

Cross-Domain Memory Systems: A Unified Analytical Framework

Yuyang Wang, Tianyi Huang, Boyang Sun

March 22, 2025

Abstract

We present a comprehensive theoretical framework for analyzing memory systems across different AI domains, including LLM-agents, vision-language models, and video understanding systems. This framework enables systematic comparison of memory architectures based on their retrieval mechanisms, memory structures, and update schemas. By examining these systems through the lenses of efficiency, scalability, and privacy preservation, we identify universal patterns and domain-specific optimizations. Our analysis reveals important trade-offs in memory system design and highlights opportunities for cross-domain knowledge transfer to build more effective AI memory architectures.

1 Introduction

Memory systems have become an essential component in modern AI architectures, enabling models to store, retrieve, and update information efficiently. As these systems grow in sophistication, understanding their underlying mechanisms becomes increasingly important. This framework provides a systematic approach to analyzing memory systems across different AI domains:

- Language Agents (LLM-based systems)
- Vision-Language Models (VPT, prompt-based methods)
- Video Understanding Systems
- Multi-modal AI systems

By examining these diverse domains through a unified lens, we aim to identify common patterns, domain-specific innovations, and opportunities for cross-domain knowledge transfer.

2 Theoretical Analysis Framework

2.1 Retrieval Mechanism Taxonomy

Category	Key Approaches	Analysis Dimensions
Similarity-Based	<ul style="list-style-type: none">• Semantic search (LLM-mem)• Contextual similarity• Embedding-based retrieval	<ul style="list-style-type: none">• Retrieval precision• Context sensitivity• Computational complexity
Prompt-Based	<ul style="list-style-type: none">• Direct prompt selection (L2P)• Compositional prompting (CodaPrompt)• Task-specific prompt banks	<ul style="list-style-type: none">• Prompt transferability• Task adaptation capability• Memory efficiency
Temporal-Spatial	<ul style="list-style-type: none">• Temporal correlation (MemFlow)• Spatial attention (XMem)• Sequence modeling	<ul style="list-style-type: none">• Temporal consistency• Spatial coherence• Processing efficiency
Hybrid Methods	<ul style="list-style-type: none">• Multi-modal integration• Cross-attention mechanisms• Ensemble strategies	<ul style="list-style-type: none">• Integration effectiveness• Modal alignment• Robustness to domain shift

Table 1: Taxonomy of retrieval mechanisms across AI domains

2.2 Memory Structure Classification

Structure Type	Architectural Patterns	Reported Advantages/Limitations
Static Memory	<ul style="list-style-type: none"> Fixed capacity memories Pre-trained prompt repositories Immutable knowledge bases 	<ul style="list-style-type: none"> + Low maintenance overhead + Predictable performance – Limited adaptability – Scale constraints
Dynamic Memory	<ul style="list-style-type: none"> Expandable memory banks Continually updated embeddings Adaptive storage allocation 	<ul style="list-style-type: none"> + Adaptability to new data + Scalability with task growth – Higher computational overhead – Potential for memory corruption
Hierarchical Memory	<ul style="list-style-type: none"> Multi-level access structures Priority-based organization Cache-like architectures 	<ul style="list-style-type: none"> + Efficient information access + Organized knowledge storage – Complex management logic – Increased design complexity
Distributed Memory	<ul style="list-style-type: none"> Federated storage systems Shared knowledge repositories Decentralized architectures 	<ul style="list-style-type: none"> + Collaborative knowledge sharing + Enhanced privacy potential – Synchronization challenges – Consistency maintenance costs

Table 2: Classification of memory structures across AI domains

2.3 Memory Update Schema Analysis

Update Approach	Core Mechanisms	Theoretical Implications
Frequency-Based	<ul style="list-style-type: none"> • Usage tracking • Popularity-based retention • LFU/LRU-inspired approaches 	<ul style="list-style-type: none"> • Potential bias toward common patterns • Efficient for repetitive tasks • Performance in long-tail scenarios
Recency-Based	<ul style="list-style-type: none"> • Temporal prioritization • Time-decay functions • Recent-first strategies 	<ul style="list-style-type: none"> • Adaptation to concept drift • Handling temporal dynamics • Historical information loss
Importance-Weighted	<ul style="list-style-type: none"> • Value estimation models • Critical information retention • Salience detection 	<ul style="list-style-type: none"> • Attention mechanism effectiveness • Information preservation quality • Computational overhead for value assessment
Privacy-Preserving	<ul style="list-style-type: none"> • Anonymization techniques • Differential privacy approaches • Federated updates 	<ul style="list-style-type: none"> • Privacy-utility tradeoffs • Information leakage risks • Compliance with privacy standards

Table 3: Analysis of memory update schemas across AI domains

3 Cross-Domain Comparison Methodology

3.1 Standardized Evaluation Metrics

3.1.1 Efficiency Metrics Analysis

- Computational complexity (theoretical and reported)
- Memory footprint (parameters and storage requirements)
- Scaling properties (mathematical complexity analysis)
- Update efficiency (reported or theoretically derived)

3.1.2 Performance Analysis

- Retrieval accuracy (reported precision/recall metrics)
- Domain-specific task performance (reported benchmarks)
- Generalization capability (cross-task performance)
- Adaptation ability (few-shot/zero-shot capabilities)

3.1.3 Privacy & Security Analysis

- Reported privacy guarantees
- Theoretical vulnerability assessment
- Data exposure mitigation strategies
- Compliance mechanisms

3.2 Cross-Domain Comparison Framework

3.2.1 Baseline Comparison Methodology

- Standardized performance metrics extraction
- Normalization across different domains' reporting
- Identification of comparable benchmark tasks
- Development of cross-domain evaluation criteria

3.2.2 Cross-Domain Transfer Assessment

- Theoretical transferability analysis
- Adaptation requirements estimation
- Architecture compatibility evaluation
- Domain shift robustness analysis

3.2.3 Unified Analysis Approach

- Controlling for model capacity differences
- Normalizing for domain-specific challenges
- Isolating memory component contributions
- Establishing fair comparison baselines

4 Trade-off Analysis Framework

4.1 Key Trade-off Dimensions

4.1.1 Information Compression vs. Retrieval Accuracy

- Analysis of reported Pareto frontiers
- Theoretical bounds on compression-accuracy tradeoffs
- Cross-architecture comparison of operating points
- Identification of optimal configuration patterns

4.1.2 Computational Efficiency vs. Memory Capacity

- Analysis of reported inference costs against memory sizes
- Theoretical scaling properties with dataset growth
- Identification of domain-specific bottlenecks
- Cross-architecture efficiency comparisons

4.1.3 Privacy Preservation vs. Utility

- Analysis of reported privacy-utility tradeoffs
- Theoretical privacy guarantee comparisons
- Information loss quantification methodologies
- Cross-domain privacy mechanism effectiveness

4.2 Cross-Domain Pattern Identification

4.2.1 Universal Memory Patterns

- Temporal priority mechanisms (prevalence and effectiveness)
- Context-sensitive retrieval (approaches and outcomes)
- Compression-accuracy balancing strategies

4.2.2 Domain-Specific Optimizations

- Language: semantic organization techniques
- Vision: spatial relation preservation approaches
- Video: temporal consistency mechanisms
- Multi-modal: cross-modal alignment strategies

5 Research Gap Analysis

5.1 Underexplored Integration Opportunities

- Video temporal mechanisms in language systems
- Prompt-tuning applications in video understanding
- Neurobiological approaches in multi-modal systems
- LLM-agent memory techniques for vision systems

5.2 Theoretical Limitations

- Boundary conditions for memory system effectiveness
- Theoretical limits to compression-accuracy tradeoffs
- Privacy-utility fundamental constraints
- Cross-domain generalization barriers

5.3 Future Research Directions

- Unified memory architectures for cross-domain applications
- Privacy-preserving memory update mechanisms
- Efficient memory compression techniques
- Temporal-aware retrieval for dynamic environments
- Cross-modal memory sharing protocols

6 Conclusion

This theoretical framework provides a comprehensive approach for analyzing memory systems across diverse AI domains. By examining these systems through a unified lens, we can identify common patterns, unique innovations, and transferable techniques. This analysis will inform the development of more efficient, scalable, and privacy-preserving memory architectures for next-generation AI systems.

References