won	<i>The Thirteenth International Conference on Learning</i>	1433
1378	Hwang, Dongha Lee, and Jinyoung Yeo. 2024b.	<i>Representations</i> .	1434
1379	<a href="#">Ever-evolving memory by blending and refining the</a>	Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio	1435
1380	<a href="#">past</a> . <i>arXiv preprint arXiv:2403.04787</i> .	Petroni, Vladimir Karpukhin, Naman Goyal, Hein-	1436
1381	James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz,	rich Küttler, Mike Lewis, Wen-tau Yih, Tim Rock-	1437
1382	Joel Veness, Guillaume Desjardins, Andrei A Rusu,	täschel, and 1 others. 2020. Retrieval-augmented gen-	1438
1383	Kieran Milan, John Quan, Tiago Ramalho, Ag-	eration for knowledge-intensive nlp tasks. <i>Advances</i>	1439
1384	nieszka Grabska-Barwinska, and 1 others. 2017.	<i>in neural information processing systems</i> , 33:9459–	1440
1385	Overcoming catastrophic forgetting in neural net-	9474.	1441
1386	works. <i>Proceedings of the national academy of sci-</i>	Selma Leydesdorff. 2017. <i>Memory cultures: Memory,</i>	1442
1387	<i>ences</i> , 114(13):3521–3526.	<i>subjectivity and recognition</i> . Routledge.	1443
1388	Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying	Hao Li, Chenghao Yang, An Zhang, Yang Deng, Xiang	1444
1389	Sheng, Lianmin Zheng, Cody Hao Yu, Joseph Gon-	Wang, and Tat-Seng Chua. 2024a. Hello again! llm-	1445
1390	zalez, Hao Zhang, and Ion Stoica. 2023. <a href="#">Efficient</a>	powered personalized agent for long-term dialogue.	1446
1391	<a href="#">memory management for large language model serv-</a>	<i>arXiv preprint arXiv:2406.05925</i> .	1447
1392	<a href="#">ing with pagedattention</a> . In <i>Proceedings of the 29th</i>	Haoyang LI, Yiming Li, Anxin Tian, Tianhao Tang,	1448
1393	<i>Symposium on Operating Systems Principles, SOSP</i>	Zhanchao Xu, Xuejia Chen, Nicole HU, Wei Dong,	1449
1394	'23, page 611–626, New York, NY, USA. Association	Li Qing, and Lei Chen. 2025. <a href="#">A survey on large</a>	1450
1395	for Computing Machinery.	<a href="#">language model acceleration based on KV cache</a>	1451
1396	Kunfeng Lai, Zhenheng Tang, Xinglin Pan, Peijie Dong,	<a href="#">management</a> . <i>Transactions on Machine Learning</i>	1452
1397	Xiang Liu, Haolan Chen, Li Shen, Bo Li, and Xi-	<i>Research</i> .	1453
1398	aowen Chu. 2025. Mediator: Memory-efficient llm	K. Li, O. Patel, F. Viégas, H. Pfister, and M. Wattenberg.	1454
1399	merging with less parameter conflicts and uncertainty	2024b. Inference-time intervention: Eliciting truthful	1455
1400	based routing. <i>arXiv preprint arXiv:2502.04411</i> .	answers from a language model. In <i>Advances in</i>	1456
1401	John E Laird. 2019. <i>The Soar cognitive architecture</i> .	<i>Neural Information Processing Systems</i> , volume 36.	1457
1402	MIT press.	Lei Li, Yongfeng Zhang, Dugang Liu, and Li Chen.	1458
1403	Langchain. 2023. Recursively split by char-	2023a. <a href="#">Large language models for generative rec-</a>	1459
1404	acter. <a href="https://python.langchain.com/docs/modules/data_connection/document_transformers/recursive_text_splitter">https://python.langchain.com/</a>	<a href="#">ommendation: A survey and visionary discussions</a> .	1460
1405	<a href="https://python.langchain.com/docs/modules/data_connection/document_transformers/recursive_text_splitter">docs/modules/data_connection/document_</a>	<i>arXiv</i> .	1461
1406	<a href="https://python.langchain.com/docs/modules/data_connection/document_transformers/recursive_text_splitter">transformers/recursive_text_splitter</a> .	Nian Li, Chen Gao, Yong Li, and Qingmin Liao. 2023b.	1462
1407	Gibbeum Lee, Volker Hartmann, Jongho Park, Dim-	<a href="#">Large language model-empowered agents for simu-</a>	1463
1408	itris Papailiopoulos, and Kangwook Lee. 2023a.	<a href="#">lating macroeconomic activities</a> . <i>arXiv</i> .	1464
1409	<a href="#">Prompted llms as chatbot modules for long open-</a>	Shilong Li, Yancheng He, Hangyu Guo, Xingyuan Bu,	1465
1410	<a href="#">domain conversation</a> . In <i>Findings of the Association</i>	Ge Bai, Jie Liu, Jiaheng Liu, Xingwei Qu, Yang-	1466
1411	<i>for Computational Linguistics: ACL 2023, Toronto,</i>	guang Li, Wanli Ouyang, Wenbo Su, and Bo Zheng.	1467
1412	<i>Canada, July 9-14, 2023</i> , pages 4536–4554. Associa-	2024c. <a href="#">GraphReader: Building graph-based agent to</a>	1468
1413	tion for Computational Linguistics.	<a href="#">enhance long-context abilities of large language mod-</a>	1469
1414	Gibbeum Lee, Volker Hartmann, Jongho Park,	<a href="#">els</a> . In <i>Findings of the Association for Computational</i>	1470
1415	Dimitris Papailiopoulos, and Kangwook Lee.	<i>Linguistics: EMNLP 2024</i> , pages 12758–12786, Mi-	1471
1416	2023b. <a href="#">Prompted llms as chatbot modules for</a>	ami, Florida, USA. Association for Computational	1472
1417	<a href="#">long open-domain conversation</a> . <i>arXiv preprint</i>	<i>Linguistics</i> .	1473
1418	<i>arXiv:2305.04533</i> .	Weiqing Li, Guochao Jiang, Xiangyong Ding,	1474
1419	Katherine Lee, Daphne Ippolito, Andrew Nystrom,	Zhangcheng Tao, Chuzhan Hao, Chenfeng Xu,	1475
1420	Chiyuan Zhang, Douglas Eck, Chris Callison-Burch,	Yuewei Zhang, and Hao Wang. 2025a. <a href="#">Flowkv:</a>	1476
1421	and Nicholas Carlini. 2022. <a href="#">Deduplicating training</a>	<a href="#">A disaggregated inference framework with low-</a>	1477
1422	<a href="#">data makes language models better</a> . In <i>Proceedings</i>	<a href="#">latency kv cache transfer and load-aware scheduling</a> .	1478
1423	<a href="#">of the 60th Annual Meeting of the Association for</a>	<i>Preprint</i> , arXiv:2504.03775.	1479

1480	Xianming Li and Jing Li. 2023. Angle-optimized text embeddings. <i>arXiv preprint arXiv:2309.12871</i> .	1536
1481		1537
1482	Xiaoqian Li, Ercong Nie, and Sheng Liang. 2023c. From classification to generation: Insights into crosslingual retrieval augmented icl. <i>arXiv preprint arXiv:2311.06595</i> .	1538
1483		1539
1484		1540
1485		1541
1486	Yinheng Li, Shaofei Wang, Han Ding, and Hang Chen. 2023d. Large language models in finance: A survey. In <i>Proceedings of the Fourth ACM International Conference on AI in Finance</i> , pages 374–382.	1542
1487		1543
1488		1544
1489		1545
1490	Yuan Li, Yixuan Zhang, and Lichao Sun. 2023e. Metaagents: Simulating interactions of human behaviors for llm-based task-oriented coordination via collaborative generative agents. <i>arXiv preprint arXiv:2310.06500</i> .	1546
1491		1547
1492		1548
1493		1549
1494		1550
1495	Yuanchun Li, Hao Wen, Weijun Wang, Xiangyu Li, Yizhen Yuan, Guohong Liu, Jiacheng Liu, Wenxing Xu, Xiang Wang, Yi Sun, and et al. 2024d. Personal llm agents: Insights and survey about the capability, efficiency and security. <i>arXiv</i> .	1551
1496		1552
1497		1553
1498		1554
1499		1555
1500	Yuhong Li, Yingbing Huang, Bowen Yang, Bharat Venkitesh, Acyr Locatelli, Hanchen Ye, Tianle Cai, Patrick Lewis, and Deming Chen. 2024e. Snapkv: Llm knows what you are looking for before generation. <i>arXiv preprint arXiv:2404.14469</i> .	1556
1501		1557
1502		1558
1503		1559
1504		1560
1505	Zeyu Li, Chuanfu Xiao, Yang Wang, Xiang Liu, Zhenheng Tang, Baotong Lu, Mao Yang, Xinyu Chen, and Xiaowen Chu. 2025b. Antkv: Anchor token-aware sub-bit vector quantization for kv cache in large language models. <i>Preprint</i> , arXiv:2506.19505.	1561
1506		1562
1507		1563
1508		1564
1509		1565
1510	Zhiyu Li, Shichao Song, Chenyang Xi, Hanyu Wang, Chen Tang, Simin Niu, Ding Chen, Jiawei Yang, Chunyu Li, Qingchen Yu, and 1 others. 2025c. Memos: A memory os for ai system. <i>arXiv preprint arXiv:2507.03724</i> .	1566
1511		1567
1512		1568
1513		1569
1514		1570
1515	Jianghao Lin, Xinyi Dai, Yunjia Xi, Weiwen Liu, Bo Chen, Xiangyang Li, Chenxu Zhu, Huifeng Guo, Yong Yu, Ruiming Tang, and et al. 2023a. How can recommender systems benefit from large language models: A survey. <i>arXiv</i> .	1571
1516		1572
1517		1573
1518		1574
1519		1575
1520	Xi Victoria Lin, Xilun Chen, Mingda Chen, Weijia Shi, Maria Lomeli, Rich James, Pedro Rodriguez, Jacob Kahn, Gergely Szilvasy, Mike Lewis, and 1 others. 2023b. Ra-dit: Retrieval-augmented dual instruction tuning. <i>arXiv preprint arXiv:2310.01352</i> .	1576
1521		1577
1522		1578
1523		1579
1524		1580
1525	Yujun Lin, Haotian Tang, Shang Yang, Zhekai Zhang, Guangxuan Xiao, Chuang Gan, and Song Han. 2025. Qserve: W4a8kv4 quantization and system co-design for efficient llm serving. <i>Preprint</i> , arXiv:2405.04532.	1581
1526		1582
1527		1583
1528		1584
1529		1585
1530	Akide Liu, Jing Liu, Zizheng Pan, Yefei He, Gholamreza Haffari, and Bohan Zhuang. 2024a. Minicache: Kv cache compression in depth dimension for large language models. In <i>Advances in Neural Information Processing Systems</i> , volume 37, pages 139997–140031. Curran Associates, Inc.	1586
1531		1587
1532		1588
1533		1589
1534		1590
1535		1591
		1592
	Di Liu, Meng Chen, Baotong Lu, Huiqiang Jiang, Zhenhua Han, Qianxi Zhang, Qi Chen, Chengruidong Zhang, Bailu Ding, Kai Zhang, Chen Chen, Fan Yang, Yuqing Yang, and Lili Qiu. 2024b. Retrievalattention: Accelerating long-context llm inference via vector retrieval. <i>Preprint</i> , arXiv:2409.10516.	1536
		1537
		1538
		1539
		1540
		1541
	Jiahong Liu, Zexuan Qiu, Zhongyang Li, Quanyu Dai, Jieming Zhu, Minda Hu, Menglin Yang, and Irwin King. 2025a. A survey of personalized large language models: Progress and future directions. <i>arXiv preprint arXiv:2502.11528</i> .	1542
		1543
		1544
		1545
		1546
	Junyi Liu, Liangzhi Li, Tong Xiang, Bowen Wang, and Yiming Qian. 2023a. TCRA-LLM: Token compression retrieval augmented large language model for inference cost reduction. In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 9796–9810, Singapore. Association for Computational Linguistics.	1547
		1548
		1549
		1550
		1551
		1552
		1553
	Lei Liu, Xiaoyan Yang, Yue Shen, Binbin Hu, Zhiqiang Zhang, Jinjie Gu, and Guannan Zhang. 2023b. Think-in-memory: Recalling and post-thinking enable llms with long-term memory. <i>arXiv preprint arXiv:2311.08719</i> .	1554
		1555
		1556
		1557
		1558
	Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2024c. Lost in the middle: How language models use long contexts. <i>Transactions of the Association for Computational Linguistics</i> , 12:157–173.	1559
		1560
		1561
		1562
		1563
	S. Liu, H. Ye, L. Xing, and J. Zou. 2023c. In-context vectors: Making in context learning more effective and controllable through latent space steering. <i>arXiv</i> .	1564
		1565
		1566
	WenTao Liu, Ruohua Zhang, Aimin Zhou, Feng Gao, and JiaLi Liu. 2025b. Echo: A large language model with temporal episodic memory. <i>arXiv preprint arXiv:2502.16090</i> .	1567
		1568
		1569
		1570
	Xiang Liu, Hong Chen, Xuming Hu, and Xiaowen Chu. 2025c. FlowKV: Enhancing multi-turn conversational coherence in LLMs via isolated key-value cache management. In <i>First Workshop on Multi-Turn Interactions in Large Language Models</i> .	1571
		1572
		1573
		1574
		1575
	Xiang Liu, Zhenheng Tang, Peijie Dong, Zeyu Li, Yue Liu, Bo Li, Xuming Hu, and Xiaowen Chu. 2025d. Chunkkv: Semantic-preserving kv cache compression for efficient long-context llm inference. <i>Preprint</i> , arXiv:2502.00299.	1576
		1577
		1578
		1579
		1580
	Yuhan Liu, Hanchen Li, Yihua Cheng, Siddhant Ray, Yuyang Huang, Qizheng Zhang, Kuntai Du, Jiayi Yao, Shan Lu, Ganesh Ananthanarayanan, Michael Maire, Henry Hoffmann, Ari Holtzman, and Junchen Jiang. 2024d. Cachegen: Kv cache compression and streaming for fast large language model serving. <i>Preprint</i> , arXiv:2310.07240.	1581
		1582
		1583
		1584
		1585
		1586
		1587
	Zhengliang Liu, Aoxiao Zhong, Yiwei Li, Longtao Yang, Chao Ju, Zihao Wu, Chong Ma, Peng Shu, Cheng Chen, Sekeun Kim, and et al. 2023d. Radiology-gpt: A large language model for radiology. <i>arXiv preprint arXiv:2306.08666</i> .	1588
		1589
		1590
		1591
		1592

1593	Zichang Liu, Aditya Desai, Fangshuo Liao, Weitao Wang, Victor Xie, Zhaozhao Xu, Anastasios Kyrillidis, and Anshumali Shrivastava. 2023e. <a href="#">Scissorhands: Exploiting the persistence of importance hypothesis for LLM KV cache compression at test time</a> . In <i>Thirty-seventh Conference on Neural Information Processing Systems</i> .	Xinbei Ma, Yeyun Gong, Pengcheng He, Hai Zhao, and Nan Duan. 2023. <a href="#">Query rewriting for retrieval-augmented large language models</a> . <i>arXiv preprint arXiv:2305.14283</i> .	1650 1651 1652 1653
1596		A. Mack and A. Turner. 2024. <a href="#">Mechanistically eliciting latent behaviors in language models</a> .	1654 1655
1597		Adyasha Maharana, Dong-Ho Lee, Sergey Tulyakov, Mohit Bansal, Francesco Barbieri, and Yuwei Fang. 2024. <a href="#">Evaluating very long-term conversational memory of llm agents</a> . <i>arXiv preprint arXiv:2402.17753</i> .	1656 1657 1658 1659 1660
1598		Daniel Marczak, Bartłomiej Twardowski, Tomasz Trzcinski, and Sebastian Cygert. 2024. <a href="#">Magmax: Leveraging model merging for seamless continual learning</a> . In <i>ECCV</i> .	1661 1662 1663 1664
1599		Ahmed Masry and Amir Hajian. 2024. <a href="#">Longfin: A multimodal document understanding model for long financial domain documents</a> . <i>CoRR</i> , abs/2401.15050.	1665 1666 1667
1600	Zirui Liu, Jiayi Yuan, Hongye Jin, Shaochen Zhong, Zhaozhao Xu, Vladimir Braverman, Beidi Chen, and Xia Hu. 2024e. <a href="#">KIVI: A tuning-free asymmetric 2bit quantization for KV cache</a> . In <i>International Conference on Machine Learning, ICML 2024</i> , pages 32332–32344. PMLR.	Michael S Matena and Colin A Raffel. 2022. <a href="#">Merging models with fisher-weighted averaging</a> . <i>NeurIPS</i> , 35:17703–17716.	1668 1669 1670
1601		Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. 2017. <a href="#">Communication-efficient learning of deep networks from decentralized data</a> . In <i>AISTATS</i> , pages 1273–1282. PMLR.	1671 1672 1673 1674 1675
1602		Sanket Vaibhav Mehta, Jai Gupta, Yi Tay, Mostafa Dehghani, Vinh Q Tran, Jinfeng Rao, Marc Najork, Emma Strubell, and Donald Metzler. 2022. <a href="#">Dsi++: Updating transformer memory with new documents</a> . <i>arXiv preprint arXiv:2212.09744</i> .	1676 1677 1678 1679 1680
1603		mem0ai. 2024. <a href="#">mem0: The memory layer for personalized ai</a> . mem0.ai.	1681 1682
1604		Fanxu Meng, Pingzhi Tang, Xiaojuan Tang, Zengwei Yao, Xing Sun, and Muhan Zhang. 2025. <a href="#">Transmla: Multi-head latent attention is all you need</a> . <i>Preprint</i> , arXiv:2502.07864.	1683 1684 1685 1686
1605		Kevin Meng, Arnab Sen Sharma, Alex J Andonian, Yonatan Belinkov, and David Bau. 2023. <a href="#">Mass-editing memory in a transformer</a> . In <i>The Eleventh International Conference on Learning Representations</i> .	1687 1688 1689 1690 1691
1606	Anton Lozhkov, Raymond Li, Loubna Ben Allal, Federico Cassano, Joel Lamy-Poirier, Nouamane Tazi, Ao Tang, Dmytro Pykhtar, Jiawei Liu, Yuxiang Wei, Tianyang Liu, Max Tian, Denis Kocetkov, Arthur Zucker, Younes Belkada, Zijian Wang, Qian Liu, Dmitry Abulkhanov, Indraneil Paul, and 38 others. 2024. <a href="#">Starcoder 2 and the stack v2: The next generation</a> . <i>CoRR</i> , abs/2402.19173.	Sewon Min, Xinxi Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2022. <a href="#">Rethinking the role of demonstrations: What makes in-context learning work?</a> In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 11048–11064, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	1692 1693 1694 1695 1696 1697 1698 1699
1607		Mayank Mishra, Matt Stallone, Gaoyuan Zhang, Yikang Shen, Aditya Prasad, Adriana Meza Soria, Michele Merler, Parameswaran Selvam, Saptha Surendran, Shivdeep Singh, Manish Sethi, Xuan-Hong Dang,	1700 1701 1702 1703
1608			
1609			
1610			
1611			
1612			
1613			
1614	Junru Lu, Siyu An, Mingbao Lin, Gabriele Pergola, Yulan He, Di Yin, Xing Sun, and Yunsheng Wu. 2023a. <a href="#">Memochat: Tuning llms to use memos for consistent long-range open-domain conversation</a> . <i>arXiv preprint arXiv:2308.08239</i> .		
1615			
1616			
1617			
1618			
1619	Junru Lu, Siyu An, Mingbao Lin, Gabriele Pergola, Yulan He, Di Yin, Xing Sun, and Yunsheng Wu. 2023b. <a href="#">Memochat: Tuning llms to use memos for consistent long-range open-domain conversation</a> . <i>arXiv preprint arXiv:2308.08239</i> .		
1620			
1621			
1622			
1623			
1624	Zhenyi Lu, Chenghao Fan, Wei Wei, Xiaoye Qu, Danyang Chen, and Yu Cheng. 2024a. <a href="#">Twin-merging: Dynamic integration of modular expertise in model merging</a> . <i>arXiv preprint arXiv:2406.15479</i> .		
1625			
1626			
1627			
1628	Zhenyi Lu, Chenghao Fan, Wei Wei, Xiaoye Qu, Danyang Chen, and Yu Cheng. 2024b. <a href="#">Twin-merging: Dynamic integration of modular expertise in model merging</a> . <i>NIPS</i> .		
1629			
1630			
1631			
1632	Ziyang Luo, Can Xu, Pu Zhao, Xiubo Geng, Chongyang Tao, Jing Ma, Qingwei Lin, and Daxin Jiang. 2023. <a href="#">Augmented large language models with parametric knowledge guiding</a> . <i>arXiv preprint arXiv:2305.04757</i> .		
1633			
1634			
1635			
1636			
1637	Shi Luohe, Hongyi Zhang, Yao Yao, Zuchao Li, and 1 others. <a href="#">Keep the cost down: A review on methods to optimize llm’s kv-cache consumption</a> . In <i>First Conference on Language Modeling</i> .		
1638			
1639			
1640			
1641	Chenyang Lyu, Zefeng Du, Jitao Xu, Yitao Duan, Minghao Wu, Teresa Lynn, Alham Fikri Aji, Derek F. Wong, and Longyue Wang. 2024. <a href="#">A paradigm shift: The future of machine translation lies with large language models</a> . In <i>Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation, LREC/COLING 2024, 20-25 May, 2024, Torino, Italy</i> , pages 1339–1352. ELRA and ICCL.		
1642			
1643			
1644			
1645			
1646			
1647			
1648			
1649			

1704	Pengyuan Li, Kun-Lung Wu, Syed Zawad, Andrew Coleman, Matthew White, Mark Lewis, Raju Pavuluri, and 27 others. 2024. <a href="#">Granite code models: A family of open foundation models for code intelligence</a> . <i>CoRR</i> , abs/2405.04324.	Jane Pan, Tianyu Gao, Howard Chen, and Danqi Chen. 2023b. <a href="#">What in-context learning “learns” in-context: Disentangling task recognition and task learning</a> . In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 8298–8319, Toronto, Canada. Association for Computational Linguistics.	1756
1705			1757
1706			1758
1707			1759
1708			1760
1709	Ali Modarressi, Ayyoob Imani, Mohsen Fayyaz, and Hinrich Schütze. 2023. <a href="#">Ret-llm: Towards a general read-write memory for large language models</a> . <i>arXiv</i> .	Zhuoshi Pan, Qianhui Wu, Huiqiang Jiang, Xufang Luo, Hao Cheng, Dongsheng Li, Yuqing Yang, Chin-Yew Lin, H Vicky Zhao, Lili Qiu, and et al. 2025. On memory construction and retrieval for personalized conversational agents. <i>arXiv preprint arXiv:2502.05589</i> .	1762
1710			1763
1711			1764
1712	Ali Montazerlghaem, Hamed Zamani, and James Allan. 2020. A reinforcement learning framework for relevance feedback. In <i>Proceedings of the 43rd international acm sigir conference on research and development in information retrieval</i> , pages 59–68.	Joon Sung Park, Joseph O’Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. 2023. Generative agents: Interactive simulacra of human behavior. In <i>Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology</i> , pages 1–22.	1765
1713			1766
1714			1767
1715			1768
1716			1769
1717	Mohammed Muqeeth, Haokun Liu, and Colin Raffel. 2024. Soft merging of experts with adaptive routing. <i>TMLR</i> .	Wenjun Peng, Guiyang Li, Yue Jiang, Zilong Wang, Dan Ou, Xiaoyi Zeng, Enhong Chen, and 1 others. 2023. Large language model based long-tail query rewriting in taobao search. <i>arXiv preprint arXiv:2311.03758</i> .	1770
1718			1771
1719			1772
1720	Jaap MJ Murre and Joeri Dros. 2015. Replication and analysis of ebbinghaus’ forgetting curve. <i>PloS one</i> , 10(7):e0120644.	Reiner Pope, Sholto Douglas, Aakanksha Chowdhery, Jacob Devlin, James Bradbury, Anselm Levskaya, Jonathan Heek, Kefan Xiao, Shivani Agrawal, and Jeff Dean. 2022. <a href="#">Efficiently scaling transformer inference</a> . <i>Preprint</i> , arXiv:2211.05102.	1773
1721			1774
1722			1775
1723	Milad Nasr, Nicholas Carlini, Jonathan Hayase, Matthew Jagielski, A. Feder Cooper, Daphne Ippolito, Christopher A. Choquette-Choo, Eric Wallace, Florian Tramèr, and Katherine Lee. 2023. <a href="#">Scalable extraction of training data from (production) language models</a> . <i>Preprint</i> , arXiv:2311.17035.	Chen Qian, Xin Cong, Cheng Yang, Weize Chen, Yusheng Su, Juyuan Xu, Zhiyuan Liu, and Maosong Sun. 2023. Communicative agents for software development. <i>arXiv preprint arXiv:2307.07924</i> .	1776
1724			1777
1725			1778
1726			1779
1727			1780
1728			1781
1729	Yuqi Nie, Yaxuan Kong, Xiaowen Dong, John M. Mulvey, H. Vincent Poor, Qingsong Wen, and Stefan Zohren. 2024. <a href="#">A survey of large language models for financial applications: Progress, prospects and challenges</a> . <i>CoRR</i> , abs/2406.11903.	Hongjin Qian, Peitian Zhang, Zheng Liu, Kelong Mao, and Zhicheng Dou. 2024. Memorag: Moving towards next-gen rag via memory-inspired knowledge discovery. <i>arXiv preprint arXiv:2409.05591</i> .	1782
1730			1783
1731			1784
1732			1785
1733			1786
1734	Yuxiang Nie, Heyan Huang, Wei Wei, and Xian-Ling Mao. 2022. <a href="#">Capturing global structural information in long document question answering with compressive graph selector network</a> . In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 5036–5047, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	Zhangcheng Qiang, Weiqing Wang, and Kerry Taylor. 2023. <a href="#">Agent-om: Leveraging large language models for ontology matching</a> . <i>arXiv</i> .	1787
1735			1788
1736			1789
1737			1790
1738			1791
1739			1792
1740			1793
1741			1794
1742	OpenAI. 2024a. <a href="#">Memory and new controls for chatgpt</a> . Accessed: 2024-02-13.	Zhe Qu, Xingyu Li, Rui Duan, Yao Liu, Bo Tang, and Zhuo Lu. 2022. Generalized federated learning via sharpness aware minimization. In <i>International Conference on Machine Learning</i> , pages 18250–18280. PMLR.	1795
1743			1796
1744	OpenAI. 2024b. <a href="#">New embedding models and api updates</a> . Accessed: 2024-01-25.	Qwen, :, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, and 1 others. 2025. <a href="#">Qwen2.5 technical report</a> . <i>Preprint</i> , arXiv:2412.15115.	1797
1745			1798
1746	Charles Packer, Sarah Wooders, Kevin Lin, Vivian Fang, Shishir G Patil, Ion Stoica, and Joseph E Gonzalez. 2023. Memgpt: Towards llms as operating systems. <i>arXiv preprint arXiv:2310.08560</i> .	Allan Raventos, Mansheej Paul, Feng Chen, and Surya Ganguli. 2023. Pretraining task diversity and the emergence of non-bayesian in-context learning for regression. In <i>Thirty-seventh Conference on Neural Information Processing Systems</i> .	1799
1747			1800
1748			1801
1749			1802
1750	Haojie Pan, Zepeng Zhai, Hao Yuan, Yaojia Lv, Ruiji Fu, Ming Liu, Zhongyuan Wang, and Bing Qin. 2023a. <a href="#">Kwaiagents: Generalized information-seeking agent system with large language models</a> . <i>arXiv</i> .		1803
1751			1804
1752			1805
1753			1806
1754	James Pan and Guoliang Li. 2025. <a href="#">A survey of llm inference systems</a> . <i>Preprint</i> , arXiv:2506.21901.		1807
1755			1807

1808	Varshini Reddy, Rik Koncel-Kedziorski, Viet Dac Lai,	Ying Sheng, Lianmin Zheng, Binhang Yuan, Zhuohan	1862
1809	Michael Krumdick, Charles Lovering, and Chris Tan-	Li, Max Ryabinin, Beidi Chen, Percy Liang, Christo-	1863
1810	ner. 2024. <a href="#">Docfinqa: A long-context financial reason-</a>	pher Ré, Ion Stoica, and Ce Zhang. 2023. Flexgen:	1864
1811	<a href="#">ing dataset</a> . In <i>Proceedings of the 62nd Annual Meet-</i>	high-throughput generative inference of large lan-	1865
1812	<i>ing of the Association for Computational Linguistics,</i>	guage models with a single gpu. In <i>Proceedings of</i>	1866
1813	<i>ACL 2024 - Short Papers, Bangkok, Thailand, Au-</i>	<i>the 40th International Conference on Machine Learn-</i>	1867
1814	<i>gust 11-16, 2024</i> , pages 445–458. Association for	<i>ing, ICML'23</i> . JMLR.org.	1868
1815	Computational Linguistics.		
1816	Jon Saad-Falcon, Daniel Y. Fu, Simran Arora, Neel	Lauralee Sherwood, Robert Thomas Kell, and Christo-	1869
1817	Guha, and Christopher Ré. 2024. <a href="#">Benchmarking and</a>	pher Ward. 2004. <i>Human physiology: from cells to</i>	1870
1818	<a href="#">building long-context retrieval models with loco and</a>	<i>systems</i> . Thomson/Brooks/Cole.	1871
1819	<a href="#">M2-BERT</a> . In <i>Forty-first International Conference</i>		
1820	<i>on Machine Learning, ICML 2024, Vienna, Austria,</i>	Freda Shi, Xinyun Chen, Kanishka Misra, Nathan	1872
1821	<i>July 21-27, 2024</i> . OpenReview.net.	Scales, David Dohan, Ed H Chi, Nathanael Schärli,	1873
1822	Ali Safaya and Deniz Yuret. 2024. Neurocache: Effi-	and Denny Zhou. 2023. Large language models can	1874
1823	cient vector retrieval for long-range language model-	be easily distracted by irrelevant context. In <i>Inter-</i>	1875
1824	ing. <i>arXiv preprint arXiv:2407.02486</i> .	<i>national Conference on Machine Learning</i> , pages	1876
1825	Rana Salama, Jason Cai, Michelle Yuan, Anna Currey,	31210–31227. PMLR.	1877
1826	Monica Sunkara, Yi Zhang, and Yassine Benajiba.		
1827	2025. Meminsight: Autonomous memory augmenta-	Kaize Shi, Xueyao Sun, Qing Li, and Guandong Xu.	1878
1828	tion for llm agents. <i>arXiv preprint arXiv:2503.21760</i> .	2024. <a href="#">Compressing long context for enhancing</a>	1879
1829	Mohammad Reza Samsami, Artem Zholus, Janarthanan	<a href="#">RAG with amr-based concept distillation</a> . <i>CoRR</i> ,	1880
1830	Rajendran, and Sarath Chandar. 2024. <a href="#">Mastering</a>	abs/2405.03085.	1881
1831	<a href="#">memory tasks with world models</a> . In <i>The Twelfth In-</i>	Noah Shinn, Federico Cassano, Ashwin Gopinath,	1882
1832	<i>ternational Conference on Learning Representations</i> .	Karthik Narasimhan, and Shunyu Yao. 2024. Re-	1883
1833	Utkarsh Saxena, Gobinda Saha, Sakshi Choudhary, and	flexion: Language agents with verbal reinforcement	1884
1834	Kaushik Roy. 2024. <a href="#">Eigen attention: Attention in</a>	learning. <i>Advances in Neural Information Process-</i>	1885
1835	<a href="#">low-rank space for KV cache compression</a> . In <i>Find-</i>	<i>ing Systems</i> , 36.	1886
1836	<i>ings of the Association for Computational Linguistics:</i>	Noah Shinn, Federico Cassano, Ashwin Gopinath,	1887
1837	<i>EMNLP 2024</i> , pages 15332–15344, Miami, Florida,	Karthik R. Narasimhan, and Shunyu Yao. 2023. Re-	1888
1838	USA. Association for Computational Linguistics.	flexion: Language agents with verbal reinforcement	1889
1839	Lianlei Shan, Shixian Luo, Zezhou Zhu, Yu Yuan, and	learning. In <i>Thirty-seventh Conference on Neural</i>	1890
1840	Yong Wu. 2025. Cognitive memory in large language	<i>Information Processing Systems</i> .	1891
1841	models. <i>arXiv preprint arXiv:2504.02441</i> .	Haseeb Ullah Khan Shinwari and Muhammad Usama.	1892
1842	Jingbo Shang, Zai Zheng, Jiale Wei, Xiang Ying, Felix	2025. Memory-augmented architecture for long-term	1893
1843	Tao, and Mindverse Team. 2024. Ai-native memory:	context handling in large language models. <i>arXiv</i>	1894
1844	A pathway from llms towards agi. <i>arXiv preprint</i>	<i>preprint arXiv:2506.18271</i> .	1895
1845	<i>arXiv:2406.18312</i> .	Aaditya K Singh, Stephanie C.Y. Chan, Ted Moskovitz,	1896
1846	Bin Shao and Jiawei Yan. 2024. A long-context	Erin Grant, Andrew M Saxe, and Felix Hill. 2023.	1897
1847	language model for deciphering and generating	The transient nature of emergent in-context learning	1898
1848	bacteriophage genomes. <i>Nature Communications</i> ,	in transformers. In <i>Thirty-seventh Conference on</i>	1899
1849	15(1):9392.	<i>Neural Information Processing Systems</i> .	1900
1850	Yunfan Shao, Linyang Li, Junqi Dai, and Xipeng Qiu.	Robert L Solso and Jerome Kagan. 1979. <i>Cognitive</i>	1901
1851	Character-llm: A trainable agent for role-playing.	<i>psychology</i> . Houghton Mifflin Harcourt P.	1902
1852	Yunfan Shao, Linyang Li, Junqi Dai, and Xipeng Qiu.	L. R. Squire. 2009. Memory and brain systems: 1969-	1903
1853	2023. Character-llm: A trainable agent for role-	2009. <i>j neurosci</i> . 2009 oct 14;29(41):12711-6. doi:	1904
1854	playing. <i>arXiv preprint arXiv:2310.10158</i> .	10.1523/jneurosci.3575-09.2009. pmid: 19828780;	1905
1855	Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz,	pmcid: Pmc2791502. <i>J Neurosci</i> , 29(41):12711–	1906
1856	Andy Davis, Quoc Le, Geoffrey Hinton, and Jeff	12716.	1907
1857	Dean. 2017. <a href="#">Outrageously large neural networks:</a>	S Sridhar, A Khamaj, and MK Asthana. 2023. Cogni-	1908
1858	<a href="#">The sparsely-gated mixture-of-experts layer</a> .	tive neuroscience perspective on memory: overview	1909
1859	Lingfeng Shen, Aayush Mishra, and Daniel Khashabi.	and summary. <i>Frontiers in human neuroscience</i> ,	1910
1860	2024. Do pretrained transformers learn in-context by	17:1217093.	1911
1861	gradient descent? <i>arxiv</i> , arXiv:2310.08540.	N. Subramani, N. Suresh, and M. E. Peters. 2022. <a href="#">Ex-</a>	1912
		<a href="#">tracting latent steering vectors from pretrained lan-</a>	1913
		<a href="#">guage models</a> . <i>arXiv</i> .	1914

1915 Hao Sun, Hengyi Cai, Bo Wang, Yingyan Hou, Xiaochi  
1916 Wei, Shuaiqiang Wang, Yan Zhang, and Dawei Yin.  
1917 2024. [Towards verifiable text generation with evolving memory and self-reflection](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 8211–8227, Miami, Florida, USA. Association for Computational Linguistics.

1923 Ron Sun. 2001. *Duality of the mind: A bottom-up approach toward cognition*. Psychology Press.

1925 Zhiqing Sun, Xuezhi Wang, Yi Tay, Yiming Yang, and  
1926 Denny Zhou. 2022. Recitation-augmented language  
1927 models. *arXiv preprint arXiv:2210.01296*.

1928 Richard S Sutton and Andrew G Barto. 2018. *Reinforcement learning: An introduction*. MIT press.

1930 Xin Tan, Yimin Jiang, Yitao Yang, and Hong Xu. 2025.  
1931 [Teola: Towards end-to-end optimization of llm-based applications](#). *Preprint*, arXiv:2407.00326.

1933 Anke Tang, Li Shen, Yong Luo, Nan Yin, Lefei Zhang,  
1934 and Dacheng Tao. 2024a. Merging multi-task models  
1935 via weight-ensembling mixture of experts. *ICML*.

1936 Jiaming Tang, Yilong Zhao, Kan Zhu, Guangxuan  
1937 Xiao, Baris Kasikci, and Song Han. 2024b. QUEST:  
1938 Query-aware sparsity for efficient long-context LLM  
1939 inference. In *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pages 47901–47911. PMLR.

1943 Ravi Teja. 2023. Evaluating the ideal  
1944 chunk size for a rag system using llamaindex. <https://www.llamaindex.ai/blog/evaluating-the-ideal-chunk-size-for-a-rag-system-using-llamaindex>. 6207ef5d3f665

1945 ~~https://www.llamaindex.ai/blog/evaluating-the-ideal-chunk-size-for-a-rag-system-using-llamaindex~~

1946

1947 Kushal Tirumala, Aram H. Markosyan, Luke Zettle-  
1948 moyer, and Armen Aghajanyan. 2022. [Memo- rization without overfitting: Analyzing the training dynamics of large language models](#). *Preprint*, arXiv:2205.10770.

1952 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier  
1953 Martinet, Marie-Anne Lachaux, Timothée Lacroix,  
1954 Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal  
1955 Azhar, Aurelien Rodriguez, Armand Joulin, Edouard  
1956 Grave, and Guillaume Lample. 2023a. [Llama: Open and efficient foundation language models](#). *Preprint*, arXiv:2302.13971.

1959 Hugo Touvron, Louis Martin, Kevin Stone, and 1 others.  
1960 2023b. [Llama 2: Open foundation and fine-tuned chat models](#). *Preprint*, arXiv:2307.09288.

1961

1962 YunDa Tsai, Mingjie Liu, and Haoxing Ren. 2023. [Rtl- fixer: Automatically fixing rtl syntax errors with large language models](#). *arXiv preprint arXiv:2311.16543*.

1963

1964

1965 E. Tulving and W. Donaldson. 1972. *Episodic and semantic memory*. Academic Press.

1966

1967 A. Turner, M. Kurzeja, D. Orr, and D. Elson. 2025.  
1968 [Steering gemini using bidpo vectors](#).

Alexander Matt Turner, Lisa Thiergart, Gavin Leech, David Udell, Juan J Vazquez, Ulisse Mini, and Monte MacDiarmid. 2023. Steering language models with activation engineering. *arXiv preprint arXiv:2308.10248*.

Szymon Tworkowski, Konrad Staniszewski, Mikołaj Pacek, Yuhuai Wu, Henryk Michalewski, and Piotr Miłoś. 2023. [Focused transformer: Contrastive training for context scaling](#). In *Advances in Neural Information Processing Systems*, volume 36, pages 42661–42688. Curran Associates, Inc.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Proceedings of the 31st International Conference on Neural Information Processing Systems, NIPS’17*, page 6000–6010, Red Hook, NY, USA. Curran Associates Inc.

VoyageAI. 2023. Voyage’s embedding models. <https://docs.voyageai.com/embeddings/>.

Luanbo Wan and Weizhi Ma. 2025. Storybench: A dynamic benchmark for evaluating long-term memory with multi turns. *arXiv preprint arXiv:2506.13356*.

Benyou Wang, Qianqian Xie, Jiahuan Pei, Zhihong Chen, Prayag Tiwari, Zhao Li, and Jie Fu. 2023a. Pre-trained language models in biomedical domain: A systematic survey. *ACM Computing Surveys*, 56(3):1–52.

Bing Wang, Xinnian Liang, Jian Yang, Hui Huang, Shuangzhi Wu, Peihao Wu, Lu Lu, Zejun Ma, and Zhongjun Li. 2024. [Enhancing large language model with self-controlled memory framework](#). *Preprint*, arXiv:2304.13343.

Chen Wang, Xunzhuo Liu, Yuhuan Liu, Yue Zhu, Xi-angxi Mo, Junchen Jiang, and Huamin Chen. 2025a. When to reason: Semantic router for vllm. *arXiv preprint arXiv:2510.08731*.

2002

2003

2004

2005

2006 Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Man-  
2007 dlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and An-  
2008 ima Anandkumar. 2023b. [Voyager: An open-ended embodied agent with large language models](#). *arXiv preprint arXiv:2305.16291*.

2009

2010

2011 Haochun Wang, Chi Liu, Nuwa Xi, Zewen Qiang,  
2012 Sendong Zhao, Bing Qin, and Ting Liu. 2023c. Hu-  
2013 atuo: Tuning llama model with chinese medical  
2014 knowledge. *arXiv preprint arXiv:2304.06975*.

2015 Haochun Wang, Sendong Zhao, Zewen Qiang, Zi-  
2016 jian Li, Nuwa Xi, Yanrui Du, MuZhen Cai, Hao-  
2017 qiang Guo, Yuhuan Chen, Haoming Xu, and et al.  
2018 2023d. [Knowledge-tuning large language models with structured medical knowledge bases for reliable response generation in chinese](#). *arXiv preprint arXiv:2309.04175*.

2019

2020

2021

2022	Junjie Wang, Yuchao Huang, Chunyang Chen, Zhe Liu,	Qingyue Wang, Yanhe Fu, Yanan Cao, Shuai Wang,	2078
2023	Song Wang, and Qing Wang. 2024b. Software testing	Zhiliang Tian, and Liang Ding. 2025d. <a href="#">Recur-</a>	2079
2024	with large language models: Survey, landscape, and	sively summarizing enables long-term dialogue mem-	2080
2025	vision. <i>IEEE Transactions on Software Engineering</i> .	ory in large language models. <i>Neurocomputing</i> ,	2081
		639:130193.	2082
2026	Ke Wang, Nikolaos Dimitriadis, Guillermo Ortiz-	Song Wang, Yaochen Zhu, Haochen Liu, Zaiyi Zheng,	2083
2027	Jimenez, François Fleuret, and Pascal Frossard.	Chen Chen, and Jundong Li. 2024h. Knowledge	2084
2028	2024c. Localizing task information for improved	editing for large language models: A survey. <i>ACM</i>	2085
2029	model merging and compression. <i>ICML</i> .	<i>Computing Surveys</i> , 57(3):1–37.	2086
2030	Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao	Weizhi Wang, Li Dong, Hao Cheng, Xiaodong Liu,	2087
2031	Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang,	Xifeng Yan, Jianfeng Gao, and Furu Wei. 2023h.	2088
2032	Xu Chen, Yankai Lin, and et al. 2023e. <a href="#">A survey</a>	<a href="#">Augmenting language models with long-term mem-</a>	2089
2033	<a href="#">on large language model based autonomous agents.</a>	<a href="#">ory</a> . In <i>Advances in Neural Information Processing</i>	2090
2034	<i>arxiv</i> .	<i>Systems 36: Annual Conference on Neural Informa-</i>	2091
2035	Lei Wang, Jingsen Zhang, Hao Yang, Zhiyuan Chen,	<i>tion Processing Systems 2023, NeurIPS 2023, New</i>	2092
2036	Jiakai Tang, Zeyu Zhang, Xu Chen, Yankai Lin, Rui-	<i>Orleans, LA, USA, December 10 - 16, 2023</i> .	2093
2037	hua Song, Wayne Xin Zhao, Jun Xu, Zhicheng Dou,	Wenjie Wang, Xinyu Lin, Fuli Feng, Xiangnan He, and	2094
2038	Jun Wang, and Ji-Rong Wen. 2023f. When large lan-	Tat-Seng Chua. 2023i. <a href="#">Generative recommendation:</a>	2095
2039	guage model based agent meets user behavior analy-	<a href="#">Towards next-generation recommender paradigm.</a>	2096
2040	sis: A novel user simulation paradigm.	<i>arXiv</i> .	2097
2041	Liang Wang, Nan Yang, Xiaolong Huang, Linjun Yang,	Xindi Wang, Mahsa Salmani, Parsa Omidi, Xiangyu	2098
2042	Rangan Majumder, and Furu Wei. 2024d. <a href="#">Improv-</a>	Ren, Mehdi Rezagholizadeh, and Armaghan Eshaghi.	2099
2043	<a href="#">ing text embeddings with large language models</a> .	2024i. <a href="#">Beyond the limits: A survey of techniques to</a>	2100
2044	In <i>Proceedings of the 62nd Annual Meeting of the As-</i>	<a href="#">extend the context length in large language models</a> .	2101
2045	<i>sociation for Computational Linguistics (Volume 1:</i>	In <i>Proceedings of the Thirty-Third International Joint</i>	2102
2046	<i>Long Papers)</i> , ACL 2024, Bangkok, Thailand, August	<i>Conference on Artificial Intelligence, IJCAI 2024,</i>	2103
2047	11-16, 2024, pages 11897–11916. Association for	<i>Jeju, South Korea, August 3-9, 2024</i> , pages 8299–	2104
2048	Computational Linguistics.	8307. <i>ijcai.org</i> .	2105
2049	Liang Wang, Nan Yang, and Furu Wei. 2023g. Learning	Xintao Wang, Qianwen Yang, Yongting Qiu, Jiaqing	2106
2050	to retrieve in-context examples for large language	Liang, Qianyu He, Zhouhong Gu, Yanghua Xiao,	2107
2051	models. <i>arXiv preprint arXiv:2307.07164</i> .	and Wei Wang. 2023j. Knowledgept: Enhancing large	2108
2052	Liyuan Wang, Xingxing Zhang, Hang Su, and Jun Zhu.	language models with retrieval and storage access on	2109
2053	2024e. A comprehensive survey of continual learn-	knowledge bases. <i>arXiv preprint arXiv:2308.11761</i> .	2110
2054	ing: Theory, method and application. <i>IEEE Transac-</i>	Yile Wang, Peng Li, Maosong Sun, and Yang Liu.	2111
2055	<i>tions on Pattern Analysis and Machine Intelligence</i> .	2023k. Self-knowledge guided retrieval augmen-	2112
2056	Longyue Wang, Zefeng Du, Wenxiang Jiao, Chenyang	tation for large language models. <i>arXiv preprint</i>	2113
2057	Lyu, Jianhui Pang, Leyang Cui, Kaiqiang Song,	<i>arXiv:2310.05002</i> .	2114
2058	Derek F. Wong, Shuming Shi, and Zhaopeng Tu.	Yu Wang, Yifan Gao, Xiusi Chen, Haoming Jiang,	2115
2059	2024f. <a href="#">Benchmarking and improving long-text trans-</a>	Shiyang Li, Jingfeng Yang, Qingyu Yin, Zheng Li,	2116
2060	<a href="#">lation with large language models</a> . In <i>Findings of</i>	Xian Li, Bing Yin, and et al. 2024j. Memoryllm:	2117
2061	<i>the Association for Computational Linguistics, ACL</i>	<a href="#">Towards self-updatable large language models.</a>	2118
2062	<i>2024, Bangkok, Thailand and virtual meeting, Au-</i>	<i>preprint arXiv:2402.04624</i> .	2119
2063	<i>gust 11-16, 2024</i> , pages 7175–7187. Association for	Yu Wang, Chi Han, Tongtong Wu, Xiaoxin He,	2120
2064	Computational Linguistics.	Wangchunshu Zhou, Nafis Sadeq, Xiusi Chen, Zexue	2121
2065	Peng Wang, Zexi Li, Ningyu Zhang, Ziwen Xu, Yunzhi	He, Wei Wang, Gholamreza Haffari, and 1 others.	2122
2066	Yao, Yong Jiang, Pengjun Xie, Fei Huang, and Hua-	2024k. Towards lifespan cognitive systems. <i>arXiv</i>	2123
2067	jun Chen. 2024g. Wise: Rethinking the knowledge	<i>preprint arXiv:2409.13265</i> .	2124
2068	memory for lifelong model editing of large language	Yu Wang, Nedim Lipka, Ryan A Rossi, Alexa Siu, Ruiyi	2125
2069	models. <i>Advances in Neural Information Processing</i>	Zhang, and Tyler Derr. 2023l. Knowledge graph	2126
2070	<i>Systems</i> , 37:53764–53797.	prompting for multi-document question answering.	2127
2071	Qian Wang, Zhenheng Tang, and Bingsheng He. 2025b.	<i>arXiv preprint arXiv:2308.11730</i> .	2128
2072	Can LLM simulations truly reflect humanity? a deep	Yu Wang, Xinshuang Liu, Xiusi Chen, Sean O’Brien,	2129
2073	dive. In <i>The Fourth Blogpost Track at ICLR 2025</i> .	Junda Wu, and Julian McAuley. Self-updatable large	2130
2074	Qingyue Wang, Yanan Fu, Yanan Cao, Shi Wang, Zhil-	language models by integrating context into model	2131
2075	iang Tian, and Liang Ding. 2025c. <a href="#">Recursively sum-</a>	parameters. In <i>The Thirteenth International Confer-</i>	2132
2076	<a href="#">marizing enables long-term dialogue memory in large</a>	<i>ence on Learning Representations</i> .	2133
2077	<a href="#">language models</a> . <i>Neurocomputing</i> , page 130193.		

2134	Zefan Wang, Zichuan Liu, Yingying Zhang, Aoxiao	Yichen Wu, Hong Wang, Peilin Zhao, Yefeng Zheng,	2188
2135	Zhong, Lunting Fan, Lingfei Wu, and Qingsong Wen.	Ying Wei, and Long-Kai Huang. 2024b. Mitigating	2189
2136	2023m. <a href="#">Rcagent: Cloud root cause analysis by au-</a>	catastrophic forgetting in online continual learning	2190
2137	tonomous agents with tool-augmented large language	by modeling previous task interrelations via pareto	2191
2138	models. <i>arXiv</i> .	optimization. In <i>Forty-first International Conference</i>	2192
		<i>on Machine Learning</i> .	2193
2139	Zhepeng Wang, Runxue Bao, Yawen Wu, Jackson Tay-	Yuhuai Wu, Markus N Rabe, DeLesley Hutchins, and	2194
2140	lor, Cao Xiao, Feng Zheng, Weiwen Jiang, Shangqian	Christian Szegedy. 2022. Memorizing transformers.	2195
2141	Gao, and Yanfu Zhang. 2024l. <a href="#">Unlocking memo-</a>	<i>arXiv preprint arXiv:2203.08913</i> .	2196
2142	ization in large language models with dynamic soft		
2143	<a href="#">prompting</a> . In <i>Proceedings of the 2024 Conference</i>	Wm. A. Wulf and Sally A. McKee. 1995. <a href="#">Hitting</a>	2197
2144	<i>on Empirical Methods in Natural Language Process-</i>	<a href="#">the memory wall: implications of the obvious</a> .	2198
2145	<i>ing</i> , pages 9782–9796, Miami, Florida, USA. Associ-	<i>SIGARCH Comput. Archit. News</i> , 23(1):20–24.	2199
2146	ation for Computational Linguistics.		
2147	Zora Zhiruo Wang, Jiayuan Mao, Daniel Fried, and	Yunjia Xi, Weiwen Liu, Jianghao Lin, Bo Chen, Ruim-	2200
2148	Graham Neubig. 2024m. <a href="#">Agent workflow memory</a> .	Tang, Weinan Zhang, and Yong Yu. 2024. Mem-	2201
2149	<i>Preprint</i> , arXiv:2409.07429.	ocrs: Memory-enhanced sequential conversational	2202
		recommender systems with large language models.	2203
2150	Zora Zhiruo Wang, Jiayuan Mao, Daniel Fried, and	In <i>Proceedings of the 33rd ACM International Con-</i>	2204
2151	Graham Neubig. 2024n. <a href="#">Agent workflow memory</a> .	<i>ference on Information and Knowledge Management</i> ,	2205
2152	<i>arXiv preprint arXiv:2409.07429</i> .	pages 2585–2595.	2206
2153	Fanjunduo Wei, Zhenheng Tang, Rongfei Zeng,	Zhiheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwen	2207
2154	Tongliang Liu, Chengqi Zhang, Xiaowen Chu, and	Ding, Boyang Hong, Ming Zhang, Junzhe Wang,	2208
2155	Bo Han. 2025. <a href="#">JailbreakLoRA: Your downloaded</a>	Senjie Jin, Enyu Zhou, and et al. 2023. <a href="#">The rise and</a>	2209
2156	<a href="#">IoRA from sharing platforms might be unsafe</a> . In	<a href="#">potential of large language model based agents: A</a>	2210
2157	<i>ICML 2025 Workshop on Data in Generative Models</i>	<a href="#">survey</a> . <i>arxiv</i> .	2211
2158	- <i>The Bad, the Ugly, and the Greats</i> .		
2159	Lilian Weng. 2023. Llm-powered autonomous agents.	Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song	2212
2160	lilianweng.github.io.	Han, and Mike Lewis. 2023. Efficient streaming	2213
		language models with attention sinks. <i>arXiv</i> .	2214
2161	Di Wu, Hongwei Wang, Wenhao Yu, Yuwei Zhang, Kai-	Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song	2215
2162	Wei Chang, and Dong Yu. 2024a. Longmemeval:	Han, and Mike Lewis. 2024a. <a href="#">Efficient streaming</a>	2216
2163	Benchmarking chat assistants on long-term interac-	<a href="#">language models with attention sinks</a> . In <i>The Twelfth</i>	2217
2164	tive memory. <i>arXiv preprint arXiv:2410.10813</i> .	<i>International Conference on Learning Representa-</i>	2218
		<i>tions</i> .	2219
2165	Haoyi Wu and Kewei Tu. 2024. <a href="#">Layer-condensed KV</a>	Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song	2220
2166	<a href="#">cache for efficient inference of large language mod-</a>	Han, and Mike Lewis. 2024b. <a href="#">Efficient stream-</a>	2221
2167	<a href="#">els</a> . In <i>Proceedings of the 62nd Annual Meeting of</i>	<a href="#">ing language models with attention sinks</a> . <i>Preprint</i> ,	2222
2168	<i>the Association for Computational Linguistics (Vol-</i>	arXiv:2309.17453.	2223
2169	<i>ume 1: Long Papers)</i> , pages 11175–11188, Bangkok,		
2170	Thailand. Association for Computational Linguistics.	Honglin Xiong, Sheng Wang, Yitao Zhu, Zihao Zhao,	2224
		Yuxiao Liu, Qian Wang, and Dinggang Shen. 2023.	2225
2171	Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu,	Doctorglm: Fine-tuning your chinese doctor is not a	2226
2172	Shaokun Zhang, Erkang Zhu, Beibin Li, Li Jiang,	herculean task. <i>arXiv preprint arXiv:2304.01097</i> .	2227
2173	Xiaoyun Zhang, and Chi Wang. 2023. Auto-		
2174	gen: Enabling next-gen llm applications via multi-	Derong Xu, Wei Chen, Wenjun Peng, Chao Zhang, Tong	2228
2175	agent conversation framework. <i>arXiv preprint</i>	Xu, Xiangyu Zhao, Xian Wu, Yefeng Zheng, and	2229
2176	<i>arXiv:2308.08155</i> .	Enhong Chen. 2023. <a href="#">Large language models for</a>	2230
		<a href="#">generative information extraction: A survey</a> . <i>arXiv</i> .	2231
2177	Wei Wu, Zhuoshi Pan, Chao Wang, Liyi Chen, Yunchu	Jing Xu, Arthur Szlam, and Jason Weston. 2021. Be-	2232
2178	Bai, Tianfu Wang, Kun Fu, Zheng Wang, and Hui	yond goldfish memory: Long-term open-domain con-	2233
2179	Xiong. 2025a. <a href="#">Tokenselect: Efficient long-context</a>	versation. <i>arXiv preprint arXiv:2107.07567</i> .	2234
2180	<a href="#">inference and length extrapolation for llms via dy-</a>		
2181	<a href="#">namic token-level kv cache selection</a> . <i>Preprint</i> ,	Wujiang Xu, Zujie Liang, Kai Mei, Hang Gao, Jun-	2235
2182	arXiv:2411.02886.	tao Tan, and Yongfeng Zhang. 2025. <a href="#">A-mem:</a>	2236
		<a href="#">Agentic memory for llm agents</a> . <i>arXiv preprint</i>	2237
2183	Yaxiong Wu, Sheng Liang, Chen Zhang, Yichao Wang,	arXiv:2502.12110.	2238
2184	Yongyue Zhang, Huifeng Guo, Ruiming Tang, and		
2185	Yong Liu. 2025b. <a href="#">From human memory to ai mem-</a>	Prateek Yadav, Derek Tam, Leshem Choshen, Colin	2239
2186	<a href="#">ory: A survey on memory mechanisms in the era of</a>	Raffel, and Mohit Bansal. 2023. <a href="#">Ties-merging: Re-</a>	2240
2187	<a href="#">llms</a> . <i>Preprint</i> , arXiv:2504.15965.	<a href="#">solving interference when merging models</a> .	2241

2242	Shi-Qi Yan, Jia-Chen Gu, Yun Zhu, and Zhen-Hua Ling.	Weiran Yao, Shelby Heinecke, Juan Carlos Niebles,	2296
2243	2024. Corrective retrieval augmented generation.	Zhiwei Liu, Yihao Feng, Le Xue, Rithesh Murthy,	2297
2244	<i>arXiv preprint arXiv:2401.15884</i> .	Zeyuan Chen, Jianguo Zhang, Devansh Arpit, and	2298
		et al. 2023. Retroformer: Retrospective large	2299
2245	Dongjie Yang, Xiaodong Han, Yan Gao, Yao Hu, Shilin	language agents with policy gradient optimization.	2300
2246	Zhang, and Hai Zhao. 2024a. PyramidInfer: Pyramid	<i>arXiv preprint arXiv:2308.02151</i> .	2301
2247	KV cache compression for high-throughput LLM in-		
2248	ference. In <i>Findings of the Association for Computa-</i>	Yao Yao, Zuchao Li, and Hai Zhao. 2024. Sir-	2302
2249	<i>tational Linguistics: ACL 2024</i> , pages 3258–3270,	llm: Streaming infinite retentive llm. <i>Preprint</i> ,	2303
2250	Bangkok, Thailand. Association for Computational	arXiv:2405.12528.	2304
2251	Linguistics.		
2252	Enneng Yang, Zhenyi Wang, Li Shen, Shiwei Liu, Guib-	Lu Ye, Ze Tao, Yong Huang, and Yang Li. 2024.	2305
2253	ing Guo, Xingwei Wang, and Dacheng Tao. 2024b.	Chunkattention: Efficient self-attention with prefix-	2306
2254	Adamerging: Adaptive model merging for multi-task	aware kv cache and two-phase partition. <i>Preprint</i> ,	2307
2255	learning. <i>ICLR</i> .	arXiv:2402.15220.	2308
2256	Hongkang Yang, Zehao Lin, Wenjin Wang, Hao Wu,	Gyeong-In Yu, Joo Seong Jeong, Geon-Woo Kim, Soo-	2309
2257	Zhiyu Li, Bo Tang, Wenqiang Wei, Jinbo Wang,	jeong Kim, and Byung-Gon Chun. 2022a. Orca: A	2310
2258	Zeyun Tang, Shichao Song, and 1 others. 2024c.	distributed serving system for Transformer-Based	2311
2259	<i>memory</i> <sup>3</sup> : Language modeling with explicit mem-	generative models. In <i>16th USENIX Symposium</i>	2312
2260	ory. <i>arXiv preprint arXiv:2407.01178</i> .	<i>on Operating Systems Design and Implementation</i>	2313
		<i>(OSDI 22)</i> , pages 521–538, Carlsbad, CA. USENIX	2314
		Association.	2315
2261	Huan Yang, Renji Zhang, Mingzhe Huang, Weijun		
2262	Wang, Yin Tang, Yuanchun Li, Yunxin Liu, and Deyu	Le Yu, Bowen Yu, Haiyang Yu, Fei Huang, and Yongbin	2316
2263	Zhang. 2025a. Kvshare: An llm service system with	Li. 2024a. Language models are super mario: Ab-	2317
2264	efficient and effective multi-tenant kv cache reuse.	sorbing abilities from homologous models as a free	2318
2265	<i>Preprint</i> , arXiv:2503.16525.	lunch. <i>ICML</i> .	2319
2266	Jingbo Yang, Bairu Hou, Wei Wei, Yujia Bao, and	Le Yu, Bowen Yu, Haiyang Yu, Fei Huang, and Yongbin	2320
2267	Shiyu Chang. 2025b. Kvlink: Accelerating large lan-	Li. 2024b. Language models are super mario: Ab-	2321
2268	guage models via efficient kv cache reuse. <i>Preprint</i> ,	sorbing abilities from homologous models as a free	2322
2269	arXiv:2502.16002.	lunch. <i>ICML</i> .	2323
2270	June Yong Yang, Byeongwook Kim, Jeongin Bae,	Wenhao Yu, Dan Iter, Shuohang Wang, Yichong Xu,	2324
2271	Beomseok Kwon, Gunho Park, Eunho Yang, Se Jung	Mingxuan Ju, Soumya Sanyal, Chenguang Zhu,	2325
2272	Kwon, and Dongsoo Lee. 2024d. No token left be-	Michael Zeng, and Meng Jiang. 2022b. Gen-	2326
2273	hind: Reliable kv cache compression via importance-	erate rather than retrieve: Large language mod-	2327
2274	aware mixed precision quantization. <i>Preprint</i> ,	els are strong context generators. <i>arXiv preprint</i>	2328
2275	arXiv:2402.18096.	arXiv:2209.10063.	2329
2276	Ling Yang, Zhaochen Yu, Tianjun Zhang, Shiyi Cao,	Wenhao Yu, Hongming Zhang, Xiaoman Pan, Kaixin	2330
2277	Minkai Xu, Wentao Zhang, Joseph E Gonzalez,	Ma, Hongwei Wang, and Dong Yu. 2023. Chain-of-	2331
2278	and Bin Cui. 2024e. Buffer of thoughts: Thought-	note: Enhancing robustness in retrieval-augmented	2332
2279	augmented reasoning with large language models.	language models. <i>arXiv preprint arXiv:2311.09210</i> .	2333
2280	<i>arXiv preprint arXiv:2406.04271</i> .		
2281	Sophia Yang. 2023. Advanced rag 01: Small-to-big	Li Yunxiang, Li Zihan, Zhang Kai, Dan Ruilong, and	2334
2282	retrieval. <a href="https://towardsdatascience.com/advanced-rag-01-small-to-big-retrieval-172181b39644">https://towardsdatascience.com/</a>	Zhang You. 2023. Chatdoctor: A medical chat model	2335
2283	<a href="https://towardsdatascience.com/advanced-rag-01-small-to-big-retrieval-172181b39644">advanced-rag-01-small-to-big-retrieval-172181b39644</a>	fine-tuned on llama model using medical domain	2336
		knowledge. <i>arXiv preprint arXiv:2303.14070</i> .	2337
2284	Zhenjie Yang, Xiaosong Jia, Hongyang Li, and Junchi	Fanlong Zeng, Wensheng Gan, Yongheng Wang, Ning	2338
2285	Yan. 2023. A survey of large language models for	Liu, and Philip S Yu. 2023. Large language models	2339
2286	autonomous driving. <i>arXiv</i> .	for robotics: A survey. <i>arXiv</i> .	2340
2287	Jiayi Yao, Hanchen Li, Yuhan Liu, Siddhant Ray, Yihua	Ruihong Zeng, Jinyuan Fang, Siwei Liu, and Zaiqiao	2341
2288	Cheng, Qizheng Zhang, Kuntai Du, Shan Lu, and	Meng. 2024. On the structural memory of llm agents.	2342
2289	Junchen Jiang. 2025. Cacheblend: Fast large lan-	<i>arXiv preprint arXiv:2412.15266</i> .	2343
2290	guage model serving for rag with cached knowledge		
2291	fusion. <i>Preprint</i> , arXiv:2405.16444.	Liangyu Zha, Junlin Zhou, Liyao Li, Rui Wang,	2344
		Qingyi Huang, Saisai Yang, Jing Yuan, Chang-	2345
2292	Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak	bao Su, Xiang Li, Aofeng Su, and 1 others. 2023.	2346
2293	Shafran, Karthik Narasimhan, and Yuan Cao. 2022.	Tablegpt: Towards unifying tables, nature lan-	2347
2294	React: Synergizing reasoning and acting in language	guage and commands into one gpt. <i>arXiv preprint</i>	2348
2295	models. <i>arXiv preprint arXiv:2210.03629</i> .	arXiv:2307.08674.	2349



2459	language models with long-term memory. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 38, pages 19724–19731.	2515
2460		2516
2461		2517
2462	Wanjun Zhong, Lianghong Guo, Qiqi Gao, He Ye, and Yanlin Wang. 2024c. <a href="#">Memorybank: Enhancing large language models with long-term memory</a> . In <i>Thirty-Eighth AAAI Conference on Artificial Intelligence, AAAI 2024, Thirty-Sixth Conference on Innovative Applications of Artificial Intelligence, IAAI 2024, Fourteenth Symposium on Educational Advances in Artificial Intelligence, EAAI 2014, February 20-27, 2024, Vancouver, Canada</i> , pages 19724–19731. AAAI Press.	2518
2463		2519
2464		2520
2465		2521
2466		2522
2467		2523
2468		2524
2469		2525
2470		2526
2471		2527
2472	Yinmin Zhong, Shengyu Liu, Junda Chen, Jianbo Hu, Yibo Zhu, Xuanzhe Liu, Xin Jin, and Hao Zhang. 2024d. <a href="#">Distserve: disaggregating prefill and decoding for goodput-optimized large language model serving</a> . In <i>Proceedings of the 18th USENIX Conference on Operating Systems Design and Implementation, OSDI'24, USA</i> . USENIX Association.	2528
2473		2529
2474		2530
2475		2531
2476		2532
2477		2533
2478		2534
2479	Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans, Claire Cui, Olivier Bousquet, Quoc Le, and Ed Chi. 2023a. <a href="#">Least-to-most prompting enables complex reasoning in large language models</a> . <i>Preprint</i> , arXiv:2205.10625.	2535
2480		2536
2481		2537
2482		2538
2483		2539
2484		2540
2485	Hongjian Zhou, Boyang Gu, Xinyu Zou, Yiru Li, Sam S Chen, Peilin Zhou, Junling Liu, Yining Hua, Chengfeng Mao, Xian Wu, and et al. 2023b. <a href="#">A survey of large language models in medicine: Progress, application, and challenge</a> . <i>arXiv</i> .	2541
2486		2542
2487		2543
2488		2544
2489		2545
2490	Xiabin Zhou, Wenbin Wang, Minyan Zeng, Jiaxian Guo, Xuebo Liu, Li Shen, Min Zhang, and Liang Ding. 2025. <a href="#">Dynamickv: Task-aware adaptive kv cache compression for long context llms</a> . <i>Preprint</i> , arXiv:2412.14838.	2546
2491		
2492		
2493		
2494		
2495	Zixuan Zhou, Xuefei Ning, Ke Hong, Tianyu Fu, Jiaming Xu, Shiyao Li, Yuming Lou, Luning Wang, Zhihang Yuan, Xiuhong Li, and 1 others. 2024. <a href="#">A survey on efficient inference for large language models</a> . <i>arXiv preprint arXiv:2404.14294</i> .	
2496		
2497		
2498		
2499		
2500	Dawei Zhu, Liang Wang, Nan Yang, Yifan Song, Wenhao Wu, Furu Wei, and Sujian Li. 2024. <a href="#">Longembed: Extending embedding models for long context retrieval</a> . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, EMNLP 2024, Miami, FL, USA, November 12-16, 2024</i> , pages 802–816. Association for Computational Linguistics.	
2501		
2502		
2503		
2504		
2505		
2506		
2507		
2508	Qiuyu Zhu, Liang Zhang, Qianxiong Xu, Cheng Long, and Jie Zhang. 2025a. <a href="#">Subgcache: Accelerating graph-based rag with subgraph-level kv cache</a> . <i>Preprint</i> , arXiv:2505.10951.	
2509		
2510		
2511		
2512	Xizhou Zhu, Yuntao Chen, Hao Tian, Chenxin Tao, Weijie Su, Chenyu Yang, Gao Huang, Bin Li, Lewei Lu, Xiaogang Wang, and et al. 2023a. <a href="#">Ghost in the minecraft: Generally capable agents for open-world environments via large language models with text-based knowledge and memory</a> . <i>arXiv preprint arXiv:2305.17144</i> .	
2513		
2514		
	Xizhou Zhu, Yuntao Chen, Hao Tian, Chenxin Tao, Weijie Su, Chenyu Yang, Gao Huang, Bin Li, Lewei Lu, Xiaogang Wang, and et al. 2023b. <a href="#">Ghost in the minecraft: Generally capable agents for open-world environments via large language models with text-based knowledge and memory</a> . <i>arXiv preprint arXiv:2305.17144</i> .	
	Yuanbing Zhu, Zhenheng Tang, Xiang Liu, Ang Li, Bo Li, Xiaowen Chu, and Bo Han. 2025b. <a href="#">OracleKV: Oracle guidance for question-independent KV cache compression</a> . In <i>ICML 2025 Workshop on Long-Context Foundation Models</i> .	
	Yutao Zhu, Huaying Yuan, Shuting Wang, Jiongnan Liu, Wenhan Liu, Chenlong Deng, Zhicheng Dou, and Ji-Rong Wen. 2023c. <a href="#">Large language models for information retrieval: A survey</a> . <i>arXiv</i> .	
	Yutao Zhu, Huaying Yuan, Shuting Wang, Jiongnan Liu, Wenhan Liu, Chenlong Deng, Zhicheng Dou, and Ji-Rong Wen. 2023d. <a href="#">Large language models for information retrieval: A survey</a> . <i>CoRR</i> , abs/2308.07107.	
	Yuxuan Zhu, Ali Falahati, David H. Yang, and Mohammad Mohammadi Amiri. 2025c. <a href="#">Sentencekv: Efficient llm inference via sentence-level semantic kv caching</a> . <i>Preprint</i> , arXiv:2504.00970.	
	A. Zou, L. Phan, S. Chen, J. Campbell, P. Guo, R. Ren, A. Pan, X. Yin, M. Mazeika, A.-K. Dombrowski, and et al. 2023. <a href="#">Representation engineering: A top-down approach to ai transparency</a> . <i>arXiv</i> .	

## A Related Surveys

The rapid development of LLMs has triggered many surveys on LLM-based agents, which study how an LLM can perceive, act, accumulate knowledge, and adapt over time. Early reviews (e.g., (Wang et al., 2023e)) organize agent research by how agents are built, where they are used, and how they are evaluated. Later surveys expand the scope with different taxonomies and emphases (Xi et al., 2023; Zhao et al., 2023b; Cheng et al., 2024; Ge et al., 2023). There are also focused reviews on key capabilities and settings, such as multimodal agents (Durante et al., 2024), planning (Huang et al., 2024), multi-agent interaction (Guo et al., 2024), and personal assistant applications (Li et al., 2024d). These works provide useful summaries of agent pipelines, but memory is usually treated as one module among many, and is rarely analyzed as a first-class system component with a unified interface and lifecycle.

A separate line of surveys summarizes how LLMs are applied to specific domains. In information retrieval and extraction, surveys cover LLM-based query processing (Zhu et al., 2023c) and taxonomies for information extraction (Xu et al., 2023). In recommender systems, several reviews discuss how LLMs and agent-style components are used for data generation and recommendation (Li et al., 2023a; Lin et al., 2023a; Wang et al., 2023i). In software engineering, surveys summarize the use of LLMs across design, development, and testing (Fan et al., 2023; Wang et al., 2024b; Zheng et al., 2023b). Other domain surveys cover robotics (Zeng et al., 2023), autonomous driving (Cui et al., 2024; Yang et al., 2023), medicine (He et al., 2023a; Zhou et al., 2023b; Wang et al., 2023a), finance (Li et al., 2023d), and psychology (He et al., 2023b). While these surveys are valuable for understanding domain adaptation, they typically treat memory as domain-specific prompting or retrieval practice, rather than a general representation and management problem.

Surveys that target memory in LLMs and agent systems are more closely related to our work, but the current picture is still fragmented. Some surveys discuss operational aspects of memory (Zhang et al., 2024c), yet many narrow their scope to long-context modeling (Huang et al., 2023b), long-term memory (He et al., 2024d; Jiang et al., 2024a), personalization (Liu et al., 2025a), or knowledge

editing (Wang et al., 2024h). This topical split makes it hard to compare methods across settings, and it often blurs the boundary between (i) what is stored (the memory representation) and (ii) how it is used and maintained (the memory management). As a result, practical foundations such as consistent benchmarks, tools, and implementation constraints are not discussed in a unified way.

Several recent surveys propose alternative lenses for understanding memory. Some move beyond a pure time-based split (short-term vs. long-term) and categorize memory by the memory “object”, such as personal memories for user interaction and system memories for internal state (Jiang et al., 2024a; Zhang et al., 2024c; Zhong et al., 2024b). Others focus on memory mechanisms inside LLM-based agents and review their design, evaluation, and applications for self-evolving behaviors (Zhang et al., 2025b). Another direction decomposes memory into smaller operations and separates parametric and contextual forms, listing operations such as updating, indexing, retrieval, and compression (Du et al., 2025). Human-memory-inspired surveys further relate human memory categories to AI memory designs and propose multi-dimensional categorizations (Wu et al., 2025b). Empirical studies evaluate how memory structures and retrieval strategies affect agent performance, including how memory addition and deletion influence long-horizon behaviors (Zeng et al., 2024). These perspectives are informative, but they often treat the operation list as the primary organizing principle, which does not directly expose the system constraints behind different memory backends.

A complementary line of surveys studies memory from the deployment and efficiency angle. Inference system surveys summarize how to deliver high throughput and quality under large workloads (Pan and Li, 2025), and KV-cache management is reviewed as a key technique for reducing redundant computation and improving memory use during decoding (Hatalis et al., 2024; LI et al., 2025; Luohe et al.). Broader inference optimization surveys also analyze sources of inefficiency and summarize techniques at the data, model, and system levels (Zhou et al., 2024). These works are mainly organized around performance techniques. They provide less discussion on how runtime memory (e.g., KV cache) relates to other memory forms under a single representation-and-management view, and how different backends can be composed in one agent system. A related sur-

vey (Shan et al., 2025) uses a taxonomy close to ours, but it does not provide a systematic discussion of implementation details across memory backends.

In contrast, our survey treats memory as a separable system component and organizes prior work by two orthogonal axes: the representation used to store memory and the management process that constructs, updates, and queries that memory. This decouples memory mechanisms from specific learning modes such as in-context learning or weight updates, and clarifies how the same management goals can be realized through different backends (token memory, intermediate latent memory, and parametric memory) under different cost and reliability constraints.

## B Overview of Human and LLM Memory and Taxonomy

This section provides a high-level overview of human and LLM memory through a concise taxonomy of both. By decoupling human memory from LLM memory, we analyze how different categories of human memory can be instantiated using distinct LLM memory mechanisms. Building on this perspective, we present a holistic framework that demystifies LLM memory design and offers a unified intuition for understanding approaches.

### B.1 Human Memory

Human memory is a complex and multifaceted phenomenon, recognized in cognitive neuroscience as a collection of interconnected processes, including encoding, consolidation, storage, and retrieval (Baddeley and Hitch, 1974; Sridhar et al., 2023). It represents the brain’s remarkable capacity to store, retain, and recall information, serving as the foundation for learning, adapting to environments, and shaping personal identity (Sherwood et al., 2004; Weng, 2023). Memory underpins higher-order cognitive functions such as reasoning, problem-solving, and language comprehension, profoundly influencing behavior and decision-making (Budson and Kensinger, 2023; Shan et al., 2025).

Based on the *duration of information retention*, human memory is classified into **short-term** and **long-term memory**, as illustrated in Figure 2(a) (Baddeley, 2007). Short-term memory temporarily holds and processes information for seconds to minutes as working memory, which actively manipulates information for immediate tasks like reason-

ing and comprehension (Budson and Kensinger, 2023; Baddeley and Hitch, 1974). In contrast, long-term memory stores information for extended periods, ranging from minutes to years, forming a repository for enduring knowledge and experiences (Budson and Kensinger, 2023).

Based on functional roles, human memory is commonly categorized into **explicit** (declarative), **implicit** (non-declarative), and **sensory memory**<sup>1</sup>, as shown in Figure 2(b) (Budson and Kensinger, 2023). Explicit memory involves conscious recall of facts and events that can be readily articulated. It includes episodic memory, which captures personal experiences tied to specific times and contexts (e.g., recalling what one ate for lunch) (Tulving and Donaldson, 1972), and semantic memory, which stores factual knowledge independent of personal experience (e.g., knowing the Earth is round) (Begg, 1984). Implicit memory operates unconsciously and is harder to verbalize, including procedural memory that governs skills and habits acquired through repetition, such as riding a bicycle or playing a musical instrument (Squire, 2009). These memory systems interact to support information processing and storage.

From a cognitive psychology perspective, memory is a fundamental mental process critical to learning and behavior (Solso and Kagan, 1979). It enables the accumulation of knowledge, abstraction of high-level concepts, and formation of social norms through the retention of cultural values and personal experiences ( Craik and Lockhart, 1972; Leydesdorff, 2017). Memory also supports decision-making by allowing individuals to anticipate potential consequences (Johnson-Laird, 1983). These insights are invaluable for designing LLM-based agents, as memory modules that mirror human cognitive processes enhance their ability to perform complex tasks and exhibit human-like behavior (Laird, 2019; Sun, 2001).

Memory is essential for the self-evolution of LLM-based agents in dynamic environments (Sutton and Barto, 2018). It facilitates experience accumulation, enabling agents to retain past errors, inappropriate behaviors, or failed attempts to improve future performance and learning efficiency (Zheng et al., 2023a). Memory also supports environment exploration by guiding agents to prioritize less-explored actions or revisit previously unsuc-

<sup>1</sup>For LLMs, we do not discuss sensory memory in detail, as LLMs primarily operate on text.

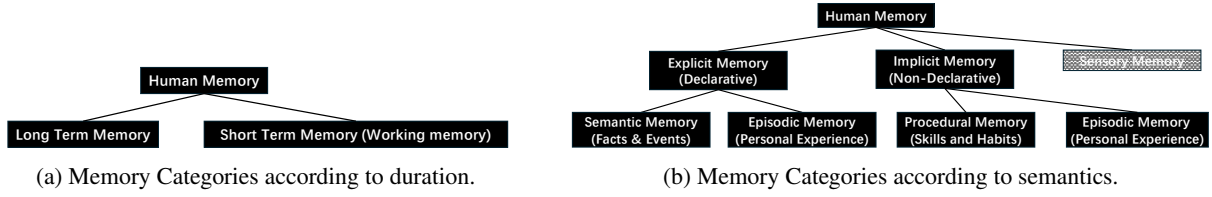


Figure 2: Human Memory Overview.

2748 successful trials, enhancing adaptability (Montazerl-  
 2749 ghaem et al., 2020; Zhu et al., 2023a). Additionally,  
 2750 memory enables knowledge abstraction, allowing  
 2751 agents to summarize raw observations into high-  
 2752 level insights, which is crucial for generalizing to  
 2753 new environments (Zhao et al., 2023a).

## 2754 B.2 LLM Memory

2755 **LLM Inference.** LLMs such as GPT (Brown  
 2756 et al., 2020b), LLaMA (Touvron et al., 2023a,b;  
 2757 Grattafiori et al., 2024), Qwen (Qwen et al., 2025;  
 2758 Bai et al., 2023), and DeepSeek (DeepSeek-AI  
 2759 et al., 2025a,b) operate under the autoregressive  
 2760 generation paradigm. In this approach, the model  
 2761 predicts the next token based on all previously gen-  
 2762 erated tokens. Given an input sequence of tokens  
 2763  $(x_1, \dots, x_n)$ , the model computes a probability dis-  
 2764 tribution over the vocabulary for the next token at  
 2765 each time step  $t$ , typically using the final token’s  
 2766 representation. This process is governed by the  
 2767 joint probability:

$$2768 P(x_1, \dots, x_n) = P(x_1) \times P(x_2 | x_1) \\
 2769 \times \dots \times P(x_n | x_1, \dots, x_{n-1}). \quad (1)$$

2770 The inference process (Equation 1) relies on the  
 2771 self-attention mechanism within the transformer  
 2772 architecture. For each token  $i$ , self-attention com-  
 2773 puts a weighted sum over the representations of  
 2774 all previous tokens  $\{1, 2, \dots, i\}$ , resulting in a time  
 2775 complexity of  $\mathcal{O}(n^2)$  for a sequence of length  $n$ .  
 2776 This quadratic complexity becomes computationally  
 2777 expensive for long sequences, as attention is  
 2778 recalculated over all preceding tokens at every step.

2779 To mitigate this inefficiency, the KV cache is  
 2780 widely used as an optimization technique. The  
 2781 KV cache divides inference into two phases: **pre-  
 2782 fill** and **decoding**. In the prefill phase, the model  
 2783 processes the entire prompt with full-sequence at-  
 2784 tention, computing and storing the key and value  
 2785 vectors for all tokens. During the decoding phase,  
 2786 as the model generates one token at a time, only the  
 key and value vectors for the new token are com-

2787 puted and appended to the cache. Attention is then  
 2788 calculated solely between the current query and the  
 2789 cached KV pairs, eliminating redundant computa-  
 2790 tions and significantly improving the efficiency of  
 2791 autoregressive generation.

2792 **Transformer Architecture.** The remarkable  
 2793 performance of transformer-based LLMs across di-  
 2794 verse tasks is primarily driven by the self-attention  
 2795 mechanism (Leviathan et al., 2025; Meng et al.,  
 2796 2025; Vaswani et al., 2017). In this mechanism,  
 2797 a sequence of input hidden states  $(h_1, \dots, h_n)$  is  
 2798 transformed through a linear projection layer to  
 2799 produce the Query ( $Q$ ), Key ( $K$ ), and Value ( $V$ )  
 2800 vectors, as defined in Equation 2.

$$2801 \text{concat}(q_i, k_i, v_i) = \text{concat}(W_q, W_k, W_v) \cdot h_i \quad (2)$$

2802 Subsequently, attention scores  $a_{ij}$  are computed  
 2803 by taking the dot product between a Query vector  
 2804  $(q_i)$  and a Key vector  $(k_j)$ , scaled by the square root  
 2805 of the dimension  $d$ . These scores are normalized  
 2806 and used as weights to sum the Value vectors  $(v_j)$ ,  
 2807 producing the output  $(o_i)$  as shown in:

$$2808 a_{ij} = \frac{\exp(q_i^\top k_j / \sqrt{d})}{\sum_{t=1}^i \exp(q_i^\top k_t / \sqrt{d})}, \quad o_i = \sum_{j=1}^i a_{ij} v_j. \quad (3)$$

2809 The output of the self-attention mechanism is  
 2810 then processed by a Feed Forward Network (FFN),  
 2811 which typically consists of multiple linear layers in-  
 2812 terleaved with activation functions, further enhanc-  
 2813 ing the model’s ability to capture complex patterns.

2814 **In Context Learning.** In-context learning en-  
 2815 ables LLMs to retrieve and utilize memory by in-  
 2816 corporating relevant information directly within  
 2817 the input prompt as natural language, without mod-  
 2818 ifying model parameters or storing intermediate  
 2819 representations explicitly. The model relies on the  
 2820 provided context—such as examples, instructions,  
 2821 or retrieved data from a pool or database—to guide  
 2822 its predictions. For a prompt with a sequence of  
 2823 tokens  $(x_1, \dots, x_n)$ , the self-attention mechanism,

Table 2: Characteristics of Human Memory.

Memory Type	Key Function/Characteristics	Duration/Capacity
Sensory Memory	Brief buffer for incoming sensory information (visual, auditory, etc.)	Milliseconds to a few seconds
Working Memory (WM)	Transient active store for manipulating information; supports complex cognitive operations (reasoning, language)	Tens of seconds to minutes; limited items
Short-Term Memory (STM)	Temporary holding of information before transfer to LTM or forgetting	Tens of seconds to minutes; limited items
Long-Term Memory (LTM)	Stores information for extended periods; large capacity and durability	Minutes to decades; Vast capacity
Declarative (Explicit)	Consciously recalled facts and events	Minutes to decades; Vast capacity
Episodic Memory	Personal experiences, specific events with contextual details	Minutes to decades
Semantic Memory	General world knowledge, facts, concepts, language	Minutes to decades
Non-Declarative (Implicit)	Unconscious learning: skills, habits, priming, conditioning	Acquired slowly, long-lasting

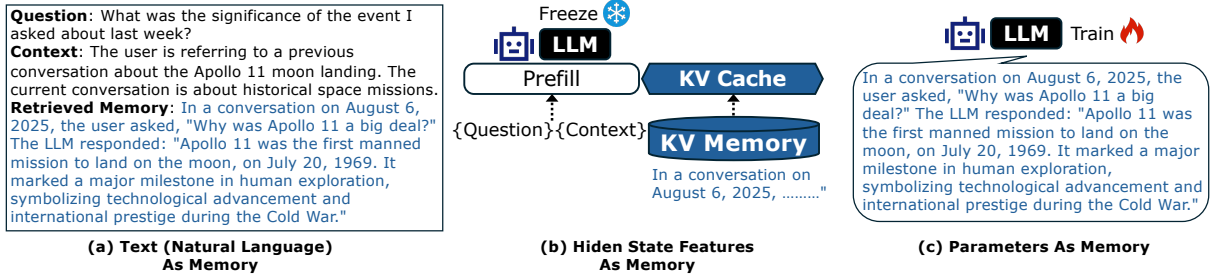


Figure 3: LLM Memory Classification.

as defined in Equation (3), weighs the relevance of each token in the context when predicting the next token, effectively mimicking memory retrieval. Research has shown that in-context learning can be viewed as implicit structure induction or meta-optimization, where the model learns task-specific patterns during inference by performing a form of gradient descent within the attention mechanism (Dai et al., 2023; Hahn and Goyal, 2023; Garg et al., 2022).

In practice, in-context learning is implemented by constructing prompts with relevant examples or retrieved documents, as seen in models like GPT (Brown et al., 2020b) or Qwen (Qwen et al., 2025; Bai et al., 2023). For instance, few-shot learning scenarios provide question-answer pairs to guide responses, with the self-attention mechanism computing scores  $a_{ij} = \frac{\exp(q_i^T k_j / \sqrt{d})}{\sum_{t=1}^i \exp(q_i^T k_t / \sqrt{d})}$  to focus on relevant context tokens. Studies suggest that in-context learning excels in task recognition and adaptation, but its effectiveness depends on the quality and structure of the provided demonstrations (Min et al., 2022; Pan et al., 2023b). Additionally, emergent in-context learning capabilities may be transient and tied to pretraining task diversity, highlighting the role of training data in enabling this mechanism (Singh et al., 2023; Raventos et al., 2023; Shen et al., 2024).

#### KV Cache or Intermediate Representations.

The KV cache stores intermediate Key ( $K$ ) and Value ( $V$ ) vectors computed during the self-attention process, as defined in Equation (2), to enhance efficiency in autoregressive generation. During the **prefill phase**, the model processes the entire prompt, generating and caching  $K$  and  $V$  vectors for all tokens. In the **decoding phase**, only the  $K$  and  $V$  vectors for each newly generated token are computed and appended to the cache, with attention calculated solely between the current Query and cached KV pairs, producing the output  $o_i = \sum_{j=1}^i a_{ij} v_j$ . This reduces the time complexity from  $\mathcal{O}(n^2)$  to  $\mathcal{O}(n)$  per token during decoding, as used in models like LLaMA (Touvron et al., 2023a,b; Grattafiori et al., 2024) and DeepSeek (DeepSeek-AI et al., 2025a,b).

Recent advancements have focused on optimizing the KV cache to handle long-context scenarios and reduce memory overhead. Techniques like CacheGen and ChunkKV employ semantic-preserving compression to reduce the cache size while maintaining performance (Liu et al., 2024d, 2025d; Zhu et al., 2025b; Liu et al., 2025c). KIVI introduces asymmetric 2-bit quantization for KV cache entries, further reducing memory footprint (Liu et al., 2024e). Methods like SnapKV and H2O selectively retain important KV pairs based on attention patterns, improving efficiency for long-context inference (Li et al., 2024e; Zhang et al.,

2882 2023b). Additionally, StreamingLLM uses atten- 2933  
2883 tion sinks to stabilize attention distributions, en- 2934  
2884 abling efficient streaming inference (Xiao et al., 2935  
2885 2023). Distributed approaches, such as KVDirect 2936  
2886 and FlowKV, optimize KV cache storage and trans- 2937  
2887 fer across multi-GPU systems (Chen et al., 2024c; 2938  
2888 Li et al., 2025a). These optimizations make the KV 2939  
2889 cache a critical memory mechanism for scalable 2940  
2890 and efficient LLM inference. 2941

### 2891 **Training-Based Knowledge Integration.**

2892 Training-based memorization embeds knowledge 2942  
2893 directly into the model’s parameters during 2943  
2894 training, effectively transforming data into a 2944  
2895 compressed, implicit memory. LLMs like GPT, 2945  
2896 LLaMA, Qwen, and DeepSeek are trained on vast 2946  
2897 datasets, encoding patterns, facts, and relation- 2947  
2898 ships within weights, particularly in the linear 2948  
2899 projection layers ( $W_q, W_k, W_v$ ) and the Feed 2949  
2900 Forward Network (FFN). The training process 2950  
2901 optimizes the joint probability  $P(x_1, \dots, x_n) =$  2951  
2902  $P(x_1) \cdot P(x_2 | x_1) \cdot \dots \cdot P(x_n | x_1, \dots, x_{n-1})$ , 2952  
2903 adjusting parameters via backpropagation to 2953  
2904 capture statistical and semantic patterns. This 2954  
2905 enables the model to recall knowledge during 2955  
2906 inference without external storage, though the 2956  
2907 memory is static unless fine-tuned or retrained. 2957

2908 **Summary.** As shown in Figure 3, In-context 2958  
2909 learning provides flexible memory through natural 2959  
2910 language prompts, leveraging self-attention to adap- 2960  
2911 tively retrieve task-specific information, with its ef- 2961  
2912 ficacy tied to demonstration quality and pretraining 2962  
2913 diversity (Min et al., 2022; Raventos et al., 2023). 2963  
2914 Examples of implementing memory with different 2964  
2915 kinds of natural language tokens are shown in Ta- 2965  
2916 ble 3. The KV cache optimizes inference by storing 2966  
2917 intermediate representations, with recent advance- 2967  
2918 ments like compression and selective retention en- 2968  
2919 hancing efficiency for long contexts (Liu et al., 2969  
2920 2024e; Li et al., 2024e; Xiao et al., 2023). Training- 2970  
2921 based knowledge integration embeds static mem- 2971  
2922 ory in model parameters, enabling generalization 2972  
2923 across tasks. Together, these mechanisms enable 2973  
2924 LLMs to balance flexibility, efficiency, and gener- 2974  
2925 alization in diverse applications. 2975

### 2926 **B.3 Taxonomy of Memory Implementations**

2927 Different from existing surveys that connect differ- 2976  
2928 ent LLM memory types with human memory with 2977  
2929 one-to-one mapping, we observe that all different 2978  
2930 LLM memory mechanisms can be used to imple- 2979  
2931 ment human-like memory with LLM. However, 2980  
2932 different LLM memory mechanisms have different 2981

2932 characteristics, advantages and limitations, which 2982  
2933 are summarized in Table 4 and 5. 2983

**In-Context Learning.** In-context learning sup- 2984  
2935 ports both *short-term memory* (*processing imme-* 2985  
2936 *diolate context*) and *long-term memory*. It relies 2986  
2937 on the self-attention mechanism to focus on rel- 2987  
2938 evant tokens, with no memory scale limitation for 2988  
2939 short-term use but high inference costs for long- 2989  
2940 term memory due to processing large contexts. 2990  
2941 In-context learning excels in episodic (context- 2991  
2942 specific events) and semantic (factual knowledge) 2992  
2943 memory, as prompts can encode task-specific exam- 2993  
2944 ples or facts. Procedural memory (task processes) 2994  
2945 is limited due to reliance on explicit instructions. 2995

2946 Since in-context learning does not modify param- 2996  
2947 eters, it avoids catastrophic forgetting, preserving 2997  
2948 pretrained knowledge. Short-term memory has no 2998  
2949 limitation, as it depends on the prompt size. Long- 2999  
2950 term memory faces high inference costs due to 3000  
2951 the  $\mathcal{O}(n^2)$  complexity of self-attention for long se- 3001  
2952 quences, and the lost-in-the-middle problem. The 3002  
2953 memory is explicit in the prompt, making it inter- 3003  
2954 pretable as the model’s output directly reflects the 3004  
2955 provided context. No additional training or stor- 3005  
2956 age is required, though inference costs increase 3006  
2957 with context length. Malicious or biased prompts 3007  
2958 can influence outputs, raising concerns about mis- 3008  
2959 use. High for episodic and semantic memory, as 3009  
2960 prompts can encode specific events or facts. Lim- 3010  
2961 ited for procedural memory, as it requires explicit 3011  
2962 task instructions, which may not generalize well. 3012

### 2963 **KV Cache or Intermediate Representations.**

2964 The KV cache primarily functions as short-term 3013  
2965 memory, storing  $\mathbf{K}$  and  $\mathbf{V}$  vectors during the pre- 3014  
2966 fill and decoding phases of inference. It retains 3015  
2967 contextual information for the current sequence, 3016  
2968 enabling efficient token generation. Limited appli- 3017  
2969 cability, as the cache is typically cleared between 3018  
2970 sessions. 3019

2971 The KV cache is a runtime mechanism that does 3020  
2972 not alter model parameters, preserving pretrained 3021  
2973 knowledge. Limited by memory constraints, as 3022  
2974 storing  $\mathbf{K}$  and  $\mathbf{V}$  vectors for long sequences can 3023  
2975 be memory-intensive. While the KV cache stores 3024  
2976 intermediate representations, interpreting their con- 3025  
2977 tent is less straightforward than in-context learning, 3026  
2978 as it involves analyzing attention weights and vec- 3027  
2979 tors. The KV cache requires additional memory to 3028  
2980 store vectors, but optimizations like compression 3029  
2981 reduce this cost. Inference efficiency is improved 3030  
2982 compared to recomputing attention. Low. The 3031  
2983 KV cache is a technical optimization with mini- 3032

Table 3: Examples of using different text memory (or can be KV cache).

<b>Question:</b> Can you remind me about the trip I mentioned planning to Paris?
<b>Memory (Vector Database):</b> Embedding match - Conversation: "User plans a trip to Paris in September 2025, interested in visiting the Louvre and Eiffel Tower." Vector Database: Enables semantic retrieval by matching the query's meaning to stored embeddings, useful for broad or vague queries.
<b>Memory (Time Index):</b> On July 15, 2025, at 14:30, the user said, "I'm planning a trip to Paris next month and want to see the Louvre." Time Index: Organizes memories chronologically, ideal for queries referencing recent or specific dates.
<b>Memory (Username Index):</b> User "JaneDoe123" discussed a Paris trip, mentioning a preference for art museums. Username Index: Ensures personalization by linking memories to a specific user, enhancing relevance.
<b>Memory (Event Name Index):</b> Event "Paris Trip 2025": User plans to visit Paris, focusing on cultural landmarks. Event Name Index: Tags memories with specific events, allowing precise retrieval for event-related queries.
<b>Memory (Story Index):</b> Story "Jane's European Adventure": Includes a chapter on planning a Paris trip, with details about booking a hotel near the Seine. Story Index: Structures memories as narratives, preserving context across related interactions.
<b>Memory (Place Index):</b> Place "Paris, France": User mentioned visiting the Eiffel Tower and dining at a café in Montmartre. Place Index: Associates memories with locations, enabling spatial queries about specific places.

Implementation Ways	Memory Type	Forgetting Pretrained Knowledge	Memory Scalability	Explainability	Serving Costs
In-context Learning	Short Term	No	High	High	Low
	Long Term	No	Weak	High	Low
Parameter by Training	Short Term	Weak	Hight	Low	Medium
	Long Term	Severe	Hight	Low	High

Table 4: Reconstructed Memory Implementation Characteristics.

mal impact on output content, though improper cache management could affect output coherence. High for episodic memory, as the cache retains sequence-specific context. Limited for semantic or procedural memory, as it does not store generalized knowledge or task processes independently of the input sequence.

**Parameter by Training.** Training embeds both short-term (immediate patterns) and long-term (generalized knowledge) memory into parameters. Short-term memory has weak forgetting, as recent patterns are retained, but long-term memory suffers from severe forgetting due to catastrophic forgetting during fine-tuning. Training excels in procedural memory, as models learn task-specific patterns during pretraining. Episodic and semantic memory are moderately supported, as specific events or facts are compressed into parameters but may be less precise. Fine-tuning can overwrite pretrained knowledge, especially for long-term memory, leading to catastrophic forgetting. The model's parameters can encode vast amounts of knowledge, limited only by model size and training data. Knowledge embedded in parameters is opaque, making it difficult to trace specific outputs to learned patterns. Training requires significant computational resources, especially for large models and datasets. None. Knowledge is fixed in parameters, reducing risks from external inputs, though biases in training data can persist. High for procedural memory, as

training optimizes task-specific patterns. Moderate for episodic and semantic memory, as specific events or facts are generalized but may lose granularity.

#### B.4 How Human Memory Benefits LLM Agentic Applications

Memory is an indispensable component in various practical LLM-based agent applications. For instance, in a conversational agent, memory stores information about historical conversations, providing the necessary context for generating coherent and relevant responses; without it, the agent cannot maintain a continuous conversation (Lu et al., 2023b). Similarly, in a simulation agent, memory is crucial for maintaining consistent role profiles, preventing the agent from deviating from its assigned character during a simulation (Wang et al., 2025b, 2023f). These examples underscore that memory is not an optional feature but a necessary component for LLM-based agents to effectively accomplish their given tasks. Thus, the cognitive basis of human memory, coupled with its necessity for agent self-evolution and practical applications, provides critical insights for designing sophisticated memory mechanisms in LLM-based systems.

**Information Retrieval and Processing.** Long-context LLMs like Longformer and LongT5 enhance response relevance and document summarization by processing larger text segments, reduc-

Implementation Ways	Memory Type	Forgetting Pretrained Knowledge	Explainability	Costs	Knowledge Match Degree
In-context Learning	Procedural	No	High	Low	Low
	Episodic	No	High	Low	High
	Semantic	No	High	Low	High
Parameter by Training	Procedural	Possible	Low	High	High
	Episodic	Possible	Low	High	Moderate
	Semantic	Possible	Low	High	Moderate

Table 5: Reconstructed Characteristics of Memory Types by Implementation.

ing reliance on external RAG tools (Jin et al., 2024a; Shi et al., 2024; Beltagy et al., 2020; Guo et al., 2022; Jin et al., 2024b). Advanced semantic vector models, such as text-embedding-3-large, jina-embeddings-v2, and BGE-M3, overcome window size limitations, improving usability in tasks like translating complex documents and entire novels (Zhu et al., 2023d; Wang et al., 2024d; OpenAI, 2024b; Günther et al., 2023; Chen et al., 2024a; Zhu et al., 2024; Saad-Falcon et al., 2024; Herold and Ney, 2023; Wang et al., 2024f; Lyu et al., 2024).

**Chatbots.** Long-context processing enhances chatbots by enabling extended memory and contextual coherence, as seen in platforms like ChatGPT, Pi, Character AI, and Talkie, which use persistent memory and techniques like prompt-based memorization, memory-augmented architectures, and context extension for style-consistent, engaging dialogues (OpenAI, 2024a; Inflection, 2023; Character AI, 2023; Ai, 2024; Lee et al., 2023a; Zhong et al., 2024c; Wang et al., 2023h, 2024i).

**Code Development.** LLMs leverage memory to store development knowledge and conversational context, with models like StarCoder2, Qwen2.5-Coder, and Granite Code Models enabling scalable code completion and predictive debugging in tools like GitHub Copilot and Anysphere Cursor (Tsai et al., 2023; Chen et al., 2023a; Qian et al., 2023; Li et al., 2023e; Zhang et al., 2024a; Lozhkov et al., 2024; Hui et al., 2024; Mishra et al., 2024; GitHub, 2022; Anysphere, 2025).

**Social Simulation.** Memory defines character traits for realistic role-playing and supports multi-agent social simulations by improving self-monitoring, maintaining economic environments, and simulating dynamic behaviors (Wang et al., 2025b; Shao et al., 2023; Kaiya et al., 2023; Gao et al., 2023; Li et al., 2023b,e; Hua et al., 2023).

**Personal Assistant.** LLM-based personal assistants rely on memory for consistent, personalized dialogues, using textual retrieval, conversation sum-

marization, and external tools to maintain conversational flow (Lu et al., 2023b; Lee et al., 2023b; Pan et al., 2023a; Wu et al., 2023).

**Application in Specific Domains.** Long-context LLMs improve coherence in news summaries, simplify legal document interpretation, enhance healthcare and financial decision-making, and advance drug discovery and scientific problem-solving by leveraging external knowledge and memory (Gao et al., 2019; Kapoor et al., 2024; Fan et al., 2024; Reddy et al., 2024; Masry and Hajian, 2024; Nie et al., 2024; Hilgert et al., 2024; Shao and Yan, 2024; Wang et al., 2023c; Xiong et al., 2023; Liu et al., 2023d; Wang et al., 2023d; Yunxiang et al., 2023; Chen et al., 2024b; Zhao et al., 2024c; Chen et al., 2023c; Wang et al., 2023m; Qiang et al., 2023).

## C Taxonomy of Different Memory

3085  
3086  
3087  
3088  
3089  
3090  
3091  
3092  
3093  
3094  
3095  
3096  
3097  
3098  
3099  
3100  
3101  
3102

Retrieval Augmented Generation §2.1		
<b>Memory Construction</b>	Unstructured Text Sources	CREAICL (Li et al., 2023c), CRAG (Yan et al., 2024),
	Knowledge Graph	TableGPT (Zha et al., 2023), PKG (Luo et al., 2023), KnowledGPT (Wang et al., 2023j), G-Retriever (He et al., 2024a)
	Internal Knowledge	SKR (Wang et al., 2023k), GenRead (Yu et al., 2022b), Selfmem (Cheng et al., 2023)
<b>Data Indexing</b>	Granularity-based	(Shi et al., 2023), CoN (Yu et al., 2023), DenseX (Chen et al., 2023b), LLMIndexer (Jin et al., 2023), LLM-R (Wang et al., 2023g)
	Chunk-based	Chunk (Teja, 2023), LangChain (Langchain, 2023), Small2Big (Yang, 2023)
	Metadata-enriched	Reverse HyDE (Gao et al., 2022)
	Graph-based	KGP (Wang et al., 2023l)
	Sparse Embedding	BM25
	Dense Embedding Hybrid	AngIE (Li and Li, 2023), Voyage (VoyageAI, 2023), BGE (BAAI, 2023)
<b>Memory Update</b>	KG-based	KGP (Wang et al., 2023l), Reverse HyDE (Gao et al., 2022), LangChain (Langchain, 2023), Small2Big (Yang, 2023)
	Retrieval Invocation	SKR (Wang et al., 2023k), GenRead (Yu et al., 2022b), Selfmem (Cheng et al., 2023)
<b>Memory Query</b>	Query Reformulation	(Zhou et al., 2023a), CoVe (Dhuliawala et al., 2023), RRR (Ma et al., 2023), BEQUE (Peng et al., 2023), HyDE (Gao et al., 2022), Take a step back (Zheng et al., 2024a)
	Query Routing	SemanticRouter (Wang et al., 2025a)
	Multi-step Querying	DSP (Khattab et al., 2022), FLARE (Jiang et al., 2023c), Self-RAG (Asai et al., 2023), BGM (Ke et al., 2024), RA-DIT (Lin et al., 2023b)

Table 6: Taxonomy of RAG.

Agentic Memory §2.2		
<b>Memory Construction</b>	Summarization	MemoryBank (Zhong et al., 2024b), RET-LLM (Modarressi et al., 2023)
	Key-value	RET-LLM (Murre and Dros, 2015), Meminsight (Salama et al., 2025), Memocrs (Xi et al., 2024),
	Semantic Representations	Memorybank (Zhong et al., 2024b), (Pan et al., 2025), Mem0 (mem0ai, 2024) (Zhong et al., 2024b; Pan et al., 2025; mem0ai, 2024)
	Relation Graph	CGSN (Nie et al., 2022), GraphReader (Li et al., 2024c), HippoRAG (Gutiérrez et al., 2024)
	Auxiliary Signals	LongMemEval (Wu et al., 2024a), THEANINE (iunn Ong et al., 2025)
	Storage and Inference Efficiency	Llmlingua (Jiang et al., 2023a), AutoCompressor (Chevalier et al., 2023), TCRA-LLM (Liu et al., 2023a), Promptcache (Gim et al., 2024)
<b>Memory Updating</b>	Summarization and Restructuring	MemoryBank (Zhong et al., 2024b), ChatGPT-RSum (Wang et al., 2025d), MemoChat (Lu et al., 2023a),
	Reasoning and Self-Reflection	ReAct (Yao et al., 2022), Reflexion (Shinn et al., 2024), BoT (Yang et al., 2024e), Agent Workflow Memory (Wang et al., 2024n)
	Reasoning with Experience	TiM (Liu et al., 2023b), GITM (Zhu et al., 2023b), Voyager (Wang et al., 2023b), Retroformer (Yao et al., 2023), ExpeL (Zhao et al., 2024a), LD-Agent (Li et al., 2024a)
	Memory Evolution	A-MEM (Xu et al., 2025), Synapse (Zheng et al., 2024c), R2I (Samsami et al., 2024), SCM (Wang et al., 2024a), Selective Editing (Bae et al., 2022), Blending (Kim et al., 2024b), Recursive Summarization (Wang et al., 2025c), Self-Reflection (Sun et al., 2024)
<b>Memory Query</b>	Query-centered	FLARE (Jiang et al., 2023b), IterCQR (Jang et al., 2024)
	Memory-centered	LongMemEval (Wu et al., 2024a), Perlta (Du et al., 2024)
	Event and Structure-aware	LoCoMo (Maharana et al., 2024), CC (Jang et al., 2023), MSC (Xu et al., 2021), HippoRAG (Gutiérrez et al., 2024), Memorag (Qian et al., 2024)

Table 7: Taxonomy of Agent Memory.

KV cache as Memory §3.1		
<b>Memory Construction</b>	KV Construction	Llama Series (Touvron et al., 2023a,b; Grattafiori et al., 2024), DeepSeek (DeepSeek-AI et al., 2025a,b).
<b>Memory Update</b>	Eviction and dropping	StreamingLLM (Xiao et al., 2024a), LM-Infinite (Han et al., 2024a), H <sub>2</sub> O (Zhang et al., 2023a), FastGen (Ge et al., 2024), Radar (Hao et al., 2025), NACL (Chen et al., 2024d)
	Attention guided elimination	Scissorhands (Liu et al., 2023e), L <sub>2</sub> Norm (Devoto et al., 2024), SirLLM (Yao et al., 2024), D-LLM (Jiang et al., 2024b), ZigZagKV (Zhong et al., 2024a), DynamicKV (Zhou et al., 2025)
	Merging and Semantic compression	MiniCache (Liu et al., 2024a), InfiniPot (Kim et al., 2024a), CHAI (Agarwal et al., 2024), Activation Beacon (Zhang et al., 2025a), CacheGen (Liu et al., 2024d), ChunkKV (Liu et al., 2025d), SentenceKV (Zhu et al., 2025c)
	Quantization and Low-rank Approximation	KIVI (Liu et al., 2024e), SKVQ (Duanmu et al., 2024), QAA (Dong et al., 2024b), KVQuant (Hooper et al., 2024), CQ (Zhang et al., 2024b), AnTKV (Li et al., 2025b), SmoothAttention (Lin et al., 2025), MiKV (Yang et al., 2024d), LESS (Dong et al., 2024a), Eigen (Saxena et al., 2024), GEAR (Kang et al., 2024), FlexGen (Sheng et al., 2023), Atom (Zhao et al., 2024b), ZipCache (He et al., 2024b)
	System and Task-aware Allocation	KVDirect (Chen et al., 2024c), FlowKV (Li et al., 2025a) PyramidInfer (Yang et al., 2024a), ChunkKV (Liu et al., 2025d), OracleKV (Zhu et al., 2025b)
<b>Memory Query</b>	KV selection	QUEST (Tang et al., 2024b), TokenSelect (Wu et al., 2025a), Selective Attention (Leviathan et al., 2025), RetrievalAttention (Liu et al., 2024b)
	KV Reuse	RadixAttention (Zheng et al., 2024b), ChunkAttention (Ye et al., 2024), KVShare (Yang et al., 2025a), Cache-craft (Agarwal et al., 2025), Teola (Tan et al., 2025), KVLink (Yang et al., 2025b), CacheBlend (Yao et al., 2025), EPIC (Hu et al., 2025), SubGCACHE (Zhu et al., 2025a), HyperRAG (An et al., 2025), RAGO (Jiang et al., 2025)

Table 8: Taxonomy of KV Cache as Memory.

Other Vectors as Memory §3.2		
<b>External Vector Memory</b>	Sentence-level Encoding	Slot-based encoding (Al Adel and Burtsev, 2021)
	Key-value Stores	kNN-LM (Khandelwal et al., 2019), Memorizing Transformer (Wu et al., 2022)
	Vector cache	MemGPT (Packer et al., 2023), Neurocache (Packer et al., 2023)
<b>Steering Vectors</b>	Structured and Associative Modules	CAMELoT (He et al., 2024c), MemOS (Li et al., 2025c), Memory3 (Yang et al., 2024c)
	Foundational method	PPLM (Dathathri et al.)
	Contrastive-based	Turner et al. (Turner et al., 2023), Liu et al. (Liu et al., 2023c), Zou et al. (Zou et al., 2023), Arditi et al. (Arditi et al., 2024)
	Optimization-based	Subramani et al. (Subramani et al., 2022), Hernandez et al. (Hernandez et al., 2023), Dunefsky et al. (Dunefsky and Cohan, 2025), Mack et al. (Mack and Turner, 2024), Li et al. (Li et al., 2024b), Turner et al. (Turner et al., 2025), Cao et al. (Cao et al., 2024)

Table 9: Taxonomy of Other Vectors as Memory.

Parameter as Memory §4		
<b>Memory Construction</b>	Data Composition	Data augmentation (Allen-Zhu and Li, 2024), Extracting (Carlini et al., 2021), Lee et al. (Lee et al., 2022), Kandpal et al. (Kandpal et al., 2022)
	Sequence length	Carlin et al. (Carlini et al., 2023), Wang et al. (Wang et al., 2024l)
	Model Scale	Mem0 (Tirumala et al., 2022), Carlin et al. (Carlini et al., 2023), Freeman et al. (Freeman et al., 2024), Geva et al. (Geva et al., 2021), Dai et al. (Dai et al., 2021)
<b>Memory Update</b>	Continual Learning	SCM (Wang et al., 2024e), EWC (Kirkpatrick et al., 2017), TaSL (Feng et al., 2024), SELF-PARAM (Wang et al.), POCL (Wu et al., 2024b), DSI++ (Mehta et al., 2022), LSCS (Wang et al., 2024k)
	PEFT	PEFT (Han et al., 2024b), Character-LLM (Shao et al.), AI-Native Memory (Shang et al., 2024), MemoRAG (Qian et al., 2024), Echo (Liu et al., 2025b)
	Model Merging	FedAvg (McMahan et al., 2017), MagMax (Marczak et al., 2024), SNIP (Lee et al., 2019), FisherMerging (Matena and Raffel, 2022), FedSAM (Qu et al., 2022), FedFisher (Jhunjhunwala et al., 2024), Gamegpt (Daheim et al., 2024), TIES (Yadav et al., 2023), DARE (Yu et al., 2024a,b), Model Breadcrumbs (Davari and Belilovsky, 2023), TALL-masks (Wang et al., 2024c), SMEAR (Muqeeth et al., 2024), Twin-Merging (Lu et al., 2024a), Weight-Ensembling MoE (Tang et al., 2024a)
	Task Arithmetic	Task Arithmetic (Ilharco et al., 2022) TIES (Yadav et al., 2023), AdaMerging (Yang et al., 2024b), TwinMerge (Lu et al., 2024b)
	Model Editing	Wang et al. (Wang et al., 2024h), MemoryLLM (Wang et al., 2024j), WISE (Wang et al., 2024g)
<b>Memory Query</b>	Exact Memorization	Carlini et al. (Carlini et al., 2021), Carlin et al. (Carlini et al., 2023), Nasr et al. (Nasr et al., 2023)
	Approximate Memorization	Ippolito et al. (Ippolito et al., 2023)
	Prompt-based Memorization	Biderman et al. (Biderman et al., 2023)

Table 10: Taxonomy of Parameter as Memory.