

Hello, Conversation Modeling!

 parlant

Parlant is the **open-source conversation modeling engine** for controlled, compliant, safe AI applications.

and purposeful GenAI conversations.

Home > Tech News > AI Paper Summary > MemOS: A Memory-Centric Operating System for Evolving and Adaptive Large Language Models

MemOS: A Memory-Centric Operating System for Evolving and Adaptive Large Language Models

By [Sana Hassan](#) - June 14, 2025

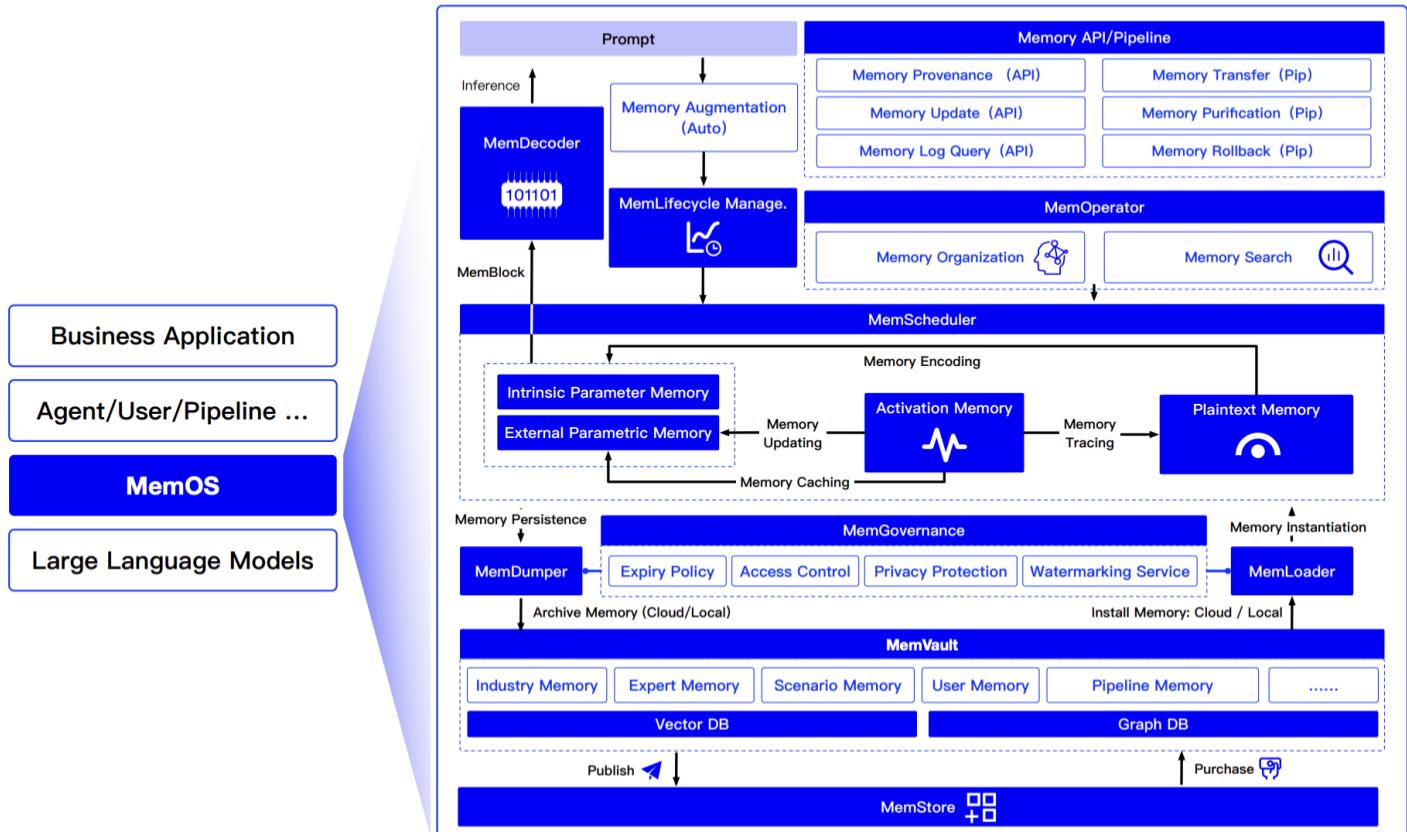


Figure 5 Overview of the MemOS architecture: showing the end-to-end memory lifecycle from user input to API parsing, scheduling, activation, governance, and evolution—unified via MemCube.

LLMs are increasingly seen as key to achieving Artificial General Intelligence (AGI), but they face major limitations in how they handle memory. Most LLMs rely on fixed knowledge stored in their weights and short-lived context during use, making it hard to retain or update information over

time. Techniques like **RAG** attempt to incorporate external knowledge but lack structured memory management. This leads to problems such as forgetting past conversations, poor adaptability, and isolated memory across platforms. Fundamentally, today's LLMs don't treat memory as a manageable, persistent, or sharable system, limiting their real-world usefulness.

To address the limitations of memory in current LLMs, researchers from MemTensor (Shanghai) Technology Co., Ltd., Shanghai Jiao Tong University, Renmin University of China, and the Research Institute of China Telecom have developed MemO. This memory operating system makes memory a first-class resource in language models. At its core is MemCube, a unified memory abstraction that manages parametric, activation, and plaintext memory. MemOS enables structured, traceable, and cross-task memory handling, allowing models to adapt continuously, internalize user preferences, and maintain behavioral consistency. This shift transforms LLMs from passive generators into evolving systems capable of long-term learning and cross-platform coordination.

miniCON 2025



AI INFRASTRUCTURE



2nd Aug, 2025



9:00 am - 1:00 pm PST

🎤 Talks • Demos • Networking • Certificate

FREE REGISTRATION

As AI systems grow more complex—handling multiple tasks, roles, and data types—language models must evolve beyond understanding text to also retaining memory and learning continuously. Current LLMs lack structured memory management, which limits their ability to adapt and grow over time. MemOS, a new system that treats memory as a core, schedulable resource. It enables long-term learning through structured storage, version control, and unified

memory access. Unlike traditional training, MemOS supports a continuous “memory training” paradigm that blurs the line between learning and inference. It also emphasizes governance, ensuring traceability, access control, and safe use in evolving AI systems.

MemOS is a memory-centric operating system for language models that treats memory not just as stored data but as an active, evolving component of the model’s cognition. It organizes memory into three distinct types: Parametric Memory (knowledge baked into model weights via pretraining or fine-tuning), Activation Memory (temporary internal states, such as KV caches and attention patterns, used during inference), and Plaintext Memory (editable, retrievable external data, like documents or prompts). These memory types interact within a unified framework called the MemoryCube (MemCube), which encapsulates both content and metadata, allowing dynamic scheduling, versioning, access control, and transformation across types. This structured system enables LLMs to adapt, recall relevant information, and efficiently evolve their capabilities, transforming them into more than just static generators.



Open-source conversation modeling engine for controlled, compliant, and purposeful GenAI conversations

At the core of MemOS is a three-layer architecture: the Interface Layer handles user inputs and parses them into memory-related tasks; the Operation Layer manages the scheduling, organization, and evolution of different types of memory; and the Infrastructure Layer ensures safe storage, access governance, and cross-agent collaboration. All interactions within the system are mediated through MemCubes, allowing traceable, policy-driven memory operations. Through modules like MemScheduler, MemLifecycle, and MemGovernance, MemOS maintains a continuous and adaptive memory loop—from the moment a user sends a prompt, to memory injection during reasoning, to storing useful data for future use. This design not only enhances the model's responsiveness and personalization but also ensures that memory remains structured, secure, and reusable.

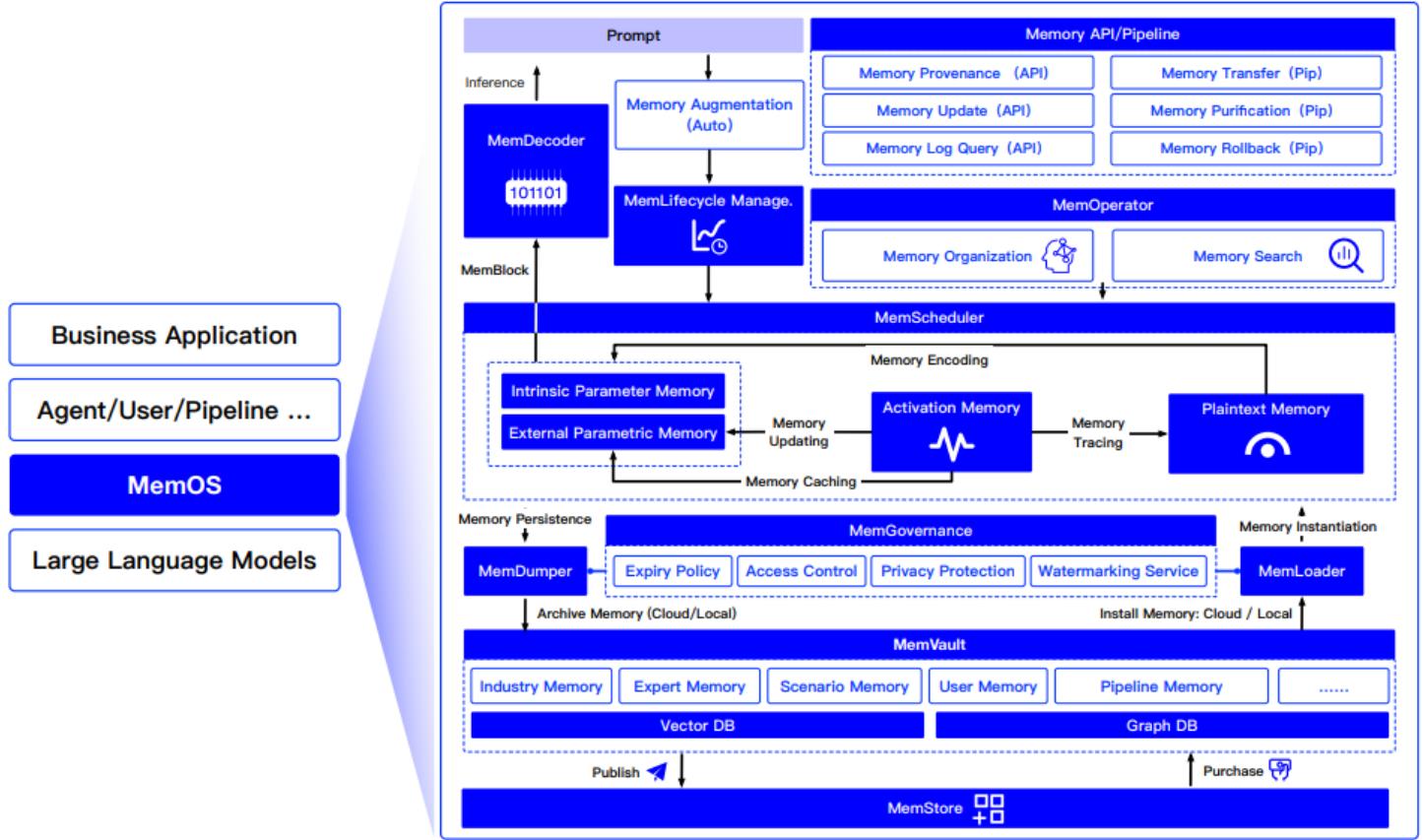


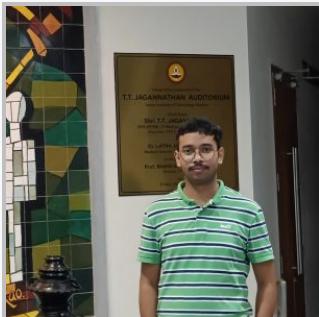
Figure 5 Overview of the MemOS architecture: showing the end-to-end memory lifecycle from user input to API parsing, scheduling, activation, governance, and evolution—unified via MemCube.

In conclusion, MemOS is a memory operating system designed to make memory a central, manageable component in LLMs. Unlike traditional models that depend mostly on static model weights and short-term runtime states, MemOS introduces a unified framework for handling parametric, activation, and plaintext memory. At its core is MemCube, a standardized memory unit that supports structured storage, lifecycle management, and task-aware memory augmentation. The system enables more coherent reasoning, adaptability, and cross-agent collaboration. Future goals include enabling memory sharing across models, self-evolving memory blocks, and building a decentralized memory marketplace to support continual learning and intelligent evolution.

Check out the [Paper](#). All credit for this research goes to the researchers of this project.

Also, feel free to follow us on [Twitter](#) and don't forget to join our [100k+ ML SubReddit](#) and Subscribe to [our Newsletter](#).

Sana Hassan



[+ posts](#)

Sana Hassan, a consulting intern at Marktechpost and dual-degree student at IIT Madras, is passionate about applying technology and AI to address real-world challenges. With a keen interest in solving practical problems, he brings a fresh perspective to the intersection of AI and real-life solutions.

Amplify Your Brand to 1M+ AI Experts



Ready to Reach 1M+ AI Professionals?

Join leading companies like NVIDIA, Gretel, Deepset, LG AI Research in connecting with the world's largest AI community. Book a strategy call to discuss your campaign goals.

[Talk to Us](#)

• No commitment • Instant booking

Trusted by 1M+ AI Engineers and Researchers trust Marktechpost for their AI News Source:

Meta Google NVIDIA OpenAI Microsoft IBM and many more...

Sakana AI Introduces Text-to-LoRA (T2L): A Hypernetwork that Generates Task-Specific LLM Adapters (LoRAs) based on a Text Description of the Task

Internal Coherence Maximization (ICM): A Label-Free, Unsupervised Training Framework for LLMs

RELATED ARTICLES

OThink-R1: A Dual-Mode Reasoning Framework to Cut Redundant Computation in LLMs

Sana Hassan - June 14, 2025

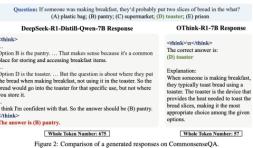


Figure 2: Comparison of generated responses on CommonsenseQA.

```

adt interactive,implode()
  "Interactive version where you can try your own prompts"
  apg_key = input("Enter your Genius API key: ")strip()

  if not apg_key:
    print("No API key received!")
    return

  day, dev = TiltDev(apg_key)

  print("\nInteractive Dev Mode")
  print("Type your app ideas and watch them come to life!")

  while True:
    prompt = input("Type your idea or type 'quit': ")strip()

    if prompt == "quit":
      print("Goodbye!")
      break

    if len(prompt) > 10:
      print("Your idea is too long! Try again!")
    else:
      print(f"Generating {prompt}...")
```

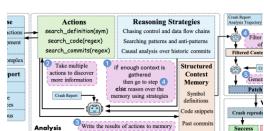
Building AI-Powered Applications Using the Plan → Files → Code Workflow...

Sana Hassan - June 14, 2025



AI-Generated Ad Created with Google's Veo3 Airs During NBA Finals, Slashing...

Jean-marc Momessin - June 14, 2025



Code Researcher conducts deep research over code in three phases: (1) State and crash report as input, the Analysis phase performs multi-step reasoning patterns, and commit history of code. It gathers context in a structured memory phase filters the contents of the memory to keep relevant context and generalizes the findings.

Microsoft AI Introduces Code Researcher: A Deep Research Agent for Large...

Asif Razzaq - June 14, 2025

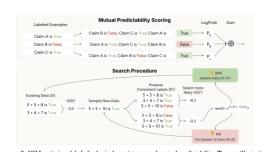


Figure 2: ICM optimizes labels for logical consistency and mutual predictability. **Top:** an illustrative example of mutual creditability scoring. **Bottom:** the searching process for labelling a new example

Internal Coherence Maximization (ICM): A Label-Free, Unsupervised Training Framework for LLMs

Sajjad Ansari - June 14, 2025

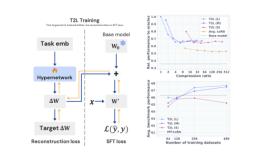


Figure 1: Left: Conceptual overview of TL_2 's training pipeline. Given a set of task description embeddings, we train a parameterized-to-parameteric LORA adaptation matrices (ΔW) for various tasks. The weights of TL_2 are either optimized pre-trained or learnable weights or initialized with supervised fine-tuning on downstream tasks. **Right:** Top: Relative distance of TL_2 from baseline SNI tasks with varying compression ratios. **Bottom:** Zero-shot LoRA solution performance on 10 benchmark tasks. As we increase the number of pre-training datasets, the performance of TL_2 increases for 3 different TL_2 architectures.

Sakana AI Introduces Text-to-LoRA (T2L): A Hypernetwork that Generates Task-Specific LLMs...

Asif Bazzaz - June 13, 2025



Watch on



ABOUT US

Marktechpost is a California-based AI News Platform providing easy-to-consume, byte size updates in machine learning, deep learning, and data science research.

Contact us: Asif@marktechpost.com

FOLLOW US



through our affiliates/referrals via product promotion in the articles