

# Generative Model of Memory Construction and Consolidation

## Summary

The paper explores how episodic memories, which combine unique sensory details with schema-based predictions, are constructed and consolidated within the brain. It proposes a computational model where hippocampal replay, which involves reactivation of neural patterns, trains generative networks like variational autoencoders (VAEs). These networks can reconstruct sensory experiences using latent variable representations, distributed across areas such as the entorhinal cortex (EC), medial prefrontal cortex (mPFC), and anterolateral temporal lobe (aITL). The model attempts to explain memory consolidation, imagination, and future episodic thinking by combining information from both hippocampal and neocortical systems, efficiently managing storage capacity by offloading familiar elements to schema-based systems while retaining new or unpredictable information in the hippocampus.

The generative model captures statistical regularities, which are gradually transferred from the hippocampus to neocortical structures. This process helps form semantic memories, relational inferences, and schema-based memory distortions. The model also explains phenomena like boundary extension, where recalled memories are not exact replicas but reconstructions that may include generalized or prototypical features.

## Research Question

The paper addresses the scientific question: *How are episodic memories constructed, consolidated, and reconstructed, and what neural mechanisms underlie these processes?* Specifically, it seeks to understand how the brain manages the trade-off between encoding detailed sensory experiences and efficiently storing generalized knowledge. The research also aims to explain the neural basis for related cognitive functions, such as imagination, semantic memory, and schema-based distortions.

## Proposed Solution

The authors propose a model where the hippocampus initially encodes sensory experiences in an autoassociative network. During replay, these experiences are used to train generative networks that capture the statistical structure of the stored events. Over time, these generative models enable the brain to recall, reconstruct, or even simulate sensory experiences by retrieving latent variable representations stored across the EC, mPFC, and aITL.

The model operates under the following framework:

1. Initial Encoding: The hippocampus rapidly encodes events using an autoassociative memory network.
2. Hippocampal Replay and Generative Networks: During periods of rest, the hippocampal network reactivates stored experiences, training the generative networks to capture the statistical properties of these events. This process is akin to

training a VAE, where the hippocampus serves as a "teacher" and the generative network as a "student."

3. Consolidation: As the generative models learn, the dependence on the hippocampal network decreases. The generative models can now reconstruct experiences based on learned schemas, allowing for efficient memory storage.
4. Schema-Based Memory Distortions: The model explains why memories might include schema-based distortions. As generative models rely more on statistical regularities, memories can become more generalized and less detailed over time, leading to phenomena like boundary extension.
5. Combination of Sensory and Conceptual Elements: The model distinguishes between predictable, conceptual components and unpredictable, sensory details, optimizing memory storage by efficiently encoding what is necessary.
6. Memory Distortions and Boundary Effects: Schema-based distortions are modeled as the generative network distorts recalled memories toward more prototypical representations, while boundary effects (extension/contraction) are observed in artificially zoomed images.
7. Semantic Memory and Imagination: Over time, as the generative networks consolidate information, they support semantic memory, imagination, and the ability to infer new information through relational inference.

## Potential Improvements

1. Enhanced Biological Plausibility: While the model uses VAEs, which are effective computationally, the learning mechanisms (such as backpropagation) are not entirely biologically plausible. Developing more biologically realistic learning rules could enhance the model's applicability to real neural processes.
2. Incorporating Sequential Memory Structures: The model focuses mainly on static memories. Adding mechanisms to handle sequential or temporal data (e.g., recalling events in a specific order) would provide a more comprehensive understanding of memory processes, especially for applications like language or complex tasks.
3. Reducing Memory Distortions: The model's reliance on schemas can lead to distortions. Introducing mechanisms that account for error correction or contextual cues during recall might help reduce these distortions and provide a more accurate reconstruction of memories.
4. Better Handling of Forgetting: The model does not address how new memories might interfere with old ones (catastrophic forgetting). Implementing mechanisms like continual learning or generative replay could stabilize latent variable representations and prevent interference.
5. Realistic Biologically-Inspired Learning Mechanisms: The generative networks currently rely on backpropagation, which is biologically implausible. Exploring biologically-inspired alternatives (e.g., Hebbian learning, predictive coding) might provide a more realistic simulation of memory consolidation.
6. Multimodal Integration: Extending the model to integrate multiple modalities (e.g., audio, visual, language) could provide a more realistic representation of how the brain processes and recalls complex, multimodal experiences.
7. Explicit Modeling of Sleep and Memory Replay: While the model assumes replay during rest, it could be enhanced by simulating different stages of sleep and their role

in memory consolidation, replay, and stabilization. This could align with existing neuroscience findings on how specific sleep stages contribute to memory processing.

8. Extending to More Complex Data: While the paper demonstrates the model using datasets like MNIST and Shapes3D, applying the approach to more complex, multimodal datasets (e.g., audio-visual or sensory-motor experiences) could test the robustness and versatility of the proposed framework.
9. Parallel Generative Networks for Multi-Task Learning: The paper hints at training multiple generative networks optimized for different tasks. Future research could explore this idea by setting up networks specialized in various cognitive functions, enhancing the understanding of how the brain might prioritize and manage different types of memory tasks simultaneously.

This paper provides a significant step forward in understanding memory consolidation by bridging the gap between computational models and neural processes, offering insights into not just memory but related functions such as imagination and semantic memory. Further exploration and improvements in the areas mentioned above could make this model an even more robust tool for understanding the intricacies of human cognition.

## **Terminology used**

**Episodic Memory:** A type of long-term memory that involves the recollection of specific events, situations, and experiences, often including details about time and place. For example, recalling a birthday party or a vacation.

**Semantic Memory:** Memory of general knowledge and facts that are not tied to a specific time or place. This includes information like knowing the capital of a country or the meaning of a word.

**Hippocampus:** A critical brain structure involved in the formation of new memories, especially episodic memories. It helps encode, store, and retrieve information by linking different sensory and conceptual elements together.

**Autoassociative Network:** A type of neural network that can store and retrieve patterns of data. When a part of a pattern is provided, it can retrieve the full pattern based on associations formed during training. In the brain, this is thought to be similar to how the hippocampus recalls memories.

**Hippocampal Replay:** The process by which the brain reactivates patterns of neural activity associated with past experiences, often during periods of rest or sleep. This reactivation is believed to play a role in consolidating memories and transferring them to other parts of the brain.

**Consolidation:** The process by which short-term memories are stabilized and transferred into long-term storage. This involves the integration of memories into the neocortex, making them more resistant to forgetting.

**Neocortex:** The part of the brain responsible for higher-order functions, including sensory perception, motor commands, spatial reasoning, and language. In the context of memory, the neocortex is where long-term memories are stored once they are consolidated.

**Schema:** A cognitive framework or pattern that helps organize and interpret information. Schemas are built from previous experiences and help predict and reconstruct future experiences. For example, having a schema for an "office" allows one to expect the presence of desks, chairs, and computers without having to see them explicitly.

**Generative Model:** A computational model that can learn the underlying probability distributions of data, allowing it to generate new data points that are consistent with what it has learned. In this paper, generative models help reconstruct and simulate past experiences based on stored memories.

**Variational Autoencoder (VAE):** A type of generative neural network that learns to encode input data (like images) into a lower-dimensional space (latent variables) and then decode it back into the original form. This process allows the network to generate new samples by manipulating the latent variables.

**Latent Variables:** Hidden variables that are not directly observed but are inferred by the model. In the context of a VAE, latent variables represent compressed, underlying features of data that can be used to generate new or reconstructed versions of the data.

**Entorhinal Cortex (EC):** A region of the brain that acts as a gateway between the hippocampus and the neocortex. It plays a crucial role in memory and navigation, particularly in forming spatial and contextual representations.

**Medial Prefrontal Cortex (mPFC):** A part of the prefrontal cortex involved in decision-making, social behavior, and the integration of information across different contexts. In memory processes, it is associated with encoding schemas and integrating semantic information.

**Anterolateral Temporal Lobe (aTL):** A region of the brain involved in processing semantic information, including the meanings of words and objects. It plays a role in storing and retrieving conceptual knowledge.

**Schema-Based Distortion:** A phenomenon where memories are not recalled exactly as they happened but are influenced by existing schemas. This can lead to the recall of prototypical or generalized versions of events rather than precise details, explaining why our memories may sometimes differ from reality.

**Boundary Extension:** A specific type of memory distortion where people remember a wider field of view than they actually saw. For example, if someone sees a photograph of a scene, they might later recall seeing more of the background than was actually visible.

**Teacher–Student Learning:** A concept where one network (the "teacher") trains another (the "student") by providing examples. In the context of the paper, the hippocampus acts as the teacher, and the generative network (like a VAE) learns to reconstruct memories based on examples replayed by the hippocampus.

**Prediction Error:** The difference between expected and actual outcomes. In memory processing, the generative network minimizes prediction error by learning to better reconstruct events, reducing the need for detailed hippocampal encoding.

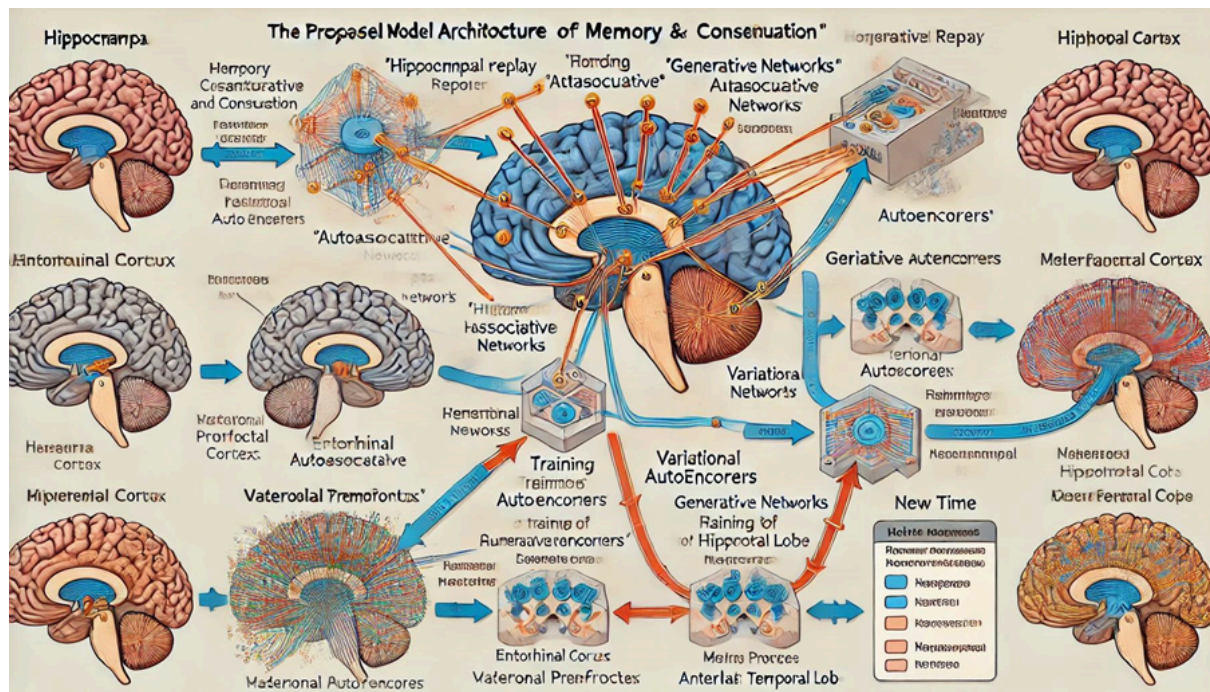
**Pattern Completion:** The ability of a neural network to retrieve a complete memory or pattern based on partial input. For instance, if you only see part of an object, your brain can use pattern completion to recognize and recall the whole object.

**Multiple Trace Theory:** A theory suggesting that every time a memory is recalled, a new trace is created in the brain, allowing multiple versions of a memory to coexist. This contrasts with the idea that there is a single, unified memory trace.

**Complementary Learning Systems (CLS):** A theory that proposes two learning systems in the brain—one for rapid, one-shot learning (hippocampus) and another for gradual, cumulative learning (neocortex). The interplay between these systems allows for both the quick capture of new information and the integration of that information over time.

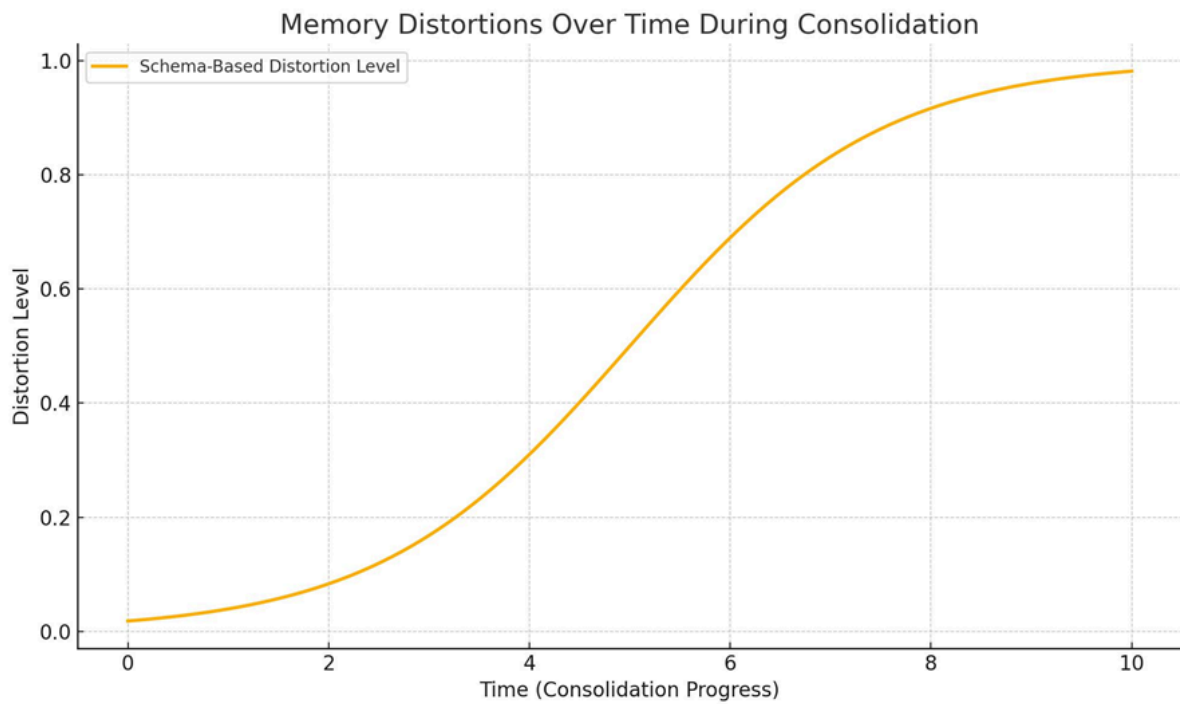
## Visualizations

Diagram 1: Proposed Model Architecture (source - generated midjourney)



This diagram illustrates the architecture of the proposed model for memory construction and consolidation. At the center is the hippocampus, which encodes new events using an autoassociative network. Arrows indicate the process of "hippocampal replay," where these encoded memories are reactivated and used to train generative models (like variational autoencoders). These generative networks are distributed across different brain regions, such as the entorhinal cortex, medial prefrontal cortex, and anterolateral temporal lobe. Over time, these networks take over the task of reconstructing memories, reducing the need for direct hippocampal involvement, showing a gradual shift from detailed sensory encoding to more abstract, schema-based representations.

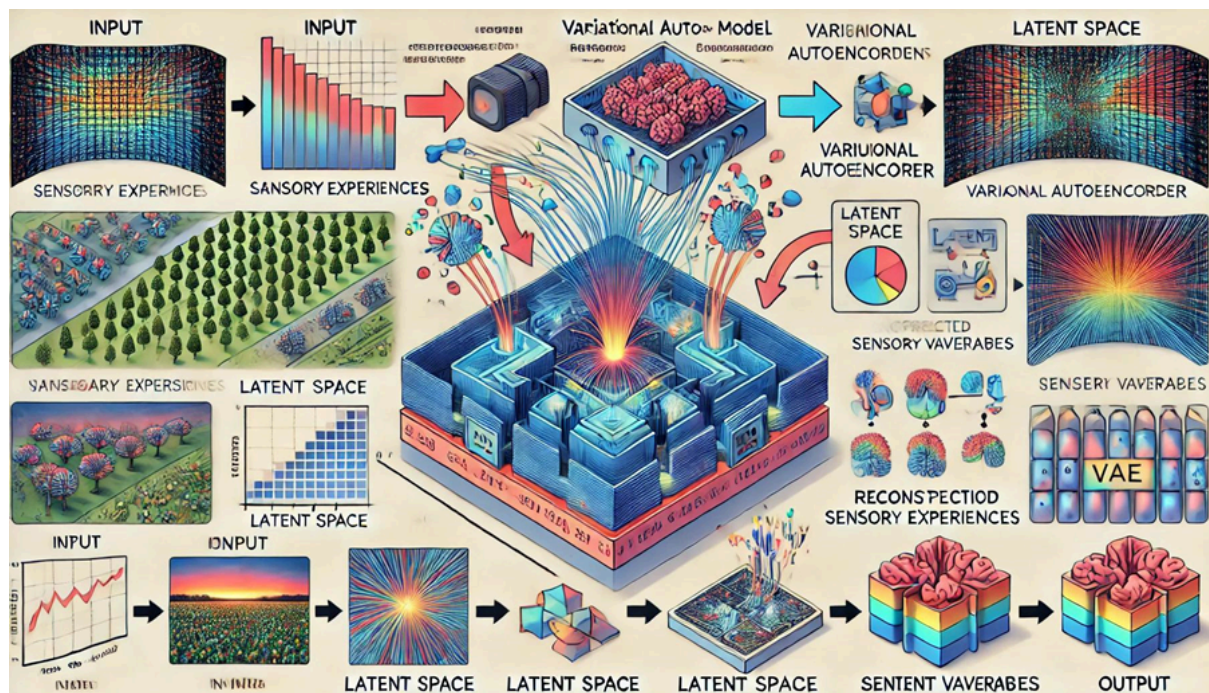
Diagram 2: Memory Distortions Over Time (source - generated midjourney)



This graph depicts the increase in schema-based distortions in memory over time as the consolidation process progresses. The line represents how reliance on generative models for memory recall leads to more generalized, prototypical reconstructions, rather than precise replicas of the original experiences. The curve shows that as memories become more consolidated, they are increasingly influenced by existing schemas, which can lead to distortions, reflecting the brain's tendency to optimize storage by generalizing common elements.



Diagram 3: Generative Model's Learning Process (source - generated midjourney)



This visualization explains how a generative model (like a variational autoencoder) learns during memory consolidation. The process is divided into three stages: input, latent space, and output. The left side represents the input phase, where raw sensory experiences are provided to the model. In the middle, the latent space compresses and encodes these inputs into latent variables, capturing the essential features of the experience. Finally, the right side shows the output phase, where the latent variables are decoded back into reconstructed sensory experiences. This diagram highlights the role of latent variables in efficiently reconstructing and simulating memories, showing how the model learns to generate experiences from compressed representations.