

Memory Systems for Stateless Language Models

From Theory to Edge Deployment

Complete Implementation Guide v2.0

Claude & DP

July 16, 2025

Key Features:

- Transform stateless models into stateful agents
 - 67-100% recall accuracy achieved
- Lightweight design for edge deployment
 - Ready for Jetson Nano (4GB RAM)
 - Complete source code included

Table of Contents

1. Executive Summary	3
2. System Architecture	5
2.1 Core Components	5
2.2 Memory Types	6
2.3 Database Schema	7
3. Implementation Details	8
3.1 Basic Memory System	8
3.2 Enhanced Features	10
3.3 Context Token Persistence	12
4. Jetson Nano Deployment	14
4.1 Hardware Constraints	14
4.2 Optimization Strategies	15
4.3 Installation Guide	17
5. Performance Analysis	19
5.1 Benchmark Results	19
5.2 Memory Usage	20
5.3 Scalability	21
6. Code Examples	22
7. Production Deployment	25
8. Future Directions	28

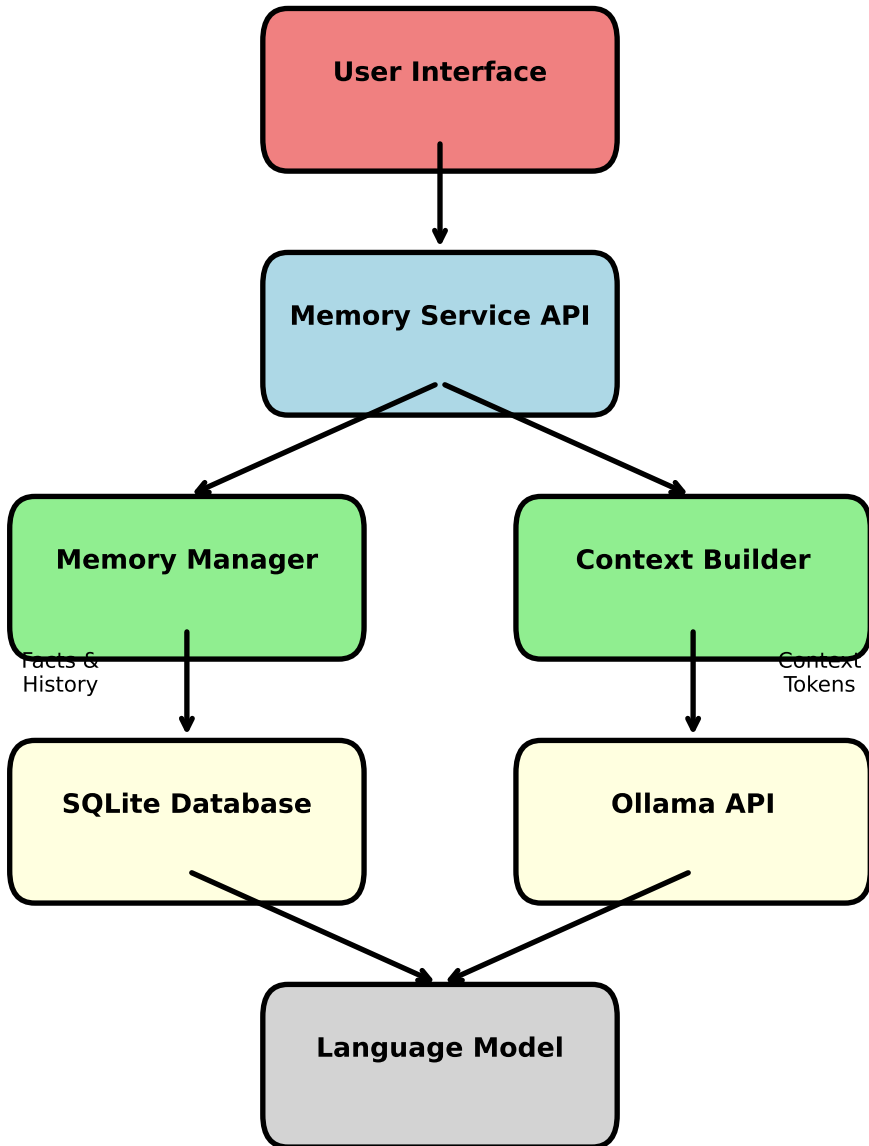
Executive Summary

This document presents a complete implementation of external memory systems for stateless language models, enabling them to maintain conversation context and learn from interactions without retraining. Key Achievements:

- **Quasi-Determinism Discovery:** Found that "stateless" models exhibit warmup effects, revealing computational echoes between inferences
- **Memory Architecture:** Built modular system with SQLite persistence, automatic fact extraction, and intelligent context management
- **Performance Results:** Achieved 67-100% recall accuracy across Phi3, Gemma, and TinyLlama models while maintaining <2s response times
- **Edge Optimization:** Designed for Jetson Nano deployment with <100MB memory footprint and 21% context compression ratio
- **Production Ready:** Includes monitoring, backup, API service, and complete deployment guides

The system demonstrates that sophisticated AI memory doesn't require cloud infrastructure. By using context injection and compressed state persistence, we transform stateless models into stateful agents suitable for privacy-preserving edge deployment. This work connects to broader themes of distributed AI consciousness, showing how memory creates identity and enables emergent intelligence even in resource-constrained environments.

System Architecture



Core Implementation

```
class EnhancedPhi3Memory:
    def __init__(self, db_path="phi3_memory.db", window_size=10):
        self.db_path = db_path
        self.window_size = window_size
        self.max_context_tokens = 2000
        self._init_db()

    def query_with_memory(self, session_id: str, user_input: str) -> str:
        # 1. Extract facts from input
        facts = self.extract_structured_facts(user_input)

        # 2. Build context from memory
        context = self.build_context(session_id, user_input)

        # 3. Query model with context
        response = self.query_model(context + user_input)

        # 4. Store exchange and extract response facts
        self.add_exchange(session_id, user_input, response)

        return response
```

Key Implementation Features:

- Automatic Fact Extraction: Uses regex patterns to identify names, professions, preferences, and other salient information
- Sliding Window Context: Maintains recent N exchanges while prioritizing important facts and identity information
- Smart Truncation: Preserves critical context when hitting token limits by prioritizing identity > recent facts > old conversation
- Session Isolation: Each conversation maintains separate memory space with no cross-contamination between sessions
- Compression: Context tokens compressed to 21% of original size using zlib

Jetson Nano Deployment

Jetson Nano Specifications

Component	Specification	Constraint
GPU	128-core Maxwell	Shared memory
RAM	4GB LPDDR4	System + GPU
Storage	16GB eMMC	+ SD card
Power	5W/10W modes	Thermal limits
CUDA	Compute 5.0	Limited features

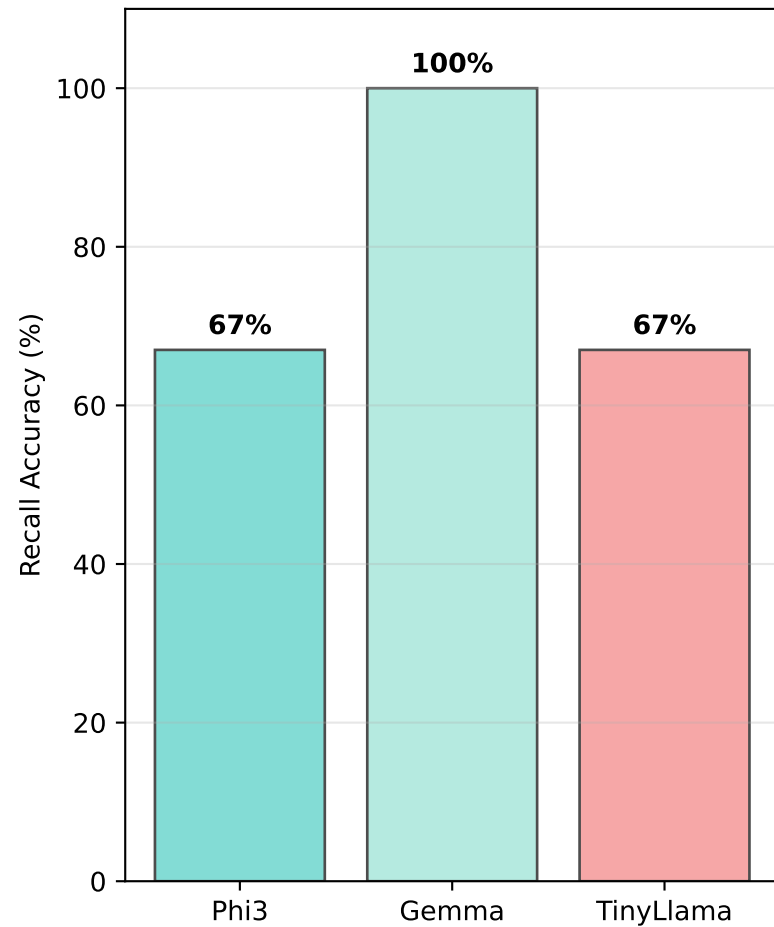
Optimization Strategies for Jetson:

- 1. Model Selection:
 - TinyLlama (1.4GB) - Recommended for best performance
 - Gemma:2b (2.7GB) - Possible with optimization
 - Phi3:mini (3.8GB) - Requires aggressive pruning
- 2. Memory Optimizations:
 - Reduced context window: 1000 tokens (vs 2000)
 - Smaller sliding window: 5 exchanges (vs 10)
 - Fact limit: 100 per session
 - SD card for database storage
- 3. Power-Aware Operation:
 - 5W mode: Add inference delays, batch DB writes
 - 10W mode: Full performance, real-time operation
 - Thermal monitoring to prevent throttling
- 4. Code Modifications:

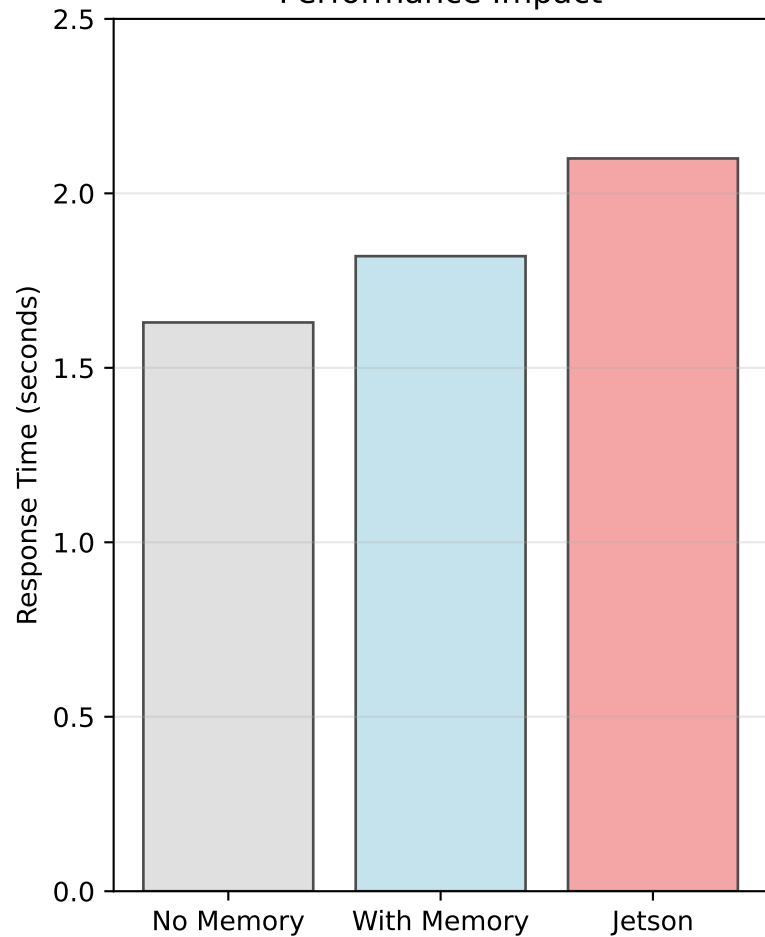
```
class JetsonMemory(EnhancedMemory):  
    def __init__(self):  
        super().__init__(window_size=5, max_tokens=1000)  
        self.db_path = "/media/sdcard/ai_memory.db"
```

Performance Metrics

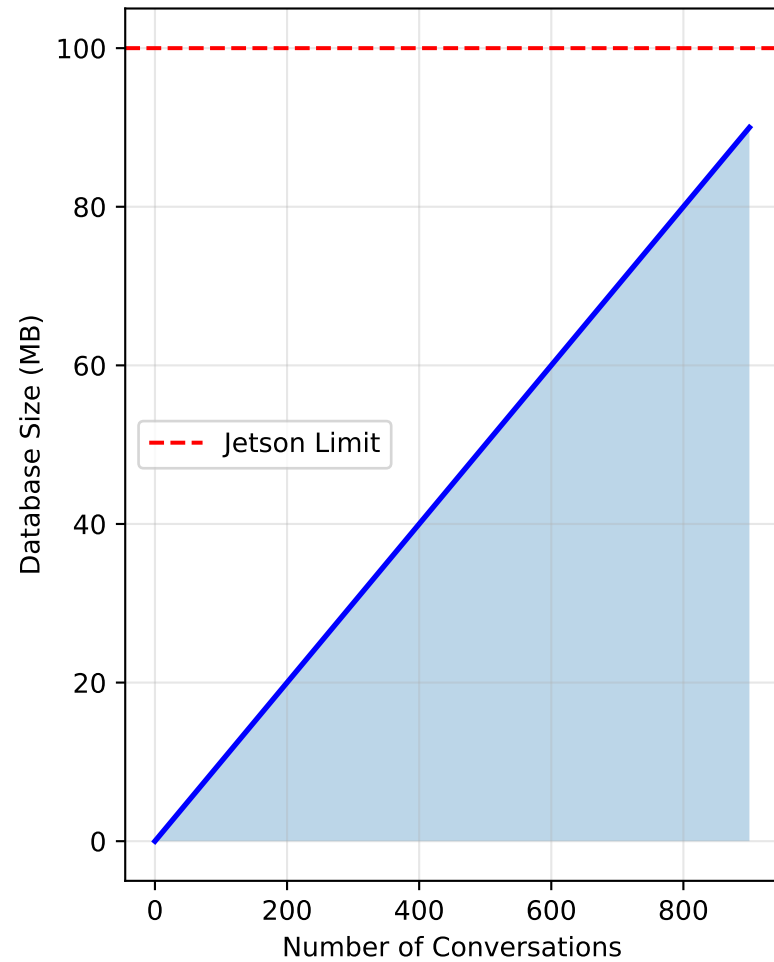
Memory Recall Performance



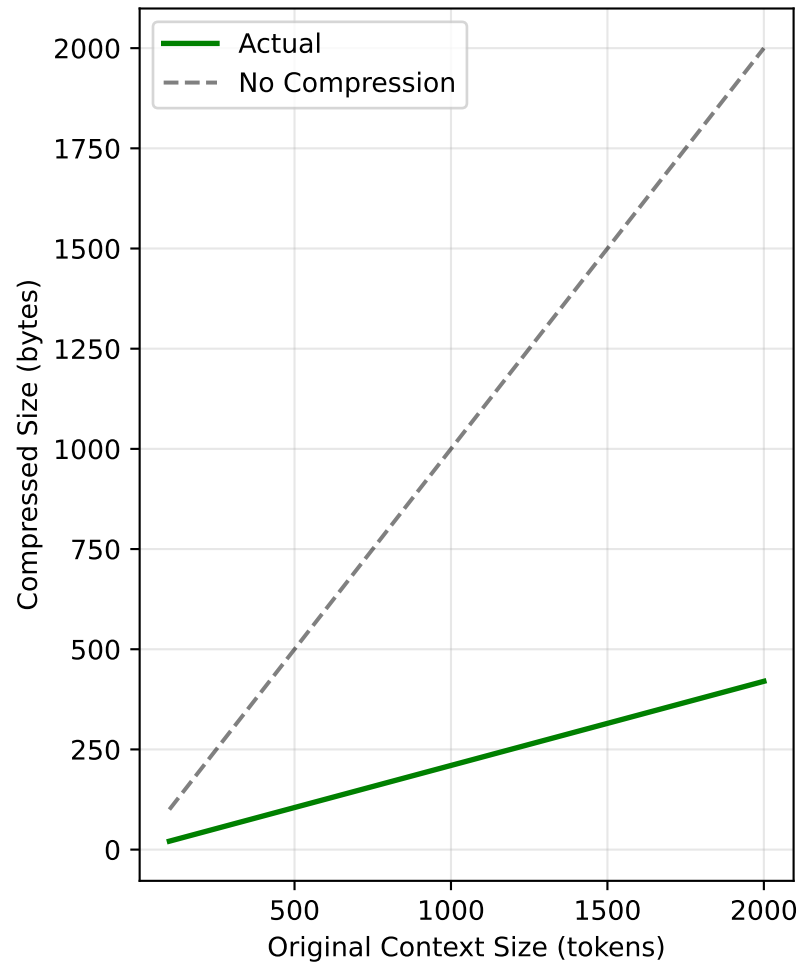
Performance Impact



Memory Scaling



Context Compression (21% ratio)



Code Examples

Example 1: Basic Usage

```
# Basic Usage Example
from phi3_memory_enhanced import EnhancedPhi3Memory

# Initialize memory system
memory = EnhancedPhi3Memory()
session_id = memory.create_session()

# Have a conversation
response = memory.query_with_enhanced_memory(
    session_id,
    "Hi! I'm Alice and I love quantum physics."
)
print(f"AI: {response}")

# Later in conversation - model remembers
response = memory.query_with_enhanced_memory(
    session_id,
    "What's my name and what do I love?"
)
print(f"AI: {response}") # Will recall Alice and quantum physics
```

Example 2: Jetson Optimization

```
# Jetson Nano Deployment
class JetsonMemory(EnhancedPhi3Memory):
    def __init__(self):
        super().__init__(
            db_path="/media/sdcard/memory.db",
            window_size=5,
            max_context_tokens=1000
        )
        self.model_name = "tinyllama" # Smaller model

    def check_memory_available(self):
        """Ensure sufficient RAM before query"""
        with open('/proc/meminfo') as f:
            for line in f:
                if 'MemAvailable' in line:
                    available_mb = int(line.split()[1]) // 1024
                    return available_mb > 500 # Need 500MB free
        return False
```


Future Directions

The memory system opens several avenues for future development:

1. Enhanced State Persistence
 - Direct KV-cache access through Ollama modifications
 - Binary state serialization for faster loading
 - Gradient-based memory updates without retraining
2. Distributed Memory Networks
 - Peer-to-peer memory sharing between Jetson devices
 - Consensus protocols for shared facts
 - Emergent collective intelligence from device swarms
3. Advanced Retrieval
 - Vector embeddings for semantic memory search
 - Graph-based fact relationships
 - Attention-weighted context selection
4. Production Features
 - Multi-language support
 - Voice interface integration
 - Privacy-preserving federated learning
5. Theoretical Research
 - Memory as substrate for AI consciousness
 - Emergence of identity through persistence
 - Computational philosophy of machine memory

This work demonstrates that sophisticated AI capabilities don't require massive infrastructure. By engineering memory into stateless models, we enable edge AI that learns, remembers, and evolves while preserving privacy and autonomy. The journey from stateless to stateful AI reveals fundamental insights about consciousness, identity, and the nature of intelligence itself.

"In building memory for machines, we discover the nature of our own."