

A neural model of short-term and intermediate-term memory using neuro-symbolic representations for semantic and spatial memory tasks

Jakeb Chouinard (jakeb.chouinard@uwaterloo.ca)^{1,2} & Chris Eliasmith^{1,2,3}

¹Department of Systems Design Engineering, University of Waterloo

²Centre for Theoretical Neuroscience, University of Waterloo

³Department of Philosophy, University of Waterloo

Abstract

Memory is a vital part of cognition. Its many parts determine how one learns, plans, and ultimately experiences the world around them. As a bridge between short-term memory and long-term memory, intermediate-term memory allows for semi-stable memories to be rapidly encoded through short-term synaptic plasticity. We present a model of short-term and intermediate-term memory based around Semantic Pointers and realize the model in a spiking neural network. Unlike past work, we characterize short-term memory as a dynamical system and realize intermediate-term memory as implicit association learning. Further, we demonstrate our model’s ability to generalize across domains by achieving human-like results on semantic and spatial memory tasks with minimal parameter changes. By examining Cohen’s h values first within and then between tasks, we find an overall small difference (Mean: $|\hat{h}| = 0.12 \pm 0.06$, Median: $|\hat{h}| = 0.16$), suggesting strong similarity between model and human data across all three tasks.

Keywords: memory; spiking neural networks; vector symbolic algebra; neural engineering framework; spatial semantic pointers; neuro-symbolic representation

Introduction

Memory plays a central role in a wide variety of cognitive tasks. Humans and animals alike depend on their memories of recent and remote events to plan their actions and anticipate the results of interactions (Kandel, 2001; Mullally & Maguire, 2014). Memory is typically divided into sensory, short-term, and long-term varieties that implicate sensory cortices, the prefrontal cortex, and medial temporal lobe and cortical structures respectively (Camina & Güell, 2017). Sensory memory is the shortest-lasting memory, decaying rapidly once the stimulus is removed. Short-term memory (STM) persists for tens of seconds but can be stretched to last longer through overt or covert rehearsal mechanisms. Lastly, long-term memory (LTM) can persist for days, years, and even decades without any form of conscious rehearsal. STM is thought to be encoded in neuron population activities across the brain (Kamiński et al., 2017), while LTM is thought to be developed through long-term synaptic plasticity (LTSP) during consolidation processes (Goto, 2022).

Recent studies suggest the existence of intermediate-term memory (ITM) mediated by short-term synaptic plasticity (STSP) in the hippocampus (Kamiński, 2017; Mongillo et al., 2008). ITM provides an alternate account of how the brain can retain recent pieces of information without constant

representation in populations’ activities.

In this work, we propose a model of short-term and intermediate-term memory realized in a spiking neural network. Furthermore, we present this model as a novel symbolic dynamic system that leverages symbol-like representations and short-term synaptic plasticity to provide a neurally plausible account of STM and ITM in both semantic and spatial memory tasks. We compare our model to human performance on three different experiments—serial recall, free recall, and spatial recall—and demonstrate its applicability across domains.

Background

Semantic Pointer Architecture

The Semantic Pointer Architecture (SPA) proposes the use of high-dimensional vectors carrying semantic meaning, or Semantic Pointers (SPs; $\psi \in \mathbb{R}^d$ s.t. $\|\psi\| = 1$ where $\|\cdot\|$ is the ℓ_2 norm) to construct models of cognition (Eliasmith, 2013). Using the Neural Engineering Framework (NEF), these models can be mapped onto spiking neural networks (SNNs) to test and analyze neuro-computational hypotheses (Eliasmith & Anderson, 2003). To impose structure and allow for compositional representations in these models, the SPA uses Holographic Reduced Representations (HRRs)—a form of Vector Symbolic Algebra (VSA; Plate, 1995).

In VSAs, *similarity measures* (\odot) are used to evaluate semantic—or feature—similarity between two representations. *Bundling* ($+$) combines a pair of vectors to produce a single vector that is similar to both of its constituent vectors. Unlike bundling, *binding* (\otimes) combines a pair of vectors to produce a single vector that is similar to neither of its constituent vectors. Lastly, *inversion* ($\tilde{\cdot}$) allows for a single vector to be recovered from a bound pair by inverting the other vector (i.e., $(\psi_A \otimes \psi_B) \otimes \psi_B \tilde{\approx} \psi_A$).

The HRR VSA in the SPA realizes similarity measures as cosine similarity, bundling as vector addition, binding as circular convolution, and inversion as the reversal of a vector’s 2nd to d^{th} entries. These operations can be written as:

$$\begin{aligned}\psi_A \odot \psi_B &= \frac{\langle \psi_A, \psi_B \rangle}{\|\psi_A\| \|\psi_B\|} \\ \psi_C &= \psi_A + \psi_B \\ \psi_D &= \psi_A \otimes \psi_B \\ \psi_A \tilde{} &= [\psi_{A,1}, \psi_{A,d}, \psi_{A,d-1}, \dots, \psi_{A,2}]\end{aligned}$$

Most commonly, SPs are generated by randomly sampling a vector’s entries from some distribution and normalizing the vector such that it is a point on the surface of the hypersphere. As a consequence of this, it is possible to accidentally generate highly-similar—or even opposite—vectors. To prevent this, newly generated SPs can be compared to the previously generated SPs to ensure that they are sufficiently orthogonal: $\langle \psi_{N+1}, \psi_i \rangle \approx 0 \forall i \in \{1, 2, \dots, N\}$. Since there is no guarantee that random SPs are unitary, circular convolution can result in an SP with non-unit length, potentially causing excessive growth as an SP is repeatedly bound.

As an extension to the SPA, Spatial Semantic Pointers (SSPs) allow for representation of continuous feature spaces as trajectories on the surface of a hypertorus. SSPs can be used in the same way as SPs within the SPA, allowing for discrete and continuous feature representations within compositional structures (Komer et al., 2019). For some n -dimensional feature space, $\mathbf{x} \in \mathbb{R}^n$, a d -dimensional SSP encoding, $\phi(\mathbf{x}) \in \mathbb{R}^d$, is defined as:

$$\phi(\mathbf{x}) = \mathcal{F}^{-1} \{e^{jA\mathbf{x}}\}$$

where $A \in \mathbb{R}^{d \times n}$ is a conjugate-symmetric phase matrix sampled from some probability distribution over the n -dimensional frequency space; $\mathcal{F}^{-1}\{\cdot\}$ is the inverse Fourier transform; and j is the imaginary unit.

SSP embeddings can also be used to generate unitary SPs by mapping an n -dimensional feature space to n SPs such that all SPs are approximately orthogonal. As an example, two SPs, $\phi[a]$ and $\phi[b]$, can be generated as $\phi[a] = \mathcal{F}^{-1} \{e^{jA[1,0]^T}\}$ and $\phi[b] = \mathcal{F}^{-1} \{e^{jA[0,1]^T}\}$ respectively, where A is sampled such that $\langle \phi[a], \phi[b] \rangle \approx 0$. This solves the problem of non-unitary vectors from traditional random SP generation. Because circular convolution in the HRR VSA can be rewritten as a Fourier-space Hadamard product, SSP binding is equivalent to Fourier phase superposition for SSPs with different phase matrices and feature space superposition for SSPs with the same phase matrix:

$$\begin{aligned} \phi_1(\mathbf{x}_1) \otimes \phi_2(\mathbf{x}_2) &= \mathcal{F}^{-1} \{e^{jA_1\mathbf{x}_1 + jA_2\mathbf{x}_2}\} \\ \phi(\mathbf{x}_1) \otimes \phi(\mathbf{x}_2) &= \phi(\mathbf{x}_1 + \mathbf{x}_2) \end{aligned}$$

Task Representation

Working memory is commonly tested through information recall experiments. In many experiments, participants are presented some set of items and are asked to recall them. The presentation of items can be from a list of words presented on a screen or as items that the participant sees in an environment. Typically, participants are asked to recall the items either in the order they were originally shown (serial recall) or in the order of their choice (free recall); however, they can also be asked for information regarding an item’s position in a list or its physical location in the environment. In this work, we compare model performance to human performance on three different working memory tasks in order to demonstrate the generality of the model’s underlying principles. Specifically, we have the model perform a serial recall task in which

participants must recall 10 items in order of their appearance (Jahnke, 1968); a free recall task in which participants must recall 12 items while orating to prevent rehearsal (Howard & Kahana, 1999); and a spatial recall task in which participants must recall the positions of 2-4 items encountered in an augmented reality (AR) environment following a chase-based distractor task (Maidenbaum et al., 2025).

In the serial and free recall tasks, we represent items within a list of length N as a set of orthogonal, unitary SPs, Φ_V , and use an SSP space to generate a set of each item’s temporal embedding to capture its serial order, Φ_T :

$$\begin{aligned} \Phi_V &= \{\phi_V[i] \mid \langle \phi_V[i], \phi_V[k] \rangle \approx 0 \forall i \neq k; i, k \in \mathbb{N}_{\leq N}\} \\ \Phi_T &= \left\{ \phi_T[i] = \mathcal{F}^{-1} \{e^{jA_T(i+b)}\} \mid i \in \mathbb{N}_{\leq N} \right\} \end{aligned}$$

where $b \gg 1$ to prevent $\phi_T[1]$ from being similar to the HRR identity vector ($\phi_I = \mathcal{F}^{-1} \{e^{\phi}\}$). Individuals are more likely to recall items near the most recently recalled item in a list; consequently, it is desirable to increase similarity between adjacent temporal embeddings. To realize this, we sample entries of A_T such that $\phi_T[i] \cdot \phi_T[i \pm 1] \approx 0.25$.

For the spatial recall task, we represent the M landmarks presented at each point in the environment as a set of orthogonal, unitary SPs, Φ_L , formed by binding their feature embeddings: colour, c , and form, f . We also use an SSP space to embed the navigable space, Φ_S , and consequently the set of positions for each landmark:

$$\begin{aligned} \Phi_L &= \{\phi_L[i] = \phi_C[c_i] \otimes \phi_F[f_i] \mid i \in \mathbb{N}_{\leq M}\} \\ \Phi_S &= \{\phi_S(\mathbf{x}) = \mathcal{F}^{-1} \{e^{jH\mathbf{x}}\} \mid \mathbf{x} \in \mathbb{R}^2\} \end{aligned}$$

where $H \in \mathbb{R}^{d \times 2}$ is a matrix formed by combining grid cell-like embeddings to form place cells as described in Dumont and Eliasmith (2020). These HexSSPs allow for biologically plausible spatial representation using VSAs that are used to facilitate self-orienting during the spatial recall task.

The Model

In this description of the model, we limit ourselves to specifying the relevant state space dynamics for capturing the effects of interest. We use the Neural Engineering Framework to embed these dynamics into a spiking neural network that we simulate, analyze, and report on in the results section. The size of the network for all tasks was just over 1 million leaky integrate-and-fire neurons.

Short-term Memory

In creating a model of short-term memory (STM; Figure 1), we consider the essential characteristics of activity-based memory and the limitations of neuron populations. In working memory experiments, STM is typically implicated in the “recency” effect of recall; that is, it plays a significant role in recalling recent information (Atkinson & Shiffrin, 1968). This is especially evident in delayed recall tasks using difficult distractor tasks where later items are significantly worse remembered than earlier items, suggesting that distractor tasks

may effectively “erase” STM. Further, STM is typically considered to be finite in its capacity and gated by attentional mechanisms (Cowan, 2001; Reeves & Sperling, 1986).

To realize these effects, we propose a dynamical model of STM that leverages finite neural resources and vector projections. The dynamics of our model are based on “loading” an input vector, $\phi(t) \in \mathbb{R}^d$, into memory, $\mathbf{m} \in \mathbb{R}^d$ until it is well represented. This behaviour can be captured by the following dynamical equation:

$$\frac{d\mathbf{m}}{dt} = \langle \phi(t), \phi(t) - \mathbf{m} \rangle \phi(t)$$

In this system, the absence of the input vector in memory is evaluated by taking the dot product between the input vector and the difference between the input vector and the memory vector. This returns the component length of the difference projected onto the input vector; a value proportional to the magnitude of its absence. This value can be used to scale the input vector being added to the memory such that changes in the value of \mathbf{m} occur only along the direction of $\phi(t)$. We also scale this value using a gain parameter, $g_m \in \mathbb{R}^+$, to ensure entire vectors can be introduced into memory in a timely manner. Lastly, we introduce a small, time-based decay parameter, $\gamma \in \mathbb{R}^-$, to account for the natural drift and decay of activity-based memory. With these changes, the full dynamic system is written as:

$$\frac{d\mathbf{m}}{dt} = -\gamma\mathbf{m} + g_m \langle \phi(t), \phi(t) - \mathbf{m} \rangle \phi(t)$$

To prevent the size of \mathbf{m} from growing indefinitely and to realize finite capacity in our memory, we leverage a known feature of spiking neurons: the saturation of spiking rates due to large inputs. In the NEF, neurons in a given population exhibit changes in spiking rates over some portion of a finite representational space with bounds based on the population’s radius, r ; typically, $r = 1$. This allows populations to represent any SSP or combination of SSPs so long as its magnitude is less than or equal to 1. Since a single SSP has a magnitude of 1, bundling multiple SSPs will result in a vector with a magnitude significantly greater than 1. Firing rate saturation due to the population’s radius will force the vector back towards a magnitude of 1; a process called soft-normalization. The ability of the memory population to represent N orthogonal SSPs simultaneously—or its capacity—can therefore be understood in terms of the representational radius:

$$r = \mathbb{E} \left[\left\| \sum_{i=1}^N \phi_i \right\| \right] = \mathbb{E} \left[\sqrt{\left\langle \sum_{i=1}^N \phi_i, \sum_{i=1}^N \phi_i \right\rangle} \right] \approx \sqrt{N\|\phi\|^2} = \sqrt{N}$$

In our experiments, we maintain the radius of the STM and the ITM_p at a value that can contain 4 “chunks” of information—consistent with theorized STM capacity (Cowan, 2001)—such that $r = \sqrt{4\|\mathbf{c}\|^2}$ where $\mathbf{c} \in \mathbb{R}^d$ is a chunk. Once STM exceeds 4 chunks, older chunks degrade due to the soft-normalization process while chunks being presented are maintained. Unlike previous STM models, this allows us to realize forgetting

without depending on careful scaling or noising of items in memory (Choo, 2010; Reimann, 2025). For the serial and free recall tasks, a chunk for the STM, \mathbf{c}_m , is the bundle of an item with the bound product of the item with its temporal embedding such that $\|\mathbf{c}_m\| = \sqrt{2}$ and $r_m = \sqrt{8}$. For the spatial recall task, \mathbf{c}_m is the bound product of a landmark’s feature SP and a position SSP such that $\|\mathbf{c}_m\| = 1$ and $r_m = \sqrt{4}$. For the ITM_p in all experiments, a chunk is a single SSP (either a temporal embedding or a position SSP; $r_p = \sqrt{4}$).

Intermediate-term Memory

Unlike STM, intermediate-term memory (ITM) is typically implicated in the “primacy” effect of recall (Talmi et al., 2005). The persistence of the ability to recall early items in a list when distracted highlights ITM’s role in memory; however, this primacy bias is noticeably impacted when rehearsal mechanisms are sufficiently interrupted (Marshall & Werder, 1972). Our model of ITM (Figure 1) is based on STSP. A variety of Hebbian learning paradigms allow for STSP learning in SNNs; we use Associative Matrix Learning (AML; Gosmann and Eliasmith, 2021). AML reduces catastrophic forgetting by updating the weights between neuron populations in such a way that they avoid overwriting previously learned mappings.

The AML achieves this by modifying the weights between two populations of neurons to learn an outer-product matrix that facilitates input-to-target mappings. That is, for N items, the weight matrix approximates the auto-association matrix, L , as the sum of outer products:

$$L = \sum_{i=1}^N (\phi_V[i] + \phi_T[i])(\phi_V[i] + \phi_T[i])^\top$$

Since bundling preserves similarity to its constituent vectors, the learned outer product matrix transformation allow for recovery of an approximation of a bundle when only one of its constituent vectors is presented (i.e., $L\phi_V \approx \phi_V + \phi_T$). Unlike previous association-based models, such as the Context-Unified Encoding (Gosmann & Eliasmith, 2021) and the Context Maintenance and Retrieval models (Polyn et al., 2009), we learn mappings that implicitly associate context and feature vectors rather than explicit feature-to-context-to-feature mappings. Further, we perform all learning online and test over continuous spaces, unlike the graph-based tasks performed by the Tolman-Eichenbaum Machine (Whittington et al., 2020).

To query the auto-associative ITM network (ITM_L), we introduce a network that integrates temporal embeddings. By duplicating the projection-based memory network used in STM, reducing its gain parameter, and introducing a non-negative filter to the vector projection, we implement a network that can integrate earlier embeddings preferentially and resist newer items being added into the memory. This integrating ITM network (ITM_p) realizes an activity-based memory vector, \mathbf{p} , based on the dynamical system:

$$\frac{d\mathbf{p}}{dt} = -\gamma\mathbf{p} + g_p f(\langle \phi(t), \phi(t) - \mathbf{p} \rangle) \phi(t)$$

where $f(x) = x$ if $x \geq 0$, otherwise $f(x) = 0$.

In all tasks, the ITM_p is used as an input to the ITM_L during recall. ITM_p serves to capture the timing of the task as a whole, with the balance of integrating new embeddings and reinforcing early embeddings determined by g_p . An increased g_p reflects increased preference to new embeddings over early embeddings. ITM_p 's role is especially important in the free recall task, where there is no external query vector provided.

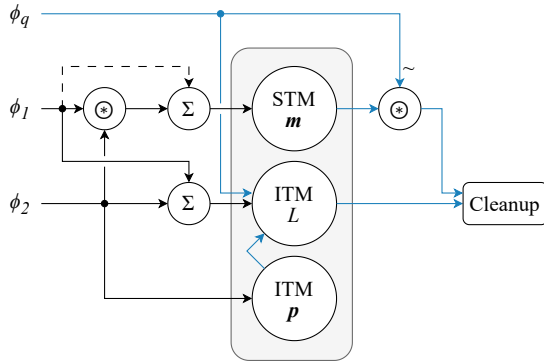


Figure 1: A high-level diagram of the memory model; $\phi_1 \in \Phi_V \cup \Phi_L$, $\phi_2 \in \Phi_T \cup \Phi_S$, and $\phi_q \in \Phi_T \cup \Phi_L$. Black lines show presentation-time model interactions with dashed lines being specific to serial and free recall tasks. Blue lines show recall-time model interactions.

Cleanup

To facilitate recall in each task, we use a neural cleanup network. A neural cleanup network is an auto-associative network that helps reduce noise by ‘cleaning up’ input vectors to known good representations. For serial and free recall tasks, Independent Accumulator networks (IAs) with normally distributed noise, $\sigma_{IA} = 0.009$, and minimum evidence thresholds, η_{IA} , are used as cleanups for temporal embeddings and list items (Gosmann et al., 2017). In the IAs, minimum evidence corresponds to the minimum dot product between an input and a ground truth vector to allow recall. The temporal embedding IA is inhibited in serial recall, as the model instead cycles back through self-generated query vectors. For free recall, temporal embedding recall is an input to ITM_L , letting the network use a remembered temporal embedding to query for its corresponding item (or vice versa).

In the spatial recall task, the model is queried with the SP representing a landmark to acquire its position in the navigable space as an SSP. This allocentric SSP is cleaned up using an argmax of dot products over a 150×150 grid of place cell-like representations. To report the location of the landmark relative to the agent, the agent binds the allocentric SSP with the inverse of its own positional SSP, ultimately outputting an estimate of the egocentric SSP. Once the simulation is complete, these egocentric SSPs are bound with ground truth positional SSPs and decoded using the same grid-based cleanup.

In all cases, the contributions of the STM to recall, of the ITM_p to the ITM_L , and of the query SP are weighted by

parameters θ_m , θ_p and θ_q respectively. For the serial and free recall tasks, these values balance inputs to recall such that their magnitudes are not exceedingly above 1. For the spatial recall task, θ_m and θ_p are decreased to reflect a participant’s ability to inhibit recollection of non-landmark positions.

Results

To capture human-like limitations, there are no structural changes and limited parameter changes made to the model between tasks (see Table 1). Initial parameter values were selected to achieve model components’ intended behaviours—such as loading an SP into the STM or learning a new ITM_L mapping within 300ms—and were incrementally tuned across tasks based on the model’s results. Differences in parameter values between tasks reflect differences in experimental paradigms: decreased relative contribution of the STM and the ITM_p in spatial recall due to the distractor task; reduction of the input gain for the ITM_p in free recall due to rehearsal interruption; an increase to SSP dimensionality in spatial recall due to increased confusability in SSPs; an increase to the learning rate in the ITM_L due to increased dimensionality; and increased minimum evidence in the IAs for serial recall compared to free recall due to desire for positional accuracy.

We calculate 95% confidence intervals (CIs) using the Clopper-Pearson Exact Method (Clopper & Pearson, 1934) to determine if there are significant differences in results’ distributions. We also calculate Cohen’s h to quantify differences in response proportions between human and model behaviour and report it to two decimal places (Cohen, 2013). Based on Cohen’s interpretations, we group small effect sizes as $|h| \leq 0.35$, moderate effect sizes as $0.35 < |h| \leq 0.65$, and large effect sizes as $0.65 < |h|$. In doing so, we allow for interpretable comparisons where effect sizes indicate magnitudes of difference between model and human task results.

Serial Recall

As per Jahnke (1968), the model was presented 10 items at a rate of 1 item per second, following which it was tasked with recalling them in order. Temporal embeddings were generated within the model as new items were presented and then reused during recall. The system to distinguish between presentation and recall phases of the task was realized internal to the model, and no runtime interventions were necessary. For 100 randomly initialized trials, we analyzed the probability of recalling an item (Figure 2a) and the probability of an item’s misplacement during recall (Figure 2b)—a transposition error.

Consistent with human data, our model exhibits primacy and recency biases as well as comparable positional accuracy—overlapping CIs for 11 of 11 data points. For serial position, only one point shows a moderate effect size ($|h_9| = 0.40$) while the remainder show small effect sizes ($0.04 < |h_{i \neq 9}| < 0.30$), indicating a small overall difference between model performance and human behavioural data. While Jahnke (1968) does not provide specific transposition error data, we can compare the probability of correct

Table 1: Model and Task Parameters for Different Recall Tasks

Parameter	Serial	Free	Spatial	Description	Reason for Change
d	256	256	338	(S)SP Dimensionality	SSP Confusability
l	10	10	30	AML Learning Rate, ITM_L	Δd
$\gamma_{m,p}$	-0.0228	-0.0228	-0.0228	Natural Decay Parameter, STM & ITM_p	—
g_m	5.0	5.0	5.0	Input Gain, STM	—
g_p	0.2	1.0	0.2	Input Gain, ITM_p	Rehearsal Interruption
$N_{m,p}$	4	4	4	Capacity (chunks), STM & ITM_p	—
θ_m	1.0	1.0	0.1	Contribution to Recall, STM	Distractor Task
θ_p	0.707	0.707	0.1	Contribution to Recall, ITM_p	Distractor Task
θ_q	1.0	1.0	1.0	Contribution to Recall, ϕ_q	—
η_{IA}	0.375	0.30	—	Minimum Evidence, IAs	Serial Position Specificity

placement for recalled items between model and human results. We find no significant difference, and a small effect size ($|h_{T0}| = 0.12$) further suggests human-like positional accuracy in the model. The model’s transposition error results are otherwise qualitatively consistent with a variety of other serial recall experiments where transpositions are unlikely, localized, and slightly biased in the forward direction (Farrell & Lewandowsky, 2004; Henson et al., 1996).

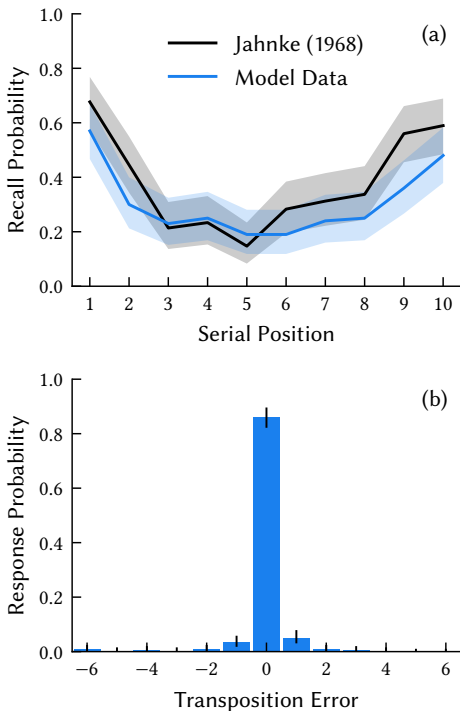


Figure 2: Serial recall task results. (a) Recall probability for list items. (b) Transposition error probability. Shaded regions and error bars denote 95% confidence intervals.

Free Recall

As per Experiment 1 of Howard and Kahana (1999), the model was presented 12 items at a rate of one item per second and

immediately tasked with recalling them in any order. For 100 randomly initialized trials, we analyzed the probability of recalling an item (Figure 3a); the probability of an item being the first recalled item (Figure 3b); and the probability of recalling an item as a function of the previously recalled item’s position (Figure 3c; Cumulative Recall Probability, CRP).

Our model exhibits significant similarity to human behavioural data, showing significant overlap of CIs for 30 of 36 data points. For recall probability, 2 of the 12 items showed significant differences between the model responses and human responses with moderate effect sizes ($|h_{11}| = 0.49$, $|h_{12}| = 0.65$). Despite this, the remainder of list items show small effect sizes ($0.06 < |h_{i \neq 11,12}| < 0.23$). For probability of first recall, no items show significant differences and all show small effect sizes ($0.01 < |h_i| < 0.22$). For CRP, 4 of the 12 show significant differences; however, only one data point has an effect size that isn’t small ($|h_6| = 0.43$; $0.003 < |h_{i \neq 6}| < 0.31$). Qualitatively, our model is highly consistent with the human behavioural data. Our model demonstrates a degraded primacy bias; an intact recency bias; higher probability of recalling terminal items first; and CRP asymmetry typical of human behaviour in free recall tasks (Golomb et al., 2008; Howard & Kahana, 1999).

Spatial Recall

To demonstrate the applicability of the model to other memory tasks, it was embedded into a neural Simultaneous Localization and Mapping (SLAM) agent (Dumont et al., 2023). In each trial, the model constructs a semantic map by forming auto-associative mappings between items and where the agent believes them to be using its estimate of self-position. The semantic map is maintained in the weights of the ITM_L , while the ITM_p integrates landmark positions and the STM remembers recently visited landmarks and their positions.

As per Maidenbaum et al. (2025), the SLAM agent was placed in the centre of a 7×11 m room and followed some provided 2D velocity profile, visiting 4 randomly located landmarks for approximately 1500ms each. Each stationary landmark was visible only during its corresponding 1500ms window. Before beginning recall, the agent followed a moving

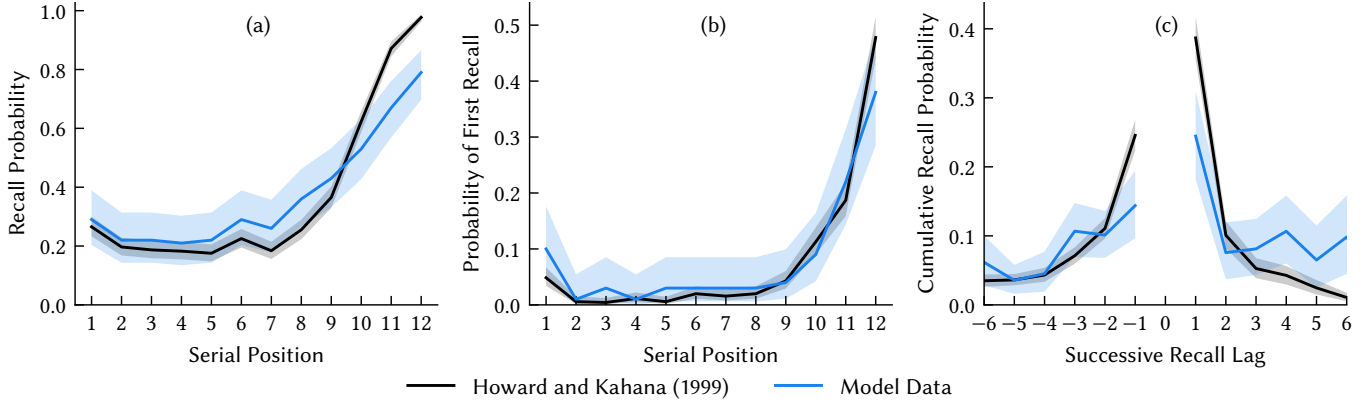


Figure 3: Free recall task results. **(a)** Recall probability for list items. **(b)** Probability of an item to be recalled first. **(c)** Probability of the next recalled item having a given positional lag. Shaded regions denote 95% confidence intervals.

landmark for at least 5m to reflect the experimental distractor task. During this distractor task, ITM_L representations were degraded, and STM and ITM_p activities were compromised. For 100 randomly initialized trials, we analyzed the Mean of the Corrected Error Distance as in the original experiment (MCED; Figure 4). There was no significant difference between the human MCED (0.08 ± 0.01) and our model’s MCED (0.08 ± 0.02)—both of which correspond to an average response in the 92nd percentile of all possible responses. Consistent with the model exhibiting near-human behaviour, there was a near-zero effect size ($|h| = 0.01$).

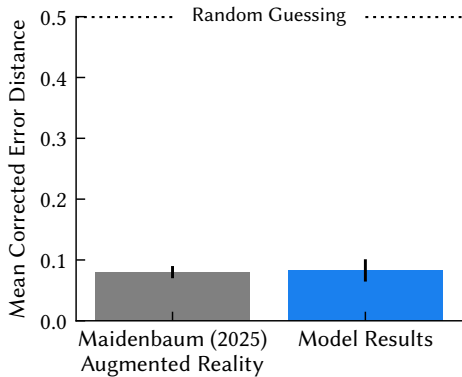


Figure 4: Spatial recall task results. MCED across trials. Error bars denote 95% confidence intervals.

Across all tasks, our model showed significant differences for 6 of 48 data points when compared to human data. Of these 6 data points, 3 showed moderate effect sizes ($0.40 < |h_M| < 0.65$) while the remaining 3 showed small effect sizes. Notably, our model shows small effect sizes for 44 of 48 data points ($0.003 < |h_S| < 0.31$) and no large effect sizes. By averaging $|h|$ values within each task and then comparing between tasks, we find a mean $|\hat{h}| = 0.12 \pm 0.06$ and a median $|\hat{h}| = 0.16$, suggesting an overall strong similarity between the model and human data across all three tasks.

Discussion

In this paper, we have presented a novel model of short and intermediate-term memory and demonstrated its ability to capture human behaviour in diverse experimental paradigms. Our model allows for a realization of short-term memory capacity and recency effects by using a novel combination of vector projection and neural saturation. Additionally, the use of an auto-associative outer product network demonstrates a novel application of local learning rules to learn implicit associations in a spiking neural network. We applied this model to semantic and spatial memory tasks, achieving performance comparable to task participants and demonstrating domain generalization—a necessity for understanding how large-scale networks can realize human-like behaviour.

While this model makes strides in generalizable approaches to memory, significant work remains. Our model could be improved by further optimizing hyperparameters and neuron parameter sampling in each simulation. Notably, the relatively low dimensionality for (S)SPs used in our experiments may cause violations of the approximate orthogonality assumption in particularly unlucky cases (e.g., the bound product of two sampled, orthogonal vectors may be semi-similar to another sampled vector or the bound product of other vectors in the sampled set). We addressed this by running a large set of trials for our model; however, increased dimensionality could also play a significant role in future work to avoid issues with orthogonality. Lastly, this model notably leaves out long-term memory and its role in memory tasks. Towards this, the ITM_L allows for semi-stable representations that could be stabilized through a recall-based consolidation process that uses ITM_p to solidify long-term memories in distal synapses. Introducing consolidation would also allow for slower learning processes with access to more neurons and modality specific filtering.

In sum, our model provides a novel, general, and scalable way to implement both short-term and intermediate-term memory using dynamics realized in biologically-plausible spiking neurons.

Acknowledgments

The authors would like to thank the reviewers for their feedback as well as members of the Centre for Theoretical Neuroscience for their discussions—all of which helped to improve this paper. This work was supported by CFI (52479-10006) and OIT (35768) infrastructure funding as well as the Canada Research Chairs program, NSERC Discovery grant 261453, and AFOSR grant FA9550-17-1-0644.

References

- Atkinson, R., & Shiffrin, R. (1968). Human Memory: A Proposed System and its Control Processes. In *Psychology of Learning and Motivation* (pp. 89–195, Vol. 2). Elsevier. [https://doi.org/10.1016/S0079-7421\(08\)60422-3](https://doi.org/10.1016/S0079-7421(08)60422-3)
- Camina, E., & Güell, F. (2017). The Neuroanatomical, Neurophysiological and Psychological Basis of Memory: Current Models and Their Origins. *Frontiers in Pharmacology*, 8, 438. <https://doi.org/10.3389/fphar.2017.00438>
- Choo, F.-X. (2010, August). *The Ordinal Serial Encoding Model: Serial Memory in Spiking Neurons* [Master's thesis, University of Waterloo]. <https://uwspace.uwaterloo.ca/items/d8c13ae2-daac-4cd1-885f-041fa94aa478>
- Clopper, C. J., & Pearson, E. S. (1934). The Use of Confidence or Fiducial Limits Illustrated in the Case of the Binomial. *Biometrika*, 26(4), 404–413. <https://doi.org/10.1093/biomet/26.4.404>
- Cohen, J. (2013, May). *Statistical Power Analysis for the Behavioral Sciences* (0th ed.). Routledge. <https://doi.org/10.4324/9780203771587>
- Cowan, N. (2001). The magical number 4 in short-term memory: A reconsideration of mental storage capacity. *The Behavioral and Brain Sciences*, 24(1), 87–114, discussion 114–185. <https://doi.org/10.1017/s0140525x01003922>
- Dumont, N. S.-Y., & Eliasmith, C. (2020). Accurate representation for spatial cognition using grid cells. *Proceedings of the Annual Meeting of the Cognitive Science Society*, 42, 2367–2373. <https://escholarship.org/uc/item/8720b88v>
- Dumont, N. S.-Y., Furlong, P. M., Orchard, J., & Eliasmith, C. (2023). Exploiting semantic information in a spiking neural SLAM system. *Frontiers in Neuroscience*, 17, 1190515. <https://doi.org/10.3389/fnins.2023.1190515>
- Eliasmith, C., & Anderson, C. (2003). *Neural Engineering: Computation, Representation, and Dynamics in Neurobiological Systems*. MIT Press. <https://books.google.ca/books?id=vDhzzD6oc60C>
- Eliasmith, C. (2013, June). *How to Build a Brain: A Neural Architecture for Biological Cognition*. Oxford University Press. <https://doi.org/10.1093/acprof:oso/9780199794546.001.0001>
- Farrell, S., & Lewandowsky, S. (2004). Modelling transposition latencies: Constraints for theories of serial order memory. *Journal of Memory and Language*, 51(1), 115–135. <https://doi.org/10.1016/j.jml.2004.03.007>
- Golomb, J. D., Peelle, J. E., Addis, K. M., Kahana, M. J., & Wingfield, A. (2008). Effects of adult aging on utilization of temporal and semantic associations during free and serial recall. *Memory & Cognition*, 36(5), 947–956. <https://doi.org/10.3758/MC.36.5.947>
- Gosmann, J., & Eliasmith, C. (2021). CUE: A unified spiking neuron model of short-term and long-term memory. *Psychological Review*, 128(1), 104–124. <https://doi.org/10.1037/rev0000250>
- Gosmann, J., Voelker, A. R., & Eliasmith, C. (2017). A Spiking Independent Accumulator Model for Winner-Take-All Computation. *Proceedings of the Annual Meeting of the Cognitive Science Society*, 39, 2125–2130. <https://escholarship.org/uc/item/01z51564#main>
- Goto, A. (2022). Synaptic plasticity during systems memory consolidation. *Neuroscience Research*, 183, 1–6. <https://doi.org/10.1016/j.neures.2022.05.008>
- Henson, R. N. A., Norris, D. G., Page, M. P. A., & Baddeley, A. D. (1996). Unchained Memory: Error Patterns Rule out Chaining Models of Immediate Serial Recall. *The Quarterly Journal of Experimental Psychology A*, 49(1), 80–115. <https://doi.org/10.1080/027249896392810>
- Howard, M. W., & Kahana, M. J. (1999). Contextual variability and serial position effects in free recall. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 25(4), 923–941. <https://doi.org/10.1037/0278-7393.25.4.923>
- Jahnke, J. C. (1968). Delayed recall and the serial-position effect of short-term memory. *Journal of Experimental Psychology*, 76(4, Pt.1), 618–622. <https://doi.org/10.1037/h0025692>
- Kamiński, J. (2017). Intermediate-Term Memory as a Bridge between Working and Long-Term Memory. *The Journal of Neuroscience*, 37(20), 5045–5047. <https://doi.org/10.1523/JNEUROSCI.0604-17.2017>
- Kamiński, J., Sullivan, S., Chung, J. M., Ross, I. B., Mamelak, A. N., & Rutishauser, U. (2017). Persistently active neurons in human medial frontal and medial temporal lobe support working memory. *Nature Neuroscience*, 20(4), 590–601. <https://doi.org/10.1038/nn.4509>
- Kandel, E. R. (2001). The Molecular Biology of Memory Storage: A Dialogue Between Genes and Synapses. *Science*, 294(5544), 1030–1038. <https://doi.org/10.1126/science.1067020>
- Komer, B., Stewart, T. C., Voelker, A. R., & Eliasmith, C. (2019). A neural representation of continuous space using fractional binding. *Proceedings of the Annual Meeting of the Cognitive Science Society*, 41, 2038–2043. <https://escholarship.org/uc/item/3zz346g1>
- Maidenbaum, S., Kremen, V., Sladky, V., Miller, K., Gompel, J. V., Worrell, G. A., & Jacobs, J. (2025). Improved spatial memory for physical versus virtual navigation. *Journal of Neural Engineering*, 22(4), 046014. <https://doi.org/10.1088/1741-2552/ade6aa>

- Marshall, P. H., & Werder, P. R. (1972). The effects of the elimination of rehearsal on primacy and recency. *Journal of Verbal Learning and Verbal Behavior*, *11*(5), 649–653. [https://doi.org/https://doi.org/10.1016/S0022-5371\(72\)80049-5](https://doi.org/https://doi.org/10.1016/S0022-5371(72)80049-5)
- Mongillo, G., Barak, O., & Tsodyks, M. (2008). Synaptic Theory of Working Memory. *Science*, *319*(5869), 1543–1546. <https://doi.org/10.1126/science.1150769>
- Mullally, S. L., & Maguire, E. A. (2014). Memory, Imagination, and Predicting the Future: A Common Brain Mechanism? *The Neuroscientist*, *20*(3), 220–234. <https://doi.org/10.1177/1073858413495091>
- Plate, T. (1995). Holographic reduced representations. *IEEE Transactions on Neural Networks*, *6*(3), 623–641. <https://doi.org/10.1109/72.377968>
- Polyn, S. M., Norman, K. A., & Kahana, M. J. (2009). A context maintenance and retrieval model of organizational processes in free recall. *Psychological Review*, *116*(1), 129–156. <https://doi.org/10.1037/a0014420>
- Reeves, A., & Sperling, G. (1986). Attention gating in short-term visual memory. *Psychological Review*, *93*(2), 180–206. <https://doi.org/10.1037/0033-295X.93.2.180>
- Reimann, S. (2025). Memory States From Almost Nothing: Representing and Computing in a Nonassociative Algebra. *Neural Computation*, *37*(6), 1154–1170. https://doi.org/10.1162/neco_a_01755
- Talmi, D., Grady, C. L., Goshen-Gottstein, Y., & Moscovitch, M. (2005). Neuroimaging the Serial Position Curve: A Test of Single-Store Versus Dual-Store Models. *Psychological Science*, *16*(9), 716–723. <https://doi.org/10.1111/j.1467-9280.2005.01601.x>
- Whittington, J. C. R., Muller, T. H., Mark, S., Chen, G., Barry, C., Burgess, N., & Behrens, T. E. J. (2020). The Tolman-Eichenbaum Machine: Unifying Space and Relational Memory through Generalization in the Hippocampal Formation. *Cell*, *183*(5), 1249–1263.e23. <https://doi.org/https://doi.org/10.1016/j.cell.2020.10.024>