

Compositional Memory in Attractor Neural Networks with One-Step Learning

Gregory P. Davis^{a,*}, Garrett E. Katz^b, Rodolphe J. Gentili^c, James A. Reggia^a

^a*Department of Computer Science, University of Maryland, College Park, MD, USA*

^b*Department of Elec. Engr. and Comp. Sci., Syracuse University, Syracuse, NY, USA*

^c*Department of Kinesiology, University of Maryland, College Park, MD, USA*

Abstract

Compositionality refers to the ability of an intelligent system to construct models out of reusable parts. This is critical for the productivity and generalization of human reasoning, and is considered a necessary ingredient for human-level artificial intelligence. While traditional symbolic methods have proven effective for modeling compositionality, artificial neural networks struggle to learn systematic rules for encoding generalizable structured models. We suggest that this is due in part to short-term memory that is based on persistent maintenance of activity patterns without fast weight changes. We present a recurrent neural network that encodes structured representations as systems of contextually-gated dynamical attractors called attractor graphs. This network implements a functionally compositional working memory that is manipulated using top-down gating and fast local learning. We evaluate this approach with empirical experiments on storage and retrieval of graph-based data structures, as well as an automated hierarchical planning task. Our results demonstrate that compositional structures can be stored in and retrieved from neural working memory without persistent maintenance of multiple activity patterns. Further, memory capacity is improved by the use of a fast store-erase learning rule that permits controlled erasure and mutation of previously learned associations. We conclude that the combination of top-down gating and fast associative learning provides recurrent neural networks with a robust functional mechanism for compositional working memory.

Keywords: Compositionality, Working memory, Itinerant attractor dynamics, Multiplicative gating, One-step learning, Programmable neural networks

1

2 1. Introduction

3 Compositionality refers to the ability of an intelligent system to construct representations out of reusable parts.
4 The principle of compositionality states that “the meaning of a complex expression is determined by its structure
5 and the meanings of its constituents” (Fodor and Pylyshyn, 1988; Nefdt, 2020; Szabó, 2012). Despite widespread
6 debate over the precise definition and interpretation of compositionality, there is substantial evidence that structured
7 representation plays a critical role in human reasoning. Compositional representations are useful because they can be
8 systematically generalized and reorganized to facilitate rapid comprehension in novel circumstances. For example, a
9 person who knows how to prepare tea can learn to prepare coffee with minimal difficulty by decomposing the process
10 and recognizing familiar behaviors (e.g., boiling water). Compositional reasoning is considered crucial in domains as
11 diverse as language comprehension (Baayen, 1994; Pelletier, 1994; Aizawa, 2003), behavioral planning and imitation
12 (Baayen, 1994; Reverberi et al., 2012; Botvinick, 2008), visual perception (Riesenhuber and Poggio, 1999; Witkin
13 and Tenenbaum, 1983; Bienenstock et al., 1997), and concept learning (Kamp and Partee, 1995; Barsalou, 1993;
14 Piantadosi et al., 2016).

*Corresponding author.

Email addresses: grpdavis@umd.edu (Gregory P. Davis), gkatz01@syr.edu (Garrett E. Katz), rodolphe@umd.edu (Rodolphe J. Gentili), reggia@umd.edu (James A. Reggia)

15 Compositionality is readily achieved in cognitive systems capable of symbolic manipulation, but is much more
16 challenging for sub-symbolic systems such as artificial neural networks. This has fueled a long-standing controversy
17 over whether neural networks can represent compositional structures (Fodor and Pylyshyn, 1988; Marcus, 2018; Lake
18 et al., 2017). Recent advances in neural machine translation and sequence-to-sequence modeling have demonstrated
19 remarkable progress toward compositional learning in neural systems. This is due to several innovations that improve
20 short term memory in neural networks, including recurrent processing units (Hochreiter and Schmidhuber, 1997;
21 Chung et al., 2014), attention mechanisms (Vaswani et al., 2017; Hupkes et al., 2018), and external memory resources
22 (Graves et al., 2014; Sukhbaatar et al., 2015; Pham et al., 2018). These techniques are often combined because
23 they provide complementary benefits, allowing neural networks to maintain activation states over time and model
24 dependencies between distant representations (Bahdanau et al., 2014; Lake, 2019).

25 Despite this progress, and although there are disagreements about how compositionality should be evaluated in
26 neural models (Hupkes et al., 2020; Nefdt, 2020), empirical studies demonstrate that state-of-the-art deep neural
27 networks struggle to learn systematic rules that permit generalization beyond training data (Lake and Baroni, 2018;
28 Loula et al., 2018; Hupkes et al., 2020). This is in stark contrast to the ease with which such rules can be implemented
29 in symbolic programs. In our view, this discrepancy is due to a lack of *compositional working memory* in neural
30 networks and an inability to encode structured representations. Working memory is a form of short term memory
31 containing information that is actively manipulated by cognitive processes (Baars, 2007; Oberauer, 2009). Although
32 its capacity is limited, working memory is capable of storing structured or “chunked” representations that provide
33 access to a broad range of information (Cowan, 2010; Campitelli et al., 2007).

34 One way to provide neural networks with compositional working memory is to integrate them into hybrid sys-
35 tems with non-neural symbolic algorithms that manipulate compositional data structures (Besold et al., 2015; Marcus,
36 2020; Kipf et al., 2019; Burke et al., 2019; Andreas et al., 2016). For example, some neural-guided search algorithms
37 maintain compositional data structures in non-neural symbolic memory and use neural processing to inform their
38 construction (Bunel et al., 2018; Silver et al., 2016; Kalyan et al., 2018). From an engineering standpoint, this is a
39 reasonable and effective approach, but it does not address how compositional structures can be encoded in purely neu-
40 ral models. An alternative approach is to develop purely neural computers with general-purpose memory arrays that
41 incorporate key aspects of symbolic computation (Graves et al., 2016, 2014). While these systems have an impressive
42 ability to learn algorithmic procedures from data, they require unconstrained access to large sets of activation patterns
43 maintained in external memory, which is considered highly implausible from a biological perspective.

44 Many purely neural models of working memory employ spatially localized representations, such as localized
45 attractors of Dynamic Field Theory (Sandamirskaya et al., 2013; Erlhagen and Schöner, 2002), and minimally over-
46 lapping cell assemblies in neural blackboard architectures (Van der Velde and Kamps, 2006). Localist representations
47 introduce an undesirable correspondence between memory and architecture that limits representational flexibility, and
48 often requires task-specific circuitry. In contrast, vector-space approaches employ fully distributed representations that
49 can be composed using superposition and binding operations (Gayler, 2003; Plate, 1995). For example, the Sema-
50 ntic Pointer Architecture uses symbol-like vector representations that can be recursively composed to store structured
51 information (Blouw et al., 2016; Eliasmith et al., 2012). Notably, operations for composing semantic pointers can
52 be learned using biologically-plausible learning rules in spiking neural networks (Stewart et al., 2011). However, a
53 significant disadvantage to this approach is that the encodings of structured representations are semantically related to
54 their constituent elements. This means that additions to a compositional data structure involve changes to the sema-
55 ntic pointer that encodes it, as well as any semantic pointers for super-structures that it is contained in. For example,
56 adding a leaf node to a tree would require reconstruction of the semantic pointers for each ancestor of the new node.
57 This makes semantic pointer encodings of data structures effectively immutable, as modifications involve constructing
58 new encodings.

59 Performance on working memory tasks is correlated with measures of general intelligence (Conway et al., 2003;
60 Colom et al., 2004; Jaeggi et al., 2008), and working memory operations are considered to be consciously reportable
61 (Baars and Franklin, 2003; Baddeley, 1993). For these reasons, biologically-plausible models of compositional work-
62 ing memory may lead to significant advances in AI systems and contribute to a deeper understanding of consciousness
63 and cognition (Reggia et al., 2020, 2019), including cognitive-motor control (Hauge et al., 2020, 2019). Substantial
64 evidence indicates that activity-silent mechanisms such as rapid synaptic plasticity play a critical role in working
65 memory (Manohar et al., 2019; Mongillo et al., 2008; Rose et al., 2016; Stokes, 2015; Barbosa et al., 2019). This
66 suggests that persistent activity maintenance alone is not sufficient for modeling the complexities of human working

67 memory, and may explain in part why compositional learning is difficult for artificial neural networks.

68 Contemporary neural networks typically undergo a training phase during which weights are updated via iterative
69 gradient descent and fixed during task performance. Recent exceptions to this demonstrate that fast associative learn-
70 ing greatly improves short-term memory in neural networks because it permits storage without active maintenance
71 (Ba et al., 2016; Danihelka et al., 2016; Miconi et al., 2018). Furthermore, models based entirely on fast associative
72 learning and itinerant dynamical attractors can learn to perform complex working memory tasks (Sylvester and Reg-
73 gia, 2016; Sylvester et al., 2013), and can simulate Turing machines without the need for consistent maintenance of
74 memory activation (Katz et al., 2019). However, such models have focused on storage of individual memories and
75 temporal sequences of memories. To our knowledge, fast associative learning has not yet been applied to explicit
76 encoding of hierarchical compositional structures (e.g., trees) in working memory as we do here.

77 In this paper, we introduce a neural model of compositional memory based on context-dependent itinerant attrac-
78 tors in recurrent neural networks. We refer to such models as *attractor graph networks* or AGNs. Attractor graphs
79 are composed of fixed-point dynamical attractors (vertices) and transitions between attractors (edges) that are learned
80 using a combination of multiplicative gating and one-step associative learning. The use of multiplicative gating is
81 motivated by evidence that it contributes to functional grouping in neural populations (Baan et al., 2019; Vecoven
82 et al., 2020; Masse et al., 2018; Chabuk and Reggia, 2013; Rikhye et al., 2018). In an AGN, contextual gating signals
83 select learned associations to govern the dynamics of the network over time, directing a traversal through the attractor
84 graph that is analogous to the operation of a finite state machine. We show how compositional data structures such
85 as associative arrays (i.e., dictionaries or maps), linked lists, and trees can be encoded in AGNs and retrieved by pro-
86 grammatic procedures that control sequential iteration through attractor graphs. In contrast to semantic pointers and
87 other vector space representations, structured representations in AGNs are made up of learned associations between
88 activity states, and can therefore be freely modified by changing synaptic weights without changing activity states.

89 AGNs are inspired by neurobiological studies of working memory and are based on several biologically-inspired
90 principles. Representations in attractor graphs are composed of distributed activation patterns (dynamical attractors)
91 that are supported by learned connection weights, and can be reactivated as needed without persistent activity mainte-
92 nance. This means that the capacity of working memory is not directly limited by network architecture, and is instead
93 determined by the organization of both the underlying attractor model and the control signals used for memory re-
94 trieval. We propose that the phenomenon of “chunking” is supported by compositional structuring, which organizes
95 the contents of working memory according to systematic procedures for top-down control.

96 2. Methods

97 In this section we describe our method for representing compositional data structures as systems of dynamical
98 attractors in recurrent neural networks. We illustrate this with a small multi-region model shown in Figure 1. The
99 core memory region (*mem*) is a recurrent neural network with attractor states that represent general-purpose elements
100 in short-term working memory. We refer to these elements as *memory states*. Memory states are linked together with
101 context-dependent transitions that are controlled by inputs from the context region (*ctx*). Patterns of activity in the
102 lexicon region (*lex*) represent symbolic tokens that can be “stored” in memory states and used as names for variable
103 pointers to memory states.

104 Compositional data structures are encoded in the *mem* region as systems of itinerant dynamical attractors called
105 *attractor graphs*. Vertices in these graphs are fixed-point dynamical attractors, and edges are context-dependent
106 transitions between attractor states. We refer to recurrent networks with attractor graph dynamics as *attractor graph*
107 *networks*, or AGNs. Section 2.1 describes the structure and dynamics of AGNs in detail. Section 2.2 shows how
108 compositional data structures can be represented as attractor graphs, and describes the interactions between the *mem*,
109 *ctx*, and *lex* regions of the model. Section 2.3 shows how the model shown in Figure 1 can be embedded into a larger
110 programmable neural network that constructs and manipulates data structures according to learned programmatic
111 procedures.

112 2.1. Attractor Graph Networks

113 Attractor graphs are systems of itinerant dynamical attractors with context-dependent transitions. Attractor itiner-
114 ancy refers to a temporal process in which a dynamical system undergoes transitions through a sequence of attractor

144 cumulative synaptic activity. These signals enable functional branching in attractor transitions by modulating hetero-
 145 associative dynamics, as described below, and are distinct from connection gates (e.g., g_r^A).

146 The following equations describe the dynamics of an AGN region such as *mem* over time. In Section 2.2, we
 147 return to the neural circuit shown in Figure 1 to explain how these dynamics contribute to compositional working
 148 memory. Here we use subscript r to refer to a generic region. First, synaptic input is aggregated from gated recurrent
 149 connectivity and external inputs:

$$\mathbf{s}_r(t) = \underbrace{g_r^S(t) \omega_r \mathbf{v}_r(t)}_{\text{saturation}} + \underbrace{g_r^A(t) A_r \mathbf{v}_r(t)}_{\text{convergence}} + \underbrace{g_r^H(t) H_r \mathbf{v}_r(t)}_{\text{transition}} + \mathbf{I}_r(t) \quad (1)$$

150 where

- 151 • $\mathbf{s}_r(t)$ is a vector of cumulative synaptic input to region r at time t .
- 152 • $\mathbf{v}_r(t)$ is a vector of neural activity of region r at time t .
- 153 • ω_r is a scalar self-weight that causes saturation and maintenance of neural activity in region r when $g_r^S(t) = 1$.
- 154 • A_r is an auto-associative weight matrix for region r that causes convergence to a nearby fixed-point attractor
 155 when $g_r^A(t) = 1$.
- 156 • H_r is a hetero-associative weight matrix for region r that causes a transition between attractor states when
 157 $g_r^H(t) = 1$.
- 158 • $\mathbf{I}_r(t)$ is a vector of external (non-recurrent) synaptic input to region r at time t . This input may be provided by
 159 other neural regions, as described in Section 2.2, or from outside the model for initialization purposes, as shown
 160 in Section 2.1.1.

161 Next, multiplicative activation is determined based on gated contextual input:

$$\mathbf{x}_r(t) = \begin{cases} \mathbf{c}_r(t), & \text{if } g_r^C(t) = 1 \\ \mathbf{1}, & \text{otherwise} \end{cases} \quad (2)$$

162 where

- 163 • $\mathbf{c}_r(t)$ is a vector containing the net multiplicative input to region r at time t . This input may be provided by other
 164 neural regions such as *ctx* in Figure 1.
- 165 • $\mathbf{1}$ is a vector of ones of the same size as $\mathbf{c}_r(t)$.
- 166 • $\mathbf{x}_r(t)$ is a vector of active multiplicative input to region r at time t . When $g_r^C(t)$ is enabled, $\mathbf{x}_r(t) = \mathbf{c}_r(t)$.
 167 Otherwise, $\mathbf{x}_r(t)$ defaults to $\mathbf{1}$, and synaptic input $\mathbf{s}_r(t)$ passes into the activation function regardless of $\mathbf{c}_r(t)$.

168 Finally, neural activation is computed for the next timestep by combining synaptic and multiplicative input:

$$\mathbf{v}_r(t+1) = \sigma_r(\mathbf{x}_r(t) \odot \mathbf{s}_r(t)) \quad (3)$$

169 where

- 170 • σ_r is a sign-preserving neural activation function in region r (i.e., $\text{sgn}(\sigma_r(x)) = \text{sgn}(x)$). For the experiments
 171 reported in Section 3, we use the hyperbolic tangent activation function for the *mem* AGN shown in Figure 1.
- 172 • \odot is the Hadamard (element-wise) product.

173 An attractor transition is carried out in four stages, starting with an activity pattern that may be initialized by $\mathbf{I}_r(t)$,
174 or by a previous transition. 1) The transition begins with application of a context pattern \mathbf{c} , which disables a subset of
175 neurons and “masks” the source activation pattern. This is done by enabling saturation and context gates ($g_r^C(t)$ and
176 $g_r^S(t)$) for one timestep. 2) Next, while context gate $g_r^C(t)$ remains enabled, hetero-associative gate $g_r^H(t)$ is enabled
177 for one timestep, and activity transitions to a new pattern of non-zero activation in the participating neurons. 3)
178 Once the initial transition is complete, $g_r^A(t)$ is enabled, causing auto-associative dynamics across the entire network.
179 Over several timesteps, activity converges to the nearest fixed-point attractor. Note that because context gate $g_r^C(t)$
180 is disabled, all neurons participate in auto-associative dynamics. 4) Finally, saturation gate $g_r^S(t)$ is enabled, causing
181 activity to saturate within the current orthant of activity space over several timesteps. This makes up for any errors
182 in convergence, and is successful as long as auto-associative dynamics resulted in an activity pattern in the correct
183 orthant. Note that saturation dynamics are only necessary for continuous models (e.g., when σ_r is the hyperbolic
184 tangent), and can be omitted in models with a threshold activation function (e.g., the sign/signum function). The
185 connection gate values for each stage of an attractor transition are shown in Table 1. Stages 1 and 2 take one timestep
186 each, while Stages 3 and 4 each occur over several timesteps.

Table 1: Connection gate values for each stage of an AGN attractor transition

Stage	g_r^S	g_r^C	g_r^H	g_r^A
1. Mask	1	1		
2. Transition		1	1	
3. Converge				1
4. Saturate	1			

187 Each attractor state learned by the AGN may have multiple hetero-associative transitions to several other attractor
188 states. Each of the transitions from a given state must be learned in the context of a unique multiplicative input pattern
189 $\mathbf{c}_r(t)$. During attractor transition dynamics, the choice of $\mathbf{c}_r(t)$ determines which learned association will govern the
190 transition. This process is depicted in Figure 2, which shows the activation space for an AGN with three neurons.
191 Trajectories are shown for two transitions from an attractor state $\mathbf{m}^{(0)}$ using unique context patterns $\mathbf{c}^{(1)}$ and $\mathbf{c}^{(2)}$.

192 Learned attractor states are patterns in $\{-\rho_r, +\rho_r\}^{N_r}$, where N_r is the number of neurons in the AGN, and ρ_r is a
193 parameter that determines the magnitude of learned activity states. These activation patterns are randomly generated
194 using a Bernoulli process with equal probabilities (i.e., a fair coin toss). For models using a threshold activation
195 function, activation patterns are discrete and bipolar, and $\rho_r = 1$. When the hyperbolic tangent activation function is
196 used, $0 << \rho_r < 1$ (typically 0.9999) because $\sigma_r^{-1}(1) = \infty$. Once ρ_r is set, the value of the saturation self-weight ω_r
197 is determined according to the following relation:

$$\rho_r = \sigma_r(\omega_r \rho_r)$$

198 This ensures that $\pm\rho_r$ is a stable fixed-point when saturation dynamics are enabled ($g_r^S(t) = 1$).

199 Multiplicative gating patterns $\mathbf{c}_r(t)$ that contextualize transitions are in $\{0, 1\}^N$, where N is the number of neurons
200 in the network. The density of these patterns is an important parameter that determines how many neurons participate
201 in each transition. This parameter is referred to as λ , the probability used in the Bernoulli process that generates
202 context patterns.

203 The auto-associative weight matrix for a region r (A_r) is learned using traditional Hebbian learning, starting with
204 an initial matrix with zero entries, and updating the weights for each learned pattern:

$$\Delta A_r = \frac{1}{\rho_r^2 N_r} \sigma_r^{-1}(\mathbf{v}) \mathbf{v}^\top \quad (4)$$

205 where \mathbf{v} is the pattern to be learned, σ_r is the neural activation function, ρ_r is the stable activation level of learned
206 attractor patterns, and N_r is the number of neurons in the region. When the sign/signum activation function is used
207 ($\sigma_r = \text{sgn}(x)$), its inverse is defined as:

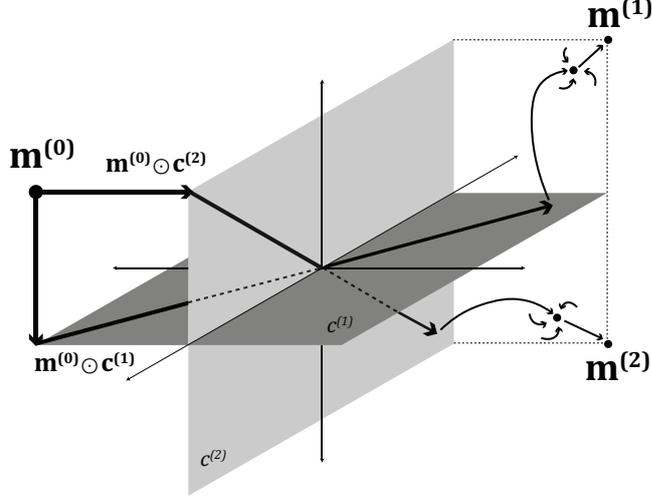


Figure 2: Visual depiction of context-dependent attractor transitions in the activation space of an AGN with three neurons. Each axis represents the activation level of one neuron (from -1 to +1). Two transitions from state $\mathbf{m}^{(0)}$ (left side) are shown, each with a distinct context pattern ($\mathbf{c}^{(1)}$ or $\mathbf{c}^{(2)}$). Each context pattern is a binary vector that selects subsets of neurons to participate in hetero-associative dynamics, and corresponds to a subspace in activation space (shaded planes labeled $c^{(1)}$ and $c^{(2)}$). Transitions are executed over several steps. 1) First, a context pattern is used to “mask” the source activation pattern, disabling a subset of neurons and collapsing activity into the context-specific subspace (arrows leaving $\mathbf{m}^{(0)}$ to $\mathbf{m}^{(0)} \odot \mathbf{c}^{(1)}$ or $\mathbf{m}^{(0)} \odot \mathbf{c}^{(2)}$, left side). 2) Next, while the context masking remains active, hetero-associative dynamics cause activation to transition into a new orthant of the subspace (center arrows within shaded planes). 3) Then, the context masking is disabled, and auto-associative dynamics cause convergence to the nearest fixed-point attractor state over several timesteps (curved arrows, right side). 4) To correct any errors in convergence, saturation dynamics push activity to the corner of the current orthant of activation space over several timesteps (straight arrows to $\mathbf{m}^{(1)}$ and $\mathbf{m}^{(2)}$, right side).

$$\text{sgn}^{-1}(x) = \begin{cases} +1, & \text{if } x > 0 \\ -1, & \text{if } x < 0 \\ 0, & \text{if } x = 0 \end{cases}$$

208 The hetero-associative weight matrix for a region r (H_r) is also learned using local one-step learning. However,
 209 because transitions are context-dependent, they must be learned with the corresponding contextual gating pattern.
 210 Each transition is learned with the following update rule, starting with a zero matrix:

$$\Delta H_r = \frac{1}{\lambda \rho_r^2 N_r} \sigma_r^{-1}(\mathbf{c} \odot \mathbf{v})(\mathbf{c} \odot \mathbf{u})^\top \quad (5)$$

211 where \mathbf{u} is the source pattern, \mathbf{v} is the target pattern, \mathbf{c} is the context pattern, λ is the density of contextual gating
 212 patterns, and σ_r , ρ_r , and N_r are as defined above. Note that the normalization factor $\lambda \rho_r^2 N$ is approximately equivalent
 213 to the squared Euclidean length of $\mathbf{c} \odot \mathbf{u}$.

214 Other associative learning rules can be used to learn A_r and H_r . To adapt a learning rule for gated hetero-associative
 215 learning, it must 1) only update weights connecting neurons that participate in the transition, as determined by con-
 216 textual gating pattern \mathbf{c} , and 2) renormalize weight updates according to λ , the density of \mathbf{c} . In this paper, we consider
 217 both traditional Hebbian learning (Equations 4 and 5) and the fast store-erase learning rule (Katz et al., 2019):

$$\Delta A_r = \frac{1}{\rho_r^2 N_r} \left(\sigma_r^{-1}(\mathbf{v}) - (A_r \mathbf{v}) \right) \mathbf{v}^\top \quad (6)$$

$$\Delta H_r = \frac{1}{\lambda \rho_r^2 N_r} \left(\underbrace{\sigma_r^{-1}(\mathbf{c} \odot \mathbf{v}) - (\mathbf{c} \odot H_r(\mathbf{c} \odot \mathbf{u}))}_{\text{masked target delta}} \right) \underbrace{(\mathbf{c} \odot \mathbf{u}^\top)}_{\text{source}} \quad (7)$$

218 where \mathbf{c} is a binary contextual gating pattern, \mathbf{u} is the initial source activity pattern, \mathbf{v} is the final target activity pattern,
 219 and λ , σ_r , ρ_r , and N_r are as defined above. The additional context masking in the target delta of Equation 7 (just
 220 before H_r) ensures that weights with deactivated post-synaptic neurons are not updated.

221 The store-erase learning rule contains an anti-Hebbian component that erases previously stored associations. This
 222 makes it particularly advantageous for neural working memory, as it can be used to overwrite relations between
 223 memories and modify learned data structures. However, this learning rule has not been evaluated with auto-associative
 224 memory, and has not been systematically compared with traditional Hebbian learning. This is addressed in Section 3,
 225 where we present empirical results that compare memory storage and retrieval in AGNs learned with either traditional
 226 Hebbian or store-erase learning.

227 2.1.1. Contextual Gating Supports Functional Branching

228 Contextual gating signals in AGNs make it possible to store multiple transitions from an attractor state using a
 229 single hetero-associative weight matrix. We refer to this branching of attractor transitions as *functional branching*
 230 because it depends on functional multiplicative inputs. This is in contrast to *structural branching*, in which each
 231 branch in the attractor sequence is stored in a distinct hetero-associative weight matrix. Structural branching imposes
 232 an undesirable correspondence between network architecture and memory because the number of outgoing transitions
 233 from any given attractor state (i.e., its out-degree or branching factor) is limited by the number of hetero-associative
 234 weight matrices. In contrast, functional branching requires only a single hetero-associative weight matrix, and does
 235 not impose constraints on the structure of learned associations. This novel aspect of AGNs makes it possible to
 236 represent arbitrary directed graphs as systems of functionally-branched itinerant attractor sequences.

237 To clarify how AGNs work, next we illustrate how contextual gating signals support functional branching with a
 238 simple toy example. For this example, we consider an AGN with $N = 4$ neurons that use the sign/signum activation
 239 function ($\rho = 1$ and $\omega = 1$). Learned attractor states are bipolar patterns in $\{-1, +1\}^4$. The network, shown in
 240 Figure 3a, learns two attractor transitions (Figure 3b) that make up a simple attractor graph (Figure 3c). The learning
 241 procedure and transition dynamics are described in detail below, with the subscript r omitted for ease of presentation.

242 Three activity patterns are learned as attractors in auto-associative matrix A_r using Equation 4: $\mathbf{v}^{(0)} = [+1, -1, +1, -1]^\top$,
 243 $\mathbf{v}^{(1)} = [+1, +1, -1, -1]^\top$, and $\mathbf{v}^{(2)} = [-1, +1, -1, +1]^\top$. For simplicity, we omit the normalization term $\frac{1}{\rho^2 N}$ from the
 244 learning rule.

$$A = \sum_i \mathbf{v}^{(i)} \mathbf{v}^{(i)\top} = \mathbf{v}^{(0)} \mathbf{v}^{(0)\top} + \mathbf{v}^{(1)} \mathbf{v}^{(1)\top} + \mathbf{v}^{(2)} \mathbf{v}^{(2)\top} = \begin{bmatrix} +3 & -1 & +1 & -3 \\ -1 & +3 & -3 & +1 \\ +1 & -3 & +3 & -1 \\ -3 & +1 & -1 & +3 \end{bmatrix}$$

245 Note that these patterns were chosen for illustration purposes, and that AGNs are not restricted to learning orthogo-
 246 nal/complementary activation patterns.

247 Two attractor transitions are learned in hetero-associative matrix H using Equation 5: $\mathbf{v}^{(0)}$ to $\mathbf{v}^{(1)}$, and $\mathbf{v}^{(0)}$ to $\mathbf{v}^{(2)}$.
 248 Each of these transitions is contextualized by a distinct context pattern that selects a subset of neurons in the network.
 249 The first transition is contextualized by $\mathbf{c}^{(1)} = [1, 1, 0, 0]^\top$, which enables the first and second neurons, while the
 250 second transition is contextualized by $\mathbf{c}^{(2)} = [0, 0, 1, 1]^\top$, which enables the third and fourth neurons. As above, we
 251 omit the renormalization term $\frac{1}{\lambda \rho^2 N}$ for simplicity.

$$H = (\mathbf{c}^{(1)} \odot \mathbf{v}^{(1)})(\mathbf{c}^{(1)} \odot \mathbf{v}^{(0)})^\top + (\mathbf{c}^{(2)} \odot \mathbf{v}^{(2)})(\mathbf{c}^{(2)} \odot \mathbf{v}^{(0)})^\top = \begin{bmatrix} +1 & +1 & 0 & 0 \\ -1 & -1 & 0 & 0 \\ 0 & 0 & -1 & +1 \\ 0 & 0 & +1 & -1 \end{bmatrix}$$

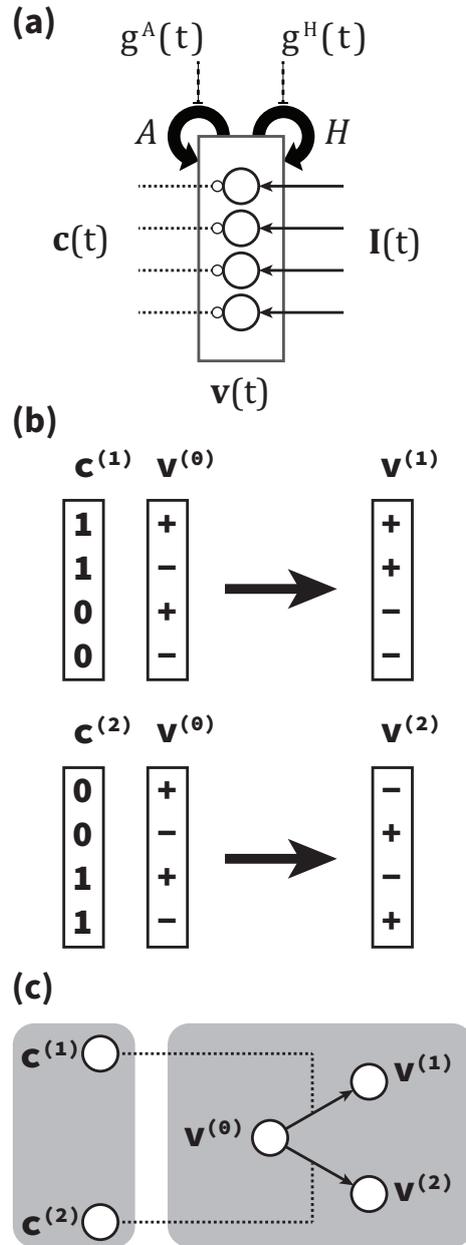


Figure 3: Toy example of transitions in an AGN (see text for details). (a) AGN with four neurons (center circles). Bold looped arrows indicate auto-associative (A) and hetero-associative (H) recurrent connectivity matrices. Dashed lines on the left indicate multiplicative gating signals ($\mathbf{c}(t)$) that contextualize attractor transitions. Solid lines on the right indicate external synaptic input ($\mathbf{I}(t)$) that initializes activation patterns ($\mathbf{v}(t)$). Binary gates ($g^A(t)$ and $g^H(t)$, top) determine which connections are active at each timestep ($g^S(t)$ and $g^C(t)$ not shown). (b) Two contextual transitions learned in the network. Each transition begins with activity pattern $\mathbf{v}^{(0)}$ and transitions to a unique successor ($\mathbf{v}^{(1)}$ or $\mathbf{v}^{(2)}$) in the context of a unique context pattern ($\mathbf{c}^{(1)}$ or $\mathbf{c}^{(2)}$). (c) Graphical depiction of the learned attractor graph. Gating patterns (circles in left rectangle) contextualize transitions between distributed activity states in the AGN (circles in right rectangle).

252 The resulting weight matrix H encodes two transitions from $\mathbf{v}^{(0)}$ that are segregated to different sub-populations of
 253 neurons and weights. The disjoint context patterns used in this example are solely for illustrative purposes; in general,
 254 context patterns can be randomly generated and can enable overlapping sets of neurons, introducing interference

255 between weight updates for different context patterns. Empirical results presented in Section 3 demonstrate successful
 256 transitions despite this interference.

257 When a transition is executed, the corresponding context pattern ($\mathbf{c}^{(1)}$ or $\mathbf{c}^{(2)}$) is presented to select the subset of
 258 neurons that learned the transition. Consequently, only one quadrant of H governs the first step of transition dynamics,
 259 updating activity in the context-specific subset of neurons associated with the transition. To complete the transition,
 260 auto-associative dynamics cause convergence to the nearest attractor, completing the target activity pattern. This
 261 process is illustrated in Table 2 for the transition from $\mathbf{v}^{(0)}$ to $\mathbf{v}^{(1)}$ using context signal $\mathbf{c}^{(1)}$ (blank cells indicate zero
 262 values).

Table 2: Timing of example attractor transition dynamics

Timestep	$g^S(t)$	$g^C(t)$	$g^H(t)$	$g^A(t)$	$\mathbf{I}(t)$	$\mathbf{c}(t)$	$\mathbf{x}(t)$	$\mathbf{s}(t)$	$\mathbf{v}(t)$
Initialization (t=0)					$\mathbf{v}^{(0)}$		$\mathbf{1}$	$\mathbf{v}^{(0)}$	$\sigma_r(\mathbf{v}^{(0)}) = \begin{bmatrix} +1 \\ -1 \\ +1 \\ -1 \end{bmatrix}$
Masking (t=1)	1	1				$\mathbf{c}^{(1)}$	$\mathbf{c}^{(1)}$	$\mathbf{v}^{(0)}$	$\sigma_r(\mathbf{c}^{(1)} \odot \mathbf{v}^{(0)}) = \begin{bmatrix} +1 \\ -1 \\ 0 \\ 0 \end{bmatrix}$
Transition (t=2)		1	1			$\mathbf{c}^{(1)}$	$\mathbf{c}^{(1)}$	$H\mathbf{v}^{(1)}$	$\sigma_r(\mathbf{c}^{(1)} \odot H\mathbf{v}^{(1)}) = \begin{bmatrix} +1 \\ +1 \\ 0 \\ 0 \end{bmatrix}$
Convergence (t=3)				1			$\mathbf{1}$	$A\mathbf{v}^{(2)}$	$\sigma_r(A\mathbf{v}^{(2)}) = \begin{bmatrix} +1 \\ +1 \\ -1 \\ -1 \end{bmatrix}$

263 This procedure is carried out in an AGN using a temporal sequence of inputs and connection gate values shown
 264 at the top of each timestep block (e.g., $\mathbf{I}(0) = \mathbf{v}^{(0)}$). The second transition from $\mathbf{v}^{(0)}$ to $\mathbf{v}^{(2)}$ can be executed by
 265 substituting context pattern $\mathbf{c}^{(2)}$ for $\mathbf{c}^{(1)}$. Note that saturation dynamics are omitted after convergence because the sign
 266 activation function has discrete outputs. The saturation connection is, however, necessary for activity maintenance
 267 during timestep 1. Although this simple example only requires one timestep of auto-associative dynamics to recover
 268 the final pattern, in general, several timesteps will be necessary.

269 2.2. Compositional Memory

270 The AGN model outlined in the previous section describes the dynamics of the *mem* region of the model shown in
 271 Figure 1. In this section, we describe the interactions between the *mem*, *ctx*, and *lex* regions of this model, and explain
 272 how they encode compositional data structures.

273 Attractor states in the *mem* region represent discrete generic memory states, and are organized into structured
 274 representations with context-dependent attractor transitions. For example, a linked list is represented by a sequence of
 275 memory states connected by transitions with a shared list-specific context. Each memory state may be associated with
 276 a pattern of activity in the lexicon region (*lex*) via the pathway from *mem* to *lex*. Patterns in *lex* represent symbolic
 277 tokens that are “stored” in memory states using associative learning. These tokens may represent words, for example,
 278 and a sentence could be represented as a linked list of memory states, each associated with the *lex* pattern representing
 279 the corresponding word in the sentence.

280 Transitions between memory states are contextualized by activity patterns in the context region (*ctx*). The *ctx*
 281 region uses the heaviside activation function, and provides the binary context signals that select subsets of *mem*
 282 neurons to participate in attractor transitions ($\mathbf{c}(t)$ in Equation 2). Patterns in *ctx* may be derived from *mem* and *lex*
 283 patterns via pathways from those regions to *ctx*. This makes it possible for memory states or symbolic tokens to
 284 indirectly contextualize memory transitions through an intermediate associated *ctx* activity pattern.

285 While the *mem* region is an AGN with multiple recurrent pathways and attractor dynamics, the *ctx* and *lex* regions
 286 have no recurrent connectivity, and only include internal saturation dynamics that maintain activity patterns over time.
 287 With the exception of the pathway from *ctx* to *mem*, which provides unweighted contextual gating signals for attractor
 288 transitions in *mem*, inter-regional pathways (solid lines with arrows in Figure 1) have dense connectivity matrices that
 289 are learned using one-step associative learning.

290 Activity in the *mem* region evolves according to Equations 1 - 3. External multiplicative input to the *mem* region
 291 $\mathbf{c}_{mem}(t)$ comes from neural activation in the *ctx* region:

$$\mathbf{c}_{mem}(t) = \mathbf{v}_{ctx}(t)$$

292 where $\mathbf{v}_{ctx}(t)$ is the neural activation of the *ctx* region at time t . When the context gate for *mem* region dynamics is
 293 enabled ($g_{mem}^C = 1$), activity in *mem* is contextualized by activity in *ctx*. All other inputs to *mem* from other regions
 294 are absorbed into the $\mathbf{I}_{mem}(t)$ term:

$$\mathbf{I}_{mem}(t) = \sum_q g_{mem,q}(t) W_{mem,q}(t) \mathbf{v}_q(t)$$

295 where $\mathbf{v}_q(t)$ is the neural activation of region q at time t , $W_{mem,q}(t)$ is the connectivity matrix from region q to *mem* at
 296 time t , and $g_{mem,q}(t)$ is a binary connection gate on the pathway.

297 Activity in the *ctx* and *lex* regions evolves according to the following equations, which are simplified versions of
 298 Equations 1 and 3 that do not include contextual gating:

$$\mathbf{s}_r(t) = g_r^S \omega_r \mathbf{v}_r(t) + \mathbf{I}_r(t) + \sum_q g_{r,q}(t) W_{r,q}(t) \mathbf{v}_q(t) \quad (8)$$

$$\mathbf{v}_r(t+1) = \sigma_r(\mathbf{s}_r(t)) \quad (9)$$

299 where

- 300 • $\mathbf{s}_r(t)$ is a vector of cumulative synaptic input to region r at time t .
- 301 • $g_r^S(t)$ is a binary gate that enables or disables activity saturation and maintenance in region r at time t . As in
 302 Equation 1, ω_r is a scalar self-weight that establishes the fixed-point of saturated neural activity.
- 303 • $\mathbf{I}_r(t)$ is a vector of external synaptic input to region r at time t .
- 304 • $W_{r,q}(t)$ is a weight matrix for the pathway from region q to region r at time t , which is controlled by a binary
 305 connection gate $g_{r,q}(t)$.
- 306 • $\mathbf{v}_r(t)$ is a vector of neural activity in region r at time t .
- 307 • σ_r is the neural activation function for region r . The *ctx* region uses the Heaviside activation function, while the
 308 *lex* region uses the sign/signum activation function.

309 Inter-regional weight matrices in the model ($W_{r,q}$) are learned with one-step associative learning. As noted in
 310 Section 2.1, we consider both traditional Hebbian learning (Equation 10) and the store-erase rule (Equation 11):

$$\Delta W_{r,q} = \frac{1}{\rho_q^2 N_q} \sigma_r^{-1}(\mathbf{v}) \mathbf{u}^\top \quad (10)$$

$$\Delta W_{r,q} = \frac{1}{\rho_q^2 N_q} (\sigma_r^{-1}(\mathbf{v}) - W_{r,q} \mathbf{u}) \mathbf{u}^\top \quad (11)$$

311 where \mathbf{u} is the source (initial) pattern in region q , \mathbf{v} is the target (final) pattern in region r , $W_{r,q}$ is the associative
 312 weight matrix for the pathway from q to r , σ_r is the neural activation function of r , ρ_q is the activation magnitude
 313 of neurons in q , and N_q is the number of neurons in q . The context region (*ctx*) must use the heaviside activation

314 function, and learn binary patterns ($\rho = 1$), while the *lex* region may use any sign-preserving activation function such
315 as the hyperbolic tangent or sign/signum function.

316 The pathway from *mem* to *lex* ($W_{lex,mem}$) learns associations between memory states and lexical symbols stored in
317 those states. This pathway is used to “read” the contents of the currently active state in the *mem* region. The pathways
318 from *mem* and *lex* to *ctx* ($W_{ctx,mem}$ and $W_{ctx,lex}$) learn associations with patterns that contextualize memory transitions.
319 Each memory state and lexical symbol corresponds to a unique randomly generated context state. These associations
320 are critical for construction of compositional data structures in attractor graphs, as described below.

321 2.2.1. Representing Compositional Data Structures

322 The graph-organized memory represented by the *mem* region attractor graph is suitable for encoding compositional
323 data structures. Here we focus on associative arrays, linked lists, and trees, each of which can be represented by a
324 particular organization of attractors. Several instances of these data structures may be encoded simultaneously as sub-
325 graphs of a single attractor graph. Individual elements of a compositional data structure are accessed by sequences
326 of attractor transitions resembling iteration through data structures in conventional computer memory. Each memory
327 state can be associated with a *lex* activity pattern that represents the symbolic token stored in that memory state.

328 Representation of associative arrays (dictionaries, maps) in attractor graphs is straightforward. A map is repre-
329 sented by a dedicated memory state, and each entry in the map is represented by a transition to a target memory state
330 (value) that is contextualized by a *ctx* pattern (key). Keys may be derived from patterns in other regions via pathways
331 into *ctx* (e.g., *lex* or *mem*). To access the value associated with a key, the key pattern is initialized in *ctx*, and the
332 memory pattern representing the map is initialized in *mem*. Then the attractor transition is executed using the *ctx*
333 pattern as context. The resulting activity pattern in *mem* represents the memory state (value) associated with the key,
334 which may be a complex data structure or a single memory state with a corresponding symbol that is retrieved via the
335 pathway from *mem* to *lex*.

336 Accessing a map with a key that has no corresponding value results in undefined behavior. When erroneous
337 lookups are possible, the data structure must be modified to allow validation. One possibility is to learn self-transitions
338 for each value that use the key as a context pattern. Thus, to validate a map lookup, an additional transition can be
339 executed after the lookup, and the resulting activity state can be compared with the value state. If they do not match,
340 the lookup was not successful.

341 A linked list is encoded as a trajectory through the attractor graph, as shown in Figure 4a. The trajectory begins
342 with a *mem* pattern that represents the list object (referred to as the *head*), and includes zero or more additional *mem*
343 patterns that represent the list elements. The end of the list is marked by a self-loop transition in the final memory
344 state in the sequence. Although a list cannot contain repeat memory states, two states in a list may be associated with
345 the same activity pattern in *lex* via the pathway from *mem* to *lex*, representing storage of the same symbolic token in
346 two positions of the list. Each transition in the trajectory is contextualized by a shared *ctx* pattern that is specific to
347 the list, and is associated with the head pattern in *mem* via the pathway from *mem* to *ctx*. To iterate through a list,
348 the head pattern is first initialized in *mem* and used to retrieve the list-specific pattern in *ctx*. This context pattern is
349 then used to contextualize a sequence of attractor transitions that terminates when the post-transition state is identical
350 to the pre-transition state indicating a complete traversal of the sequence. Note that an empty list is represented by
351 a head pattern in *mem* that transitions directly to itself (i.e., a self-loop trajectory), and therefore terminates after a
352 single transition.

353 Because transitions in a linked list are contextualized by a list-specific context pattern, a memory state may be
354 contained in several lists. In this case, the shared memory state has distinct transitions to list-specific successors. This
355 means that memory states are *reusable components* that may be used in multiple compositional data structures. Our
356 approach is similar to Borisyuk et al. (2013) in that a list-specific context signal resolves ambiguities in sequence
357 recall that occur when an element is contained in more than one sequence. However, context signals in AGNs are
358 multiplicative masks that select subsets of the neural population to participate in sequence transitions. As a result, the
359 transitions for distinct sequences are learned in distinct but overlapping sets of connection weights.

360 Trees can be encoded in attractor graphs by recursive composition of either associative arrays or linked lists, each
361 with distinct advantages. Tree nodes represented by associative arrays have direct parent-child relations that are la-
362 beled by *ctx* patterns, which must be provided during tree traversal. However, this organization permits random access
363 of child nodes during traversal. Trees represented by linked lists, on the other hand, can be traversed without external
364 provision of context patterns. This is because each node is associated with a unique *ctx* pattern that contextualizes

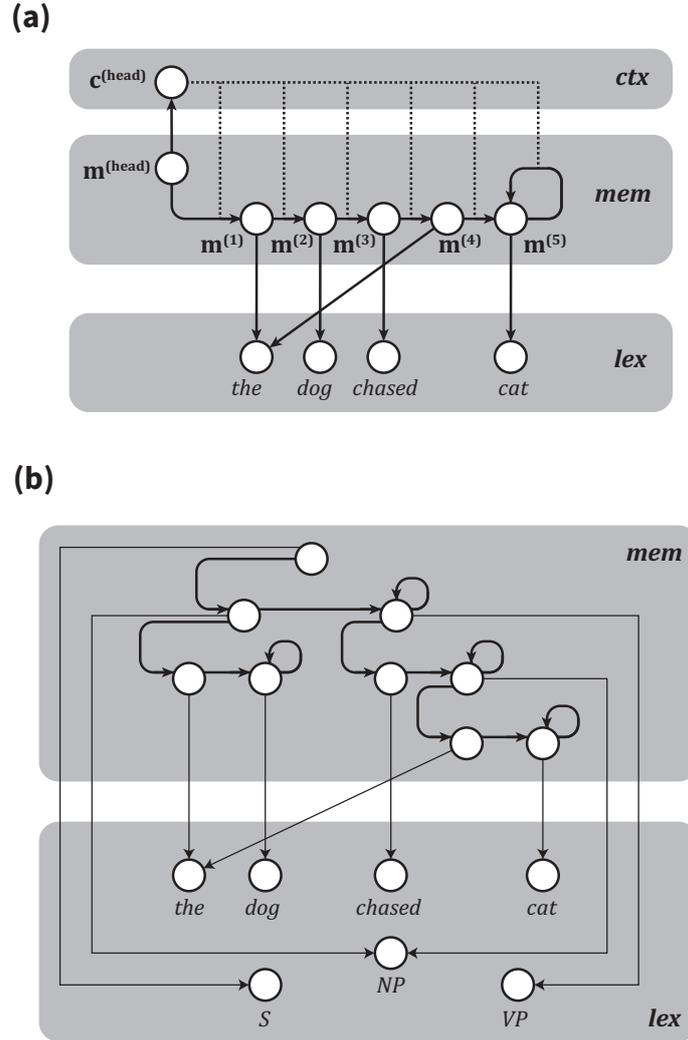


Figure 4: Graphical depiction of compositional data structures. Each gray rectangle represents the activity space of a region of the neural model (*ctx*, *mem*, and *lex*), and each circle represents a unique activity pattern (distributed representation). (a) A linked list representing the sentence “the dog chased the cat”. The list is represented by a trajectory through memory states (middle gray rectangle) that terminates with a self-loop transition. Each transition in the trajectory is contextualized by a list-specific context state ($\mathbf{c}^{(head)}$, top left) that is associated with a list head memory state ($\mathbf{m}^{(head)}$, middle left). Each element in the list is represented by a unique memory state ($\mathbf{m}^{(1)}$ through $\mathbf{m}^{(5)}$, center) that is associated with a pattern of activity in *lex* representing the corresponding word in the sentence (bottom). (b) A parse tree for “the dog chased the cat” represented as a list of lists. Each internal node of the tree is represented by a memory state that serves as the head of a list containing its child nodes (trajectories in top rectangle). The context patterns for these trajectories are omitted for clarity. Each node is associated with a pattern in *lex* representing its symbolic content (*S* for “sentence”, *NP* for “noun phrase”, *VP* for “verb phrase”, or a word in the sentence).

365 the trajectory through its children and can be retrieved during traversal. However, child nodes cannot be accessed in
 366 arbitrary order, and parent-child relations do not have associated labels. This organization is shown in Figure 4b.

367 2.3. Programmatic Control of Compositional Memory

368 The compositional data structures described in the previous section can be constructed and manipulated via top-
 369 down control of the *mem*, *ctx*, and *lex* regions over time. Figure 5 shows a programmable neural network with
 370 additional regions that provide this control. In this section we first describe the functionality of these control mecha-
 371 nisms, including the representation and execution of learned programs that make use of compositional memory. Then

372 we describe a planning task that the model learns to perform, which involves constructing and modifying complex
 373 hierarchical data structures in memory.

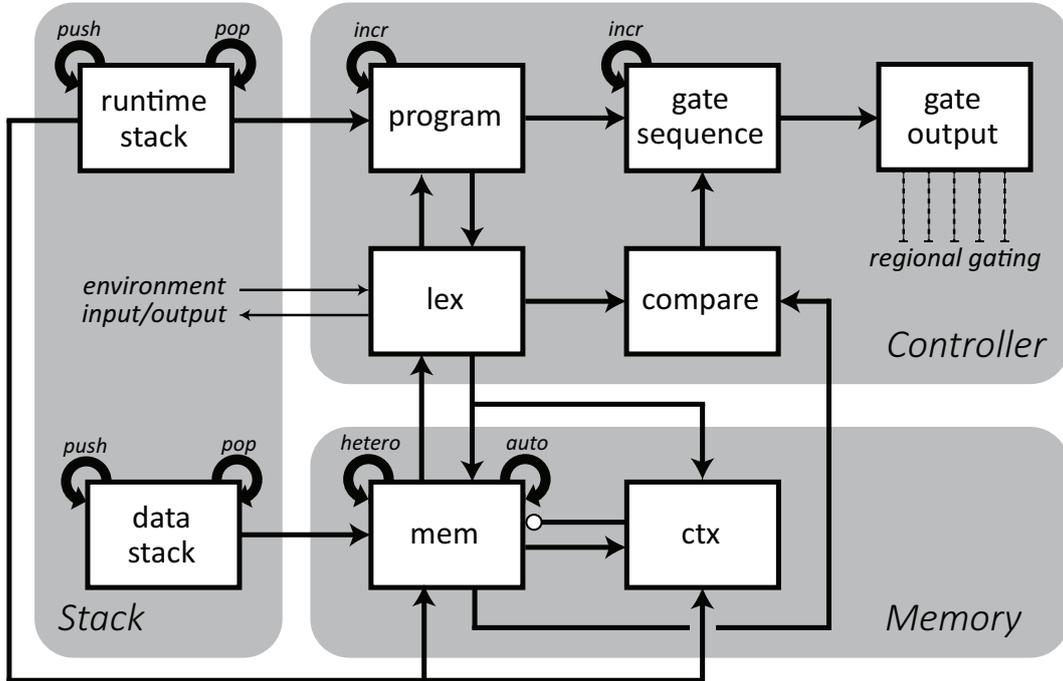


Figure 5: Programmable neural network with compositional memory. Neural regions (boxes) are interconnected with gated pathways (solid lines). The architecture of the model resembles a stack machine, and includes three subnetworks. **Controller:** The controller subnetwork (top right) controls model execution based on learned programs. Neurons in the *gate output* region (top right) determine which pathways are active during each timestep of model execution. This region is controlled by the *gate sequence* region, which encodes a sequence of gating operations that control the flow of information through the model over time. These gate sequences correspond to instruction opcodes for programs encoded in the *program* region. Instruction operands are represented in the *lex* region (center), which encodes a lexicon of recognized symbols as activity patterns. This region serves as a bridge between the controller and memory subnetworks, and is used to pass symbolic information into and out of the model (left center). **Memory:** The core region of the memory subnetwork is the *mem* region (bottom center), which is an attractor graph network with dense hetero-associative and auto-associative connectivity (bold looped arrows). Patterns of activity in *mem* represent general-purpose memory states that “store” symbolic tokens represented in *lex*. Associations between memory states and symbolic tokens are learned in the pathway from *mem* to *lex*. Transitions between attractors in *mem* are contextualized by activity in the context region (*ctx*, bottom right) via a pathway from *ctx* to *mem*. The open circle at the end of this pathway indicates that it provides multiplicative contextual gating inputs to *mem*. Activity patterns in *ctx* may be derived from *mem* or *lex* patterns, and serve as labels for relations between memory states that are used to construct compositional data structures. **Stack:** The stack subnetwork (left) contains two regions with bidirectionally associated activity patterns that represent stack frames. The *runtime stack* region stores and retrieves pointers to program instructions, memory states, and context states. These are stashed when program subroutines are called, and retrieved upon return to the caller. The *data stack* region stores pointers to memory states that are used for operations involving multiple memory states, and is used to pass arguments between subroutines.

374 The complete model depicted in Figure 5 is a programmable neural network with program-independent circuitry
 375 that is based on the Neural Virtual Machine (NVM) (Katz et al., 2019), but has several key differences. Most notably,
 376 the sequential tape-like memory of the NVM is replaced by attractor graph memory (*mem* and *ctx*), permitting storage
 377 of compositional data structures via direct context-dependent associations between memory states. In addition, our
 378 model functions as a stack machine rather than a register machine, simplifying its instruction set and program circuitry.
 379 Symbols are represented in a single region (*lex*) rather than multiple register and operand regions as in the NVM, and
 380 a *data stack* region is used for operations with multiple operands.

381 Our model includes three major subnetworks, shown as large gray rectangles in Figure 5. The memory subnetwork
 382 contains the *mem* and *ctx* regions, which implement compositional memory with attractor graphs, as described in
 383 Sections 2.1 and 2.2. The *lex* region is contained in the controller subnetwork, and serves as a bridge between several
 384 model components. Patterns in the *lex* region serve several functions, some of which were described in Section 2.2.

385 They represent symbolic tokens that can:

- 386 • be stored in memory states (*mem* to *lex*)
- 387 • be variable names that refer to data structures in memory (*lex* to *mem*)
- 388 • label transitions in *mem* attractor graphs that represent key-value relations in associative arrays (*lex* to *ctx*)
- 389 • be used as program instruction operands (*program* to *lex*)
- 390 • refer to program subroutines (*lex* to *prog*)
- 391 • be printed to or read from the environment (*lex* to/from environment)

392 The controller subnetwork also contains several regions that control connection gates based on learned programs.
393 Neurons in the *gate output* region determine which pathways in the model are active during each timestep (one neuron
394 per connection gate). This region is controlled by the *gate sequence* region, which encodes sequences of gating
395 operations that control the flow of information through the model over time, much like in conventional computer
396 architectures. These gate sequences correspond to instruction opcodes for programs encoded in the *program* region,
397 while instruction operands are represented by *lex* activity patterns.

398 The *compare* region encodes true and false patterns that are used to perform conditional jumps in programs based
399 on comparisons between memory states or lexical symbols. For example, an instruction might jump to a subroutine if
400 two memory states store the same symbol. If the comparison yields the *true* state, the model jumps to the subroutine
401 specified by the jump instruction’s operand. Otherwise, it advances to the next instruction in the sequence.

402 The stack subnetwork contains two regions with bidirectionally associated activity patterns that represent stack
403 frames. The *runtime stack* region stores a call-stack that maintains pointers to program instructions, memory states,
404 and context states. These are pushed when program subroutines are called, and popped upon return to the caller. The
405 *data stack* region maintains pointers to memory states that are used for operations involving multiple memory states.
406 For example, when a transition is learned, the target memory state is assumed to be currently active in *mem*, and a
407 pointer to the source memory state is stored on the top of the data stack. These stacks are explained further in the
408 following section.

409 2.3.1. Program Storage and Execution

410 As mentioned above, programs are represented by sequences of activity patterns in the *program* region. Each
411 activity pattern represents an individual instruction with an opcode and optional operand. Each opcode corresponds to
412 a sequence of patterns in the *gate sequence* region that implements the operation as a temporal sequence of connection
413 gates. Each instruction is associated with the first pattern in the corresponding opcode sequence via the pathway
414 from *program* to *gate sequence*. Similarly, an instruction with an operand is associated with a *lex* activity pattern
415 representing the operand value.

416 Gating operations and program sequences are established during a one-time associative learning procedure that is
417 analogous to firmware “flashing” in a non-volatile microcontroller memory. The model architecture supports opera-
418 tions that can be expressed as temporal sequences of active connection gates, such as the attractor transition procedure
419 specified in Table 1. Each timestep of an operation is represented by a randomly generated activity pattern in the
420 *gate sequence* region, and is associated with a *gate output* pattern that specifies the active connection gates. The final
421 pattern of most operation sequences is associated with a common gating sequence that advances the *program* region
422 to the next instruction, and opens the pathway from *program* to *gate sequence*, initiating the gating sequence for the
423 next instruction. The exception is jump instructions, which have conditional behavior that depends on comparison
424 operations (explained in Section 2.3.3).

425 In addition to the above associations, the “flashing” procedure also establishes associations in the stack regions
426 (*runtime stack* and *data stack*). Activity patterns in these regions represent individual stack frames that can be asso-
427 ciated with activity patterns in other regions. Each stack frame is associated with the frame above and below it in the
428 stack using distinct recurrent hetero-associative matrices (*push* and *pop* loops in Figure 5).

429 Finally, the model is “flashed” with a lexicon of recognizable symbols in the *lex* region. Each symbol pattern is
430 associated with a unique *ctx* pattern, allowing the symbol to serve as an attractor transition label (i.e., an associative
431 array key).

432 2.3.2. Online Learning

433 Program execution involves online updating of connectivity matrices. These updates are determined by distinct
434 gates that control plasticity:

$$W_{r,q}(t+1) = W_{r,q}(t) + g_{r,q}^\ell(t)\Delta W_{r,q}(t) \quad (12)$$

435 where $W_{r,q}(t)$ is a weight matrix connecting region q to region r at time t , and $g_{r,q}^\ell(t)$ is a learning gate that determines
436 when this matrix is updated. For pathways connecting distinct regions, $\Delta W_{r,q}(t)$ can be computed using the store-erase
437 learning rule with the active patterns in r and q :

$$\Delta W_{r,q}(t) = \frac{1}{\rho_q^2 N_q} \left(\sigma_r^{-1}(\mathbf{v}_r(t)) - W_{r,q}(t)\mathbf{v}_q(t) \right) \mathbf{v}_q(t)^\top \quad (13)$$

438 where $W_{r,q}(t)$ is a weight matrix connecting region q to region r at time t , $\mathbf{v}_r(t)$ and $\mathbf{v}_q(t)$ are the activation patterns of
439 r and q , σ_r is the activation function of neurons in r , ρ_q is the stable activation level of neurons in region q , and N_q is
440 the number of neurons in region q .

441 The recurrent auto-associative matrix in *mem* is also updated using the currently active pattern:

$$\Delta A_r(t) = \frac{1}{\rho_r^2 N_r} \left(\sigma_r^{-1}(\mathbf{v}_r(t)) - A_r(t)\mathbf{v}_r(t) \right) \mathbf{v}_r(t)^\top \quad (14)$$

442 where $A_r(t)$ is the auto-associative weight matrix for region r at time t , and all other terms are as defined above.

443 Online learning of recurrent hetero-associative matrices involves distinct source and target patterns that cannot be
444 simultaneously active. To address this, we introduce an eligibility trace $\epsilon_r(t)$ that stores a target activity pattern for
445 subsequent learning:

$$\epsilon_r(t+1) = \begin{cases} \mathbf{c}_r(t) \odot \mathbf{s}_r(t), & \text{if } g_r^\epsilon(t) = 1 \\ \epsilon_r(t), & \text{otherwise} \end{cases} \quad (15)$$

446 where $\epsilon_r(t)$ is the eligibility trace of region r at time t , $g_r^\epsilon(t)$ is a gate that determines when the eligibility trace is
447 updated, and $\mathbf{c}_r(t)$ and $\mathbf{s}_r(t)$ are synaptic and multiplicative inputs at time t (defined in Equations 1 and 2). When
448 $g_r^\epsilon(t) = 1$, the current gated inputs are stashed in the eligibility trace. To learn a transition to the stashed pattern, a
449 source pattern is activated, and the hetero-associative learning gate is opened:

$$\Delta H_r(t) = \frac{1}{\lambda_r \rho_r^2 N_r} \left(\epsilon_r(t) - (\mathbf{c}_r(t) \odot H_r(t)\mathbf{v}_r(t)) \right) \mathbf{v}_r(t)^\top \quad (16)$$

450 where $H_r(t)$ is the hetero-associative weight matrix for region r , $\epsilon_r(t)$ is the eligibility trace defined above, $\mathbf{c}_r(t)$ is the
451 multiplicative input to region r , λ_r is a context density for multiplicative gating of region r , and all other terms are as
452 defined above. Note that the contextual gating is already applied to the source pattern $\mathbf{v}_r(t)$ (Equation 1).

453 2.3.3. Comparisons

454 The comparison region has unique dynamics that allow it to memorize an input pattern from another region for
455 subsequent recognition. When an input pattern is memorized, it is associated with a *cmp* activity pattern representing
456 *true* (\mathbf{v}_{cmp}^{true}). Unlike other connectivity updates, this overwrites the corresponding weight matrix rather than incremen-
457 tally updating it:

$$\Delta W_{cmp,q}(t) = \frac{1}{\rho_q^2 N_q} \sigma_{cmp}^{-1}(\mathbf{v}_{cmp}^{true}) \mathbf{v}_q(t)^\top - W_{cmp,q}(t) \quad (17)$$

458 where $W_{cmp,q}(t)$ is a weight matrix connecting region q to the compare region at time t , $\mathbf{v}_q(t)$ is the activation pattern
 459 of region q , σ_{cmp} is the activation function of neurons in cmp , ρ_q is the stable activation level of neurons in region
 460 q , and N_q is the number of neurons in region q . This learning rule resembles Equation 13, except that it completely
 461 overrides the existing weights and associates the input pattern with a fixed pattern (\mathbf{v}_{cmp}^{true}).

462 The cmp region has an input bias toward a *false* activity pattern (\mathbf{v}_{cmp}^{false}) that can be overcome by the memorized
 463 input. Thus, if an input pattern is close enough to the memorized pattern, the resulting cmp activation pattern is \mathbf{v}_{cmp}^{true} ;
 464 otherwise it is \mathbf{v}_{cmp}^{false} . The following equation describes the synaptic input term for the cmp region:

$$\mathbf{s}_{cmp}(t) = \underbrace{g_{cmp}^S(t) \omega_{cmp} \mathbf{v}_{cmp}(t)}_{\text{saturation}} - \underbrace{(1 - g_{cmp}^S(t)) \theta \omega_{cmp} \mathbf{v}_{cmp}^{false}}_{\text{bias to false pattern}} + \underbrace{\sum_q (g_{cmp,q}(t) W_{cmp,q}(t) \mathbf{v}_q(t))}_{\text{inter-regional input}} \quad (18)$$

465 where \mathbf{v}_{cmp}^{false} is the *false* activity pattern in cmp , θ is a comparison similarity threshold (typically > 0.95), and all other
 466 terms are as defined for Equation 8 (with $r = cmp$). The synaptic input $\mathbf{s}_{cmp}(t)$ is transformed to neural activation
 467 $\mathbf{v}_{cmp}(t)$ by Equation 9 (again with $r = cmp$). The threshold θ determines how similar an input pattern needs to be to
 468 the memorized input pattern in order to activate the *true* pattern (\mathbf{v}_{cmp}^{true}). Specifically, a value of $\theta = 0.95$ means that
 469 an input pattern must have a cosine similarity exceeding 0.95 in order to activate \mathbf{v}_{cmp}^{true} .

470 The *true* and *false* states in cmp are associated with distinct sequences in the *gate sequence* region. When a
 471 jump instruction is executed, the pathway from cmp to *gate sequence* is opened, and the *program* region is advanced
 472 according to the result of the most recent comparison. The *true* gate sequence activates the first instruction of the
 473 subroutine indicated by the jump instruction operand. This is done by opening the pathway from *program* to *lex*,
 474 followed by the pathway from *lex* back to *program*. The *false* gate sequence simply opens the recurrent hetero-
 475 associative connection gate in the *program* region, advancing it to the next instruction in the current subroutine.

476 2.3.4. Generating Memory States

477 The model's instruction set includes operations that allocate memory for construction of data structures. Memory
 478 states are learned attractors in the *mem* region's auto-associative connectivity. We include a gated noise term in the
 479 *mem* and *ctx* regions. When a noise gate is opened, a random pattern of activation is established in the corresponding
 480 region:

$$\mathbf{v}_{mem}(t+1) = \sigma_{mem} \left(\mathbf{x}_{mem}(t) \odot \mathbf{s}_{mem}(t) + \underbrace{g_{mem}^N \omega_{mem} \mathbf{n}_{mem}(t)}_{\text{gated noise}} \right) \quad \mathbf{n}_{mem}(t) \sim \text{Bernoulli}(0.5) \quad (19)$$

$$\mathbf{v}_{ctx}(t+1) = \sigma_{ctx} \left(\mathbf{s}_{ctx}(t) + \underbrace{g_{ctx}^N \omega_{ctx} \mathbf{n}_{ctx}(t)}_{\text{gated noise}} \right) \quad \mathbf{n}_{ctx}(t) \sim \text{Bernoulli}(\lambda) \quad (20)$$

481 where g_{mem}^N and g_{ctx}^N are noise gates for the *mem* and *ctx* regions, and $\mathbf{n}_{mem}(t)$ and $\mathbf{n}_{ctx}(t)$ are random vectors. All
 482 other terms are as defined in Equation 3 (for *mem*) and Equation 9 (for *ctx*). Random vectors are generated by a
 483 Bernoulli process with probabilities 0.5 and λ (context density, defined in Section 2.1). Random patterns in *mem*
 484 are in $\{-\rho_{mem}, +\rho_{mem}\}^{N_{mem}}$, where N_{mem} is the number of neurons in *mem*, and ρ_{mem} is the steady-state magnitude of
 485 saturation dynamics in *mem* ($\rho_{mem} = \sigma_{mem}(\omega_{mem})$). Random patterns in *ctx* are binary patterns.

486 2.3.5. I/O

487 Environmental inputs to the model are provided to the *lex* region via the external input variable \mathbf{I}_{lex} (Equation 8).
 488 The model indicates when it is ready to receive input by activating a dedicated read gate g_{lex}^R . Environmental input is
 489 specified as a sequence of symbols from an alphabet that the model is pre-trained to recognize. Each alphabet symbol
 490 is mapped to a unique activity pattern. When the environment detects that $g_{lex}^R(t) = 1$, $\mathbf{I}_{lex}(t)$ is set to the activation
 491 pattern corresponding to the next unread symbol in the input sequence.

492 Similarly, a dedicated write gate g_{lex}^W indicates to the environment when the model is ready to provide output.
 493 When $g_{lex}^W(t) = 1$, the environment captures the current *lex* activity pattern $\mathbf{v}_{lex}(t)$ and translates it into a symbol by
 494 identifying the closest activity pattern in the alphabet mappings.

2.3.6. Planning Task

Planning is a high-level executive task that involves reasoning about actions and organizing them to achieve goals (Ghallab et al., 2004; Erol, 1996). Behavioral plans are often hierarchically structured and require compositional reasoning (Botvinick, 2008; Dehaene and Changeux, 1997). To further test our model’s ability to construct and maintain compositional data structures in memory, we trained it to perform an automated planning task using hierarchical task networks (HTNs) (Ghallab et al., 2004; Erol, 1996). An HTN is a tree representing the decomposition of a high-level compound task (root node) into concrete primitive actions (leaf nodes). Each internal node is broken down into sub-actions (compound or primitive) according to learned rules in a knowledge-base that may depend on environmental states. For example, opening a door may involve different motor behaviors depending on what type of door it is (e.g., pushing, pulling, sliding, etc). During planning, an agent recursively decomposes a top-level task to produce a sequence of primitive actions that is appropriate for the given environment.

We implemented a version of HTN planning with the following restrictions. The environment is represented by an associative array of named feature bindings (e.g., “door type” = “sliding”). Each compound action is decomposed according to the value of a specific environmental feature (e.g., “open door” is decomposed according to the value of “door type”). When there is no rule for a given action and feature value, the action is treated as a primitive action, and is not further decomposed. The task is performed as follows. First, the model reads in the knowledge-base rules and environmental bindings and stores them in memory. Then, it reads in a sequence of top-level actions for planning. An HTN is constructed by recursively decomposing each top-level action into primitive actions based on the knowledge-base and environmental bindings. Once the plan is complete, the model performs a pre-order traversal of the HTN and prints out the action for each node.

The knowledge-base is stored in attractor graph memory as a nested map (associative array). The keys of the top-level map are compound actions, and the values are inner maps containing decomposition rules for each action. As mentioned above, each action is decomposed according to the value of a specific environmental binding. The key for this binding is stored in the inner map memory state, and is used to query the environment for the binding’s value. This value is then used as a key for the inner map to retrieve the corresponding decomposition rule, which is represented as a linked list of sub-actions. During decomposition, knowledge-base lookups are validated as described in Section 2.2, and actions are only decomposed if a rule is successfully retrieved.

For testing purposes, we designed a planning domain to simulate a simple repair task involving a mechanical assembly unit (see Appendix for details). The unit has a door on the front, an indicator for the status of the unit, and an interaction point for performing repairs. The full task involves opening the door according to its type, performing a repair according to the status indicator, and closing the door. Each stage of the task is performed based on a set of environmental bindings describing a specific unit, including the type of door, status indicator, and repair interaction point (e.g., sliding door, LED indicator, keypad interaction point).

3. Results

The model outlined above was implemented in Python using the NumPy scientific computing library, and was tested in several stages. The first three stages evaluated the memory capacity of the AGN model, and involved learning attractor graphs in the *mem* region. Specifically, we sought to empirically determine 1) the number of attractor states that can be learned and reliably retrieved from partial patterns, 2) the number of unique context-dependent transitions from a single source pattern that can be learned (i.e., the branching factor of attractor graphs), and 3) the total number of transitions that can be learned in an attractor graph.

The fourth and fifth stages of testing evaluated storage and retrieval of compositional data structures, and included learning in the inter-regional pathways of the model. First, the model was trained with attractor graphs representing linked lists to determine whether errors in attractor transitions compound during traversal, and whether memory states can be effectively reused in multiple distinct list structures. Then, the model was trained with attractor graphs representing sentence parse trees. Each node in a tree was represented as a pattern of activity in *mem*, and each symbol stored in the node was represented by a pattern in *lex*. The associations between nodes and symbols were learned in the pathway from *mem* to *lex*.

In the final stage of testing, we evaluated the model’s ability to manipulate compositional data structures using the HTN planning task described in Section 2.3.6. Each test involved learning a knowledge-base of decomposition rules,

544 a set of environmental bindings, and a sequence of top-level compound actions. These top-level actions were decom-
545 posed into HTNs according to the knowledge-base and environmental bindings. We evaluated runtime complexity by
546 varying the top-level action sequence to vary the size of the constructed HTN. Three conditions were evaluated: 1) a
547 baseline condition using a small knowledge-base (9 rules spread across 5 compound actions) and small environment
548 (4 bindings), 2) an extended knowledge-base of 15 rules across 7 compound actions, and 3) an extended environment
549 containing 8 bindings. We validated that the model constructed the correct HTN in each test, and measured the num-
550 ber of timesteps taken to complete the task. The specific knowledge-bases, environments, and top-level sequences
551 used are listed in the Appendix.

552 In all stages of testing, performance was evaluated by comparing patterns of activity with learned target patterns.
553 For example, to evaluate a learned attractor transition from memory pattern \mathbf{m}_A to memory pattern \mathbf{m}_B , the trained
554 model would be initialized with \mathbf{m}_A in the *mem* region, the transition would be executed (including attractor conver-
555 gence), and the resulting *mem* activity pattern would be compared with \mathbf{m}_B . Because saturation dynamics can correct
556 any convergence errors within an orthant of activity space, two patterns are considered identical if they reside within
557 the same orthant (i.e., the sign of each neuron’s activation matches). For each experiment, results are reported as the
558 percentage of activity patterns that matched their corresponding targets.

559 Experiments in Sections 3.1 - 3.5 evaluated the model shown in Figure 1 using either traditional Hebbian learning
560 or the fast store-erase learning rule. Pathways with context-dependent dynamics were learned with the gated versions
561 of these rules (Equations 4 - 7). All other pathways were learned with non-gated versions of these rules (Equations
562 10 and 11. Each region of the model contained $N = 1024$ neurons. The *mem* region used the hyperbolic tangent
563 activation function and learned activity patterns with magnitude $\rho = 0.9999$. The *ctx* and *gate output* region used
564 the heaviside activation function, and learned binary patterns ($\rho = 1$). The *lex* region used the sign/signum activation
565 function ($\rho = 1$). The experiment in Section 3.6 evaluated the full programmable network shown in Figure 5. In this
566 experiment, regions were sized according to the number of patterns to be learned for the planning task described in
567 Section 2.3.6. Regions not present in Figure 1 used the sign/signum activation function ($\rho = 1$).

568 3.1. Attractor Convergence

569 Attractor transitions are carried out in multiple steps, starting with contextually-gated hetero-associative dynamics,
570 followed by auto-associative attractor convergence and activity saturation. Because the hetero-associative step results
571 in a partial pattern of activity (some neurons have zero activation), successful transitions depend on accurate auto-
572 associative pattern completion. We therefore sought to determine the number of activity patterns that can be learned
573 as attractors and successfully recovered from partial patterns.

574 The density of partial patterns (i.e., the number of non-zero elements) encountered during attractor transitions
575 depends on the parameter λ , the probability used to generate context patterns (Section 2.1). This parameter indicates
576 the number of neurons in the *mem* region that participate in hetero-associative dynamics, and consequently the number
577 of active neurons prior to attractor convergence.

578 The *mem* region AGN was trained with sets of randomly generated memory states, and evaluated for pattern
579 completion. One trial of testing involved learning a set of M memory states generated by a Bernoulli process with
580 probability 0.5 (i.e., a fair coin toss for each neuron’s activation). Each learned memory pattern was evaluated by
581 initializing the network with a partial version of the pattern, running auto-associative dynamics and saturation, and
582 comparing the resulting activity pattern with the original learned pattern. The number of timesteps for auto-associative
583 dynamics was set to 10, which was found to be sufficient for attractor convergence in preliminary testing. Partial
584 patterns were produced by randomly setting $N(1 - \lambda)$ elements to zero, simulating contextual gating. This process
585 was repeated 8 times for each memory state, and the value of λ was varied experimentally.

586 Results are shown in Figure 6. Each plot shows results for networks trained with a common learning rule (store-
587 erase or traditional Hebbian learning) and various values of λ . Each line shows accuracy of attractor convergence as
588 the number of learned memory states M increases. Accuracy deteriorated as the density of the context pattern (λ) was
589 decreased and the masked partial patterns became sparser. Results for the two learning rules were comparable, but the
590 store-erase rule showed slightly more gradual degradation with increasing numbers of stored patterns. Perfect accu-
591 racy can be achieved with the store-erase rule when $M = 64$ memory patterns are stored. In subsequent experiments,
592 we limit the number of learned attractors to 64 to avoid errors in attractor convergence.

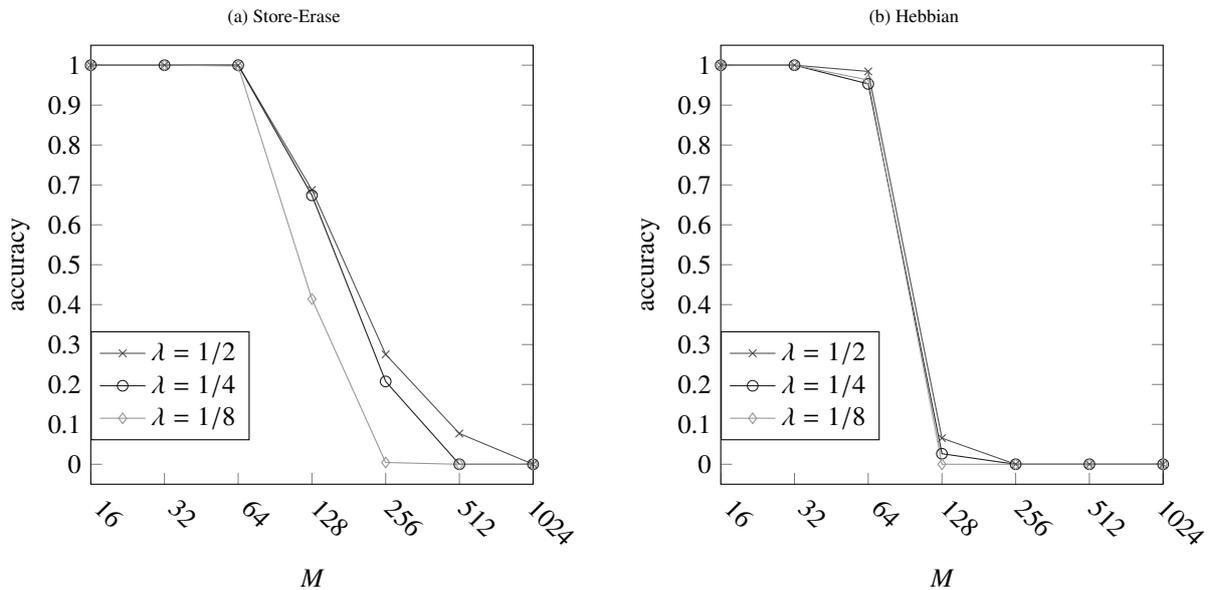


Figure 6: Accuracy of attractor convergence (pattern recovery/completion). Each plot shows the recall accuracy (y-axis) of an AGN with $N = 1024$ neurons trained with either the fast store-erase learning rule (left) or traditional Hebbian learning (right). The network was trained with sets of memory patterns of various sizes (M , x-axis). Each learned memory pattern was tested 8 times by initializing the network with a partial version of the pattern, running auto-associative dynamics for 10 timesteps, and comparing the resulting activity pattern with the learned pattern. The density of the partial version was determined by the parameter λ , and each line indicates results for a single value of λ . Each data point indicates the percentage of convergence trials resulting in perfect pattern recall, for a total of $8M$ trials per data point.

3.2. Transition Branching

Functional branching is a novel aspect of AGNs that makes it possible to learn multiple transitions from a single attractor state using a single hetero-associative weight matrix (Section 2.1.1). We refer to the number of transitions from an attractor as the attractor’s *branching factor*. We evaluated learning of attractor graphs with various branching factors to determine how many transitions can be learned from a single source attractor. Specifically, we trained AGNs with attractor graphs organized as directed stars, where one attractor serves as an internal node with transitions to several leaf nodes.

The results for attractor convergence above indicate that a network of $N = 1024$ neurons can reliably learn $M = 64$ attractors. We therefore generated attractor graphs with $M = 64$ attractors and varying numbers of transitions. One attractor was designated as the internal node, and each transition targeted an attractor randomly chosen from the remaining 63 attractors. Each transition was learned with a unique context pattern generated by a Bernoulli process with probability λ (varied experimentally). Note that there may be multiple transitions from the internal node to the same leaf node.

Results are shown in Figure 7. Each data point indicates the percentage of transitions that were successfully executed after learning. In contrast to the pattern recovery results, accuracy was higher with smaller λ . This is likely due to decreased weight sharing across contexts. The store-erase rule significantly outperformed traditional Hebbian learning, yielding high accuracy (over 97%) for branching factors up to 1024 for $\lambda = \frac{1}{4}$ and $\lambda = \frac{1}{8}$.

3.3. Random Graphs

Having established that AGNs can learn attractor graphs with high branching factors, we sought to determine how many transitions can be stored when the transitions do not share a source node. To do so, we trained the AGN with randomly generated graphs of $M = 64$ attractors and varying numbers of transitions (T). Graphs were generated by randomly selecting a source attractor (vertex), target attractor (vertex), and context pattern (edge label) for each transition (edge). The number of available context patterns was set to the maximum branching factor of nodes in the

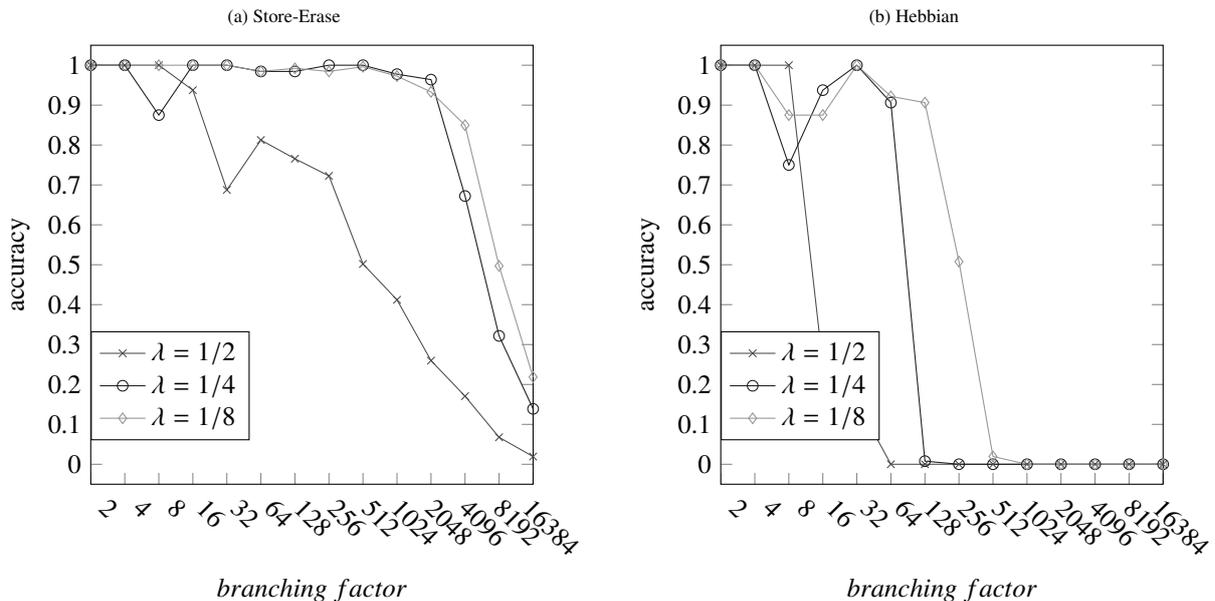


Figure 7: Accuracy of attractor graphs with various branching factors. Each plot shows the transition accuracy for an AGN with $N = 1024$ neurons trained with either the fast store-erase learning rule (left) or traditional Hebbian learning (right). The network was trained with attractor graphs containing $M = 64$ memory attractors and various numbers of transitions (x-axis). The attractor graphs were organized as directed stars, where one internal node transitions to several leaf nodes in different contexts. Multiple transitions to the same leaf node in different contexts were allowed, making it possible to test branching factors larger than the total number of leaf nodes. Each line indicates results for one value of λ , the density of context patterns used for transitions. Each data point indicates the percentage of successful transitions.

616 graph, ensuring that no two transitions shared the same source and context patterns. Transitions from different source
617 patterns were allowed to share context patterns.

618 Results are shown in Figure 8. Store-erase learning was more reliable than traditional Hebbian learning, with less
619 variable accuracy. With $\lambda = \frac{1}{4}$, both learning rules yielded high accuracy with $T = 1024$ transitions. In subsequent
620 experiments, we limit the number of learned attractor transitions to 1024.

621 3.4. Linked Lists

622 The above results evaluate the integrity of individual attractor transitions and convergence events in an AGN.
623 Iteration through attractor graphs representing compositional data structures involves sequences of transitions in which
624 errors might compound. We therefore evaluated retrieval of linked lists with itinerant traversals, allowing any errors
625 in transitions to compound.

626 The model was trained with attractor graphs of $M = 64$ attractors encoding linked lists (Section 2.2). Each
627 attractor served as the head of a unique list containing E elements (varied experimentally) drawn from the remaining
628 63 attractors. Each list is encoded as a trajectory containing $E + 1$ transitions, for a total of $64(E + 1)$ transitions in
629 the graph. Note that an attractor (memory state) may be contained in more than one list, as each list's transitions use
630 a unique list-specific context pattern. These context patterns were learned in the pathway from the *mem* to *ctx* regions
631 (Figure 1), and were retrieved at the beginning of each traversal at testing time.

632 Results are shown in Figure 9. Accuracy drops sharply after $E = 15$ elements, or a total of $T = 1024$ transitions.
633 This corresponds to the point at which accuracy declines in Figure 8. The drop in accuracy after $E = 15$ is therefore
634 likely due to limits in transition capacity. These results show that below this capacity, traversals through attractor
635 graphs can be executed without compounding errors. In addition, large numbers of linked lists can be encoded in
636 attractor graphs, and memory states can be successfully shared between lists.

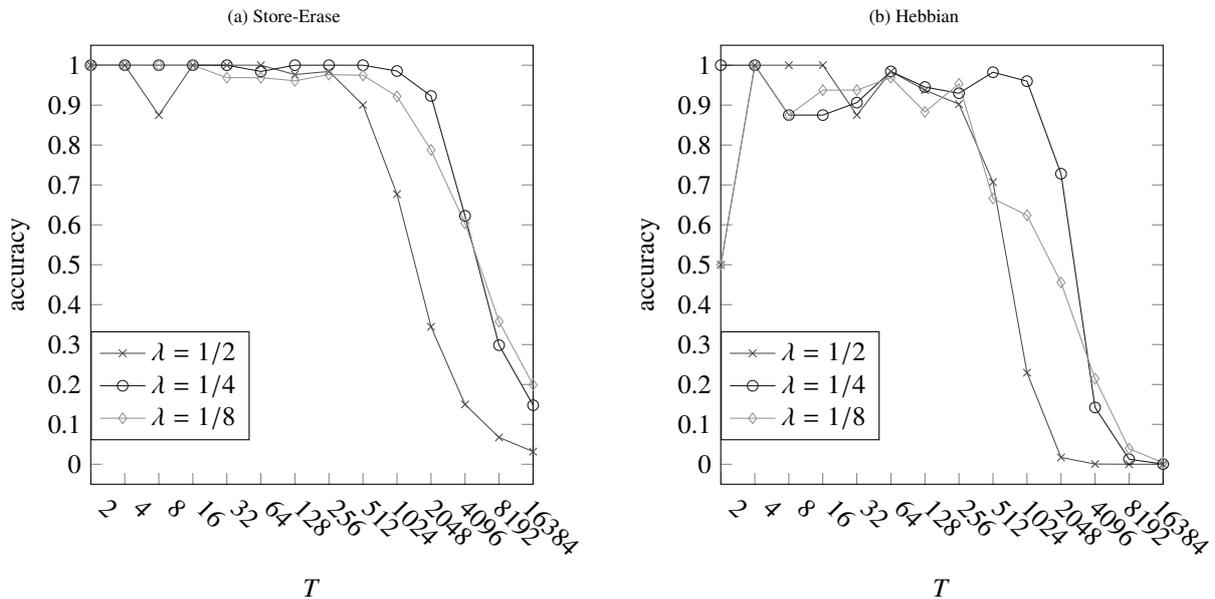


Figure 8: Accuracy of randomly generated attractor graphs. Each plot shows the transition accuracy for an AGN with $N = 1024$ neurons trained with either the fast store-erase learning rule (left) or traditional Hebbian learning (right). The network was trained with attractor graphs containing $M = 64$ memory attractors and various numbers of transitions (x-axis). Each transition in the graph connected two randomly selected attractors using a randomly generated context pattern. Each line indicates results for one value of λ , the density of context patterns used for transitions. Each data point indicates the percentage of successful transitions.

3.5. Parse Trees

The results above show that list traversals can be successfully carried out without compounding errors when the total number of attractors is limited to $M = 64$. That experiment did not evaluate storage and retrieval of symbolic information via the pathway from *mem* and *lex*, and did not involve retrieval of context patterns in *ctx* using memory states that may contain minor errors. To address these limitations, we evaluated the model with storage of parse trees with a symbol stored at each node.

Parse trees were randomly selected from the Penn Treebank corpus¹. Each tree was encoded as an attractor graph using lists of lists, where each node is represented as a list of its children. Each node in the tree was assigned to a memory attractor, and each symbol in the parse tree was assigned to a randomly generated pattern of activity in the *lex* region using a Bernoulli process with probability 0.5. Associations between nodes and symbols were learned in the pathway from *mem* to *lex* using either the store-erase learning rule or traditional Hebbian learning.

The model learned one tree at a time, and weights were reset after learning and evaluating each tree. Each tree was evaluated with a traversal starting with initialization of *mem* with the root node pattern. As with list testing, errors were allowed to propagate during traversal, and context patterns were retrieved using the pathway from *mem* to *lex*. After each transition, the symbol stored in the current node was retrieved using the pathway from *mem* to *lex*, and results are reported as the percentage of symbol patterns in *lex* that were correctly recovered. An external queue was used to store and retrieve intermediate activity patterns and perform a breadth-first traversal. Note that this queue is not considered part of the model, and is only used for evaluation purposes.

Results are shown in Figure 10. A total of 100 randomly selected trees were learned. Each mark indicates the percentage of perfectly recalled symbol patterns across all nodes in a single tree (y-axis), and the x-axis (log-scale) indicates the number of nodes in the tree. In accordance with the results in Section 3.1, accuracy begins to deteriorate with trees containing more than 64 nodes. Past this point, accuracy degrades more gradually with the store-erase rule than with traditional Hebbian learning.

¹Parse trees were retrieved from the *Penn Treebank Sample* dataset of the Python Natural Language Toolkit, found at http://www.nltk.org/nltk_data/

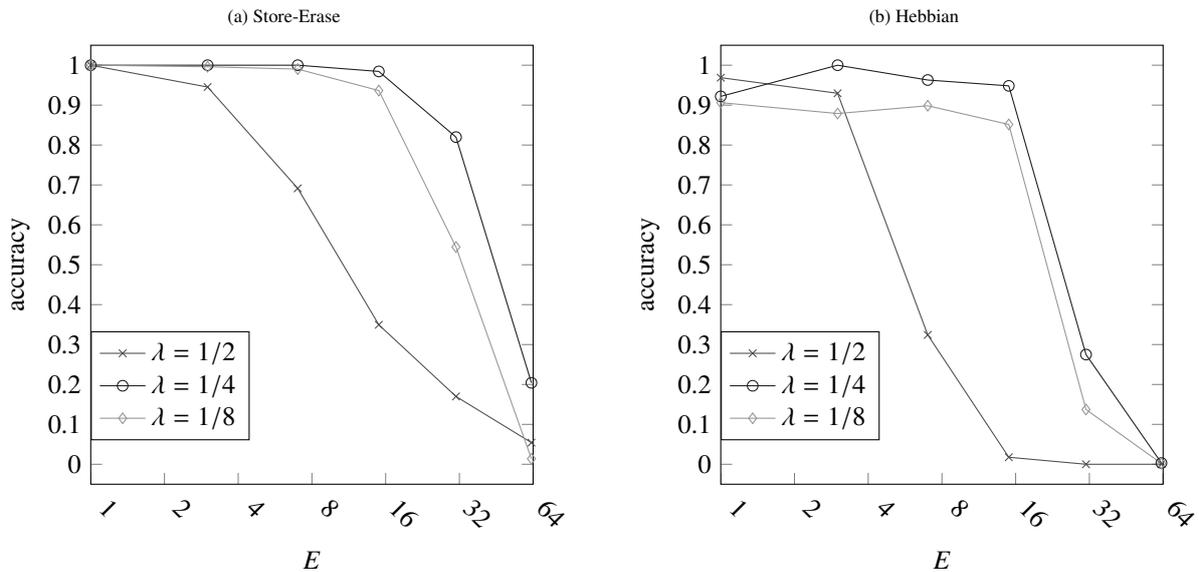


Figure 9: Accuracy of attractor graphs encoding multiple linked lists. Each plot shows the transition accuracy for an AGN with $N = 1024$ neurons trained with either the fast store-erase learning rule (left) or traditional Hebbian learning (right). The network was trained with attractor graphs containing $M = 64$ memory attractors, each representing the head of a linked list containing E elements (x-axis). Each list contained a random permutation of E attractors (not including the list’s head attractor), and was encoded as a trajectory through the attractor graph. The total number of transitions in the attractor graph is $64(E + 1)$. Performance was evaluated by executing traversals through each list starting with the head pattern, and errors in transitions were allowed to propagate during traversal. Each line corresponds to a unique context density λ , and each data point indicates the percentage of successful transitions.

660 The store-erase learning rule contains an anti-Hebbian component that erases previously learned associations. To
 661 evaluate this unique contribution, we performed a second experiment with parse trees of at most 64 nodes each. In this
 662 experiment, the model weights were not reset in between learning each tree. A set of $M = 64$ attractors was learned
 663 and made available for construction of each tree, which contained a unique set of transitions between attractors, as
 664 well as associations between attractors and context/lexicon patterns.

665 Results are shown in Figure 11 for 30 trials. The left plot shows accuracy of *lex* pattern recovery for the two
 666 learning rules with $\lambda = \frac{1}{4}$. Accuracy for Hebbian learning (dashed line) drops to zero after the first trial, as learned
 667 associations compound and interfere with one another. Because the store-erase rule allows overwriting of associations
 668 via controlled erasure, accuracy remains fairly high across the trials, but dips to as low as 75% (solid line).

669 To determine the extent of pattern deterioration, patterns in *lex* representing stored symbols were compared with a
 670 fine-grained similarity metric rather than all-or-nothing comparison. This similarity metric measured the percentage
 671 of *lex* neurons with activation matching the target symbol’s activation pattern. Results are shown in the right plot of
 672 Figure 11 for $\lambda = \frac{1}{4}$. Each data point indicates the average similarity of *lex* activity patterns retrieved from tree nodes.
 673 The solid line shows that the average similarity with the store-erase rule is nearly perfect (average similarity of 0.9998
 674 across trials) despite the dips in overall accuracy in the left plot. This indicates that errors in pattern retrieval are
 675 minimal, and involve very small numbers of neurons with activation that did not match the target pattern. In contrast,
 676 the average similarity with Hebbian learning drops rapidly to around 50%, which is the expected similarity for two
 677 randomly generated patterns.

678 3.6. Planning Task

679 To evaluate autonomous construction, access, and manipulation of compositional data structures, we tested the
 680 model using the HTN planning task outlined in Section 2.3.6. First, the controller regions of the model were “flashed”
 681 using the store-erase rule with an instruction set and set of program subroutines that implement the task (Katz et al.,
 682 2019). During testing, the model was provided with a sequence of inputs encoding a knowledge-base, environmental

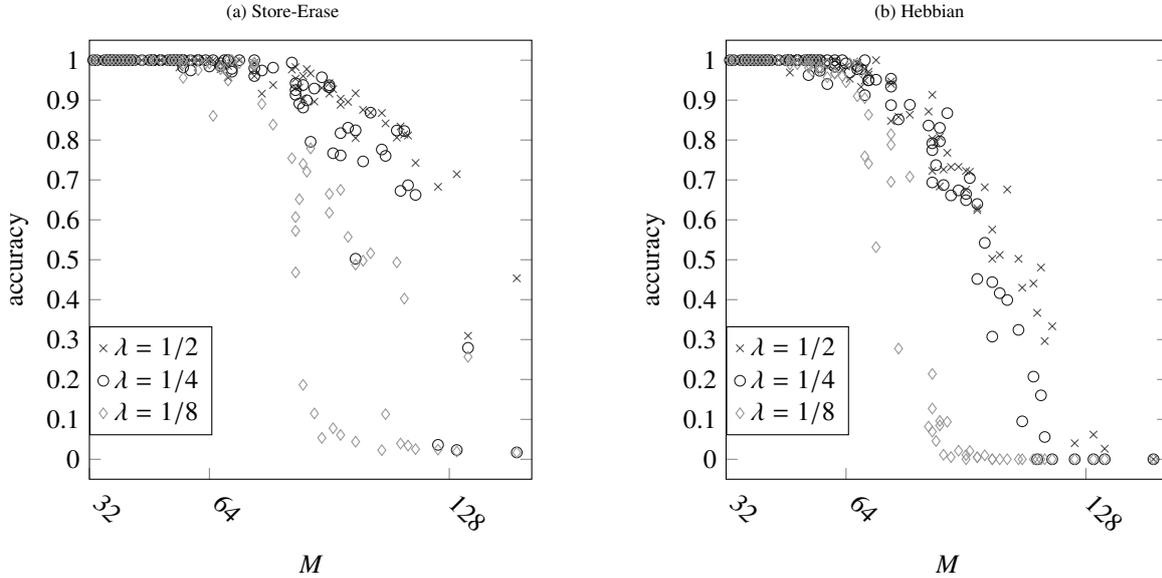


Figure 10: Accuracy of attractor graphs encoding parse trees. Each plot shows the accuracy for an AGN with $N = 1024$ neurons trained with either the fast store-erase learning rule (left) or traditional Hebbian learning (right). The network was trained with attractor graphs representing sentence parse trees drawn from the Penn Treebank. Each node was represented by an attractor in the *mem* region, and contained a symbol represented by a pattern of activity in the *lex* region. Each data point indicates the percentage of *lex* patterns successfully retrieved (y-axis) during traversal of a tree with M nodes (x-axis). Errors in attractor transitions were allowed to propagate during traversal, and an external queue was used to maintain and retrieve intermediate *mem* activation patterns to perform a breadth-first traversal.

683 bindings, and top-level actions for planning. After decomposing these actions, the model performed a pre-order
 684 traversal, printing out the resulting HTN tree.

685 The model was tested with top-level action sequences that corresponded to HTNs of various sizes. In the *baseline*
 686 condition, these tests were performed with a small knowledge-base (9 rules across 5 compound actions) and small
 687 environment (4 bindings). To determine the impact of these data structures on planning runtime, two additional
 688 conditions were considered: one with a larger knowledge-base (*extended KB*, 15 rules across 7 compound actions),
 689 and one with a larger set of environmental bindings (*extended env*, 8 bindings). We ensured that the constructed HTNs
 690 did not differ across these conditions.

691 The size of each network region was set according to the number of patterns to be represented in that region. The
 692 *mem* region contained $N_{mem} = 9216$ neurons to ensure successful storage of task-relevant data structures. The context
 693 density λ was set to 0.5, and the store-erase rule was used for online learning (Section 2.3.2).

694 The model successfully performed the task in all cases, and produced the sequence of outputs corresponding
 695 to the correct HTN tree. Figure 12 shows the number of timesteps taken during each test (multi-timestep attractor
 696 convergence events are collapsed into individual timesteps). These results show that the computational complexity
 697 of planning scales linearly with the size of the constructed HTN, independently of the size of the knowledge-base or
 698 environmental binding set. This demonstrates that associative arrays represented in attractor graph memory can be
 699 efficiently accessed, as lookups require a constant number of timesteps that is independent of the number of learned
 700 key-value pairs.

701 4. Discussion

702 We have presented a recurrent neural network model that represents compositional data structures as systems of
 703 itinerant attractors called *attractor graphs*. Our model learns context-dependent attractor transitions using a novel
 704 combination of top-down gating and one-step associative learning. Notably, this training method makes it possible to
 705 learn multiple outgoing transitions from a single attractor state using a single hetero-associative matrix. These tran-

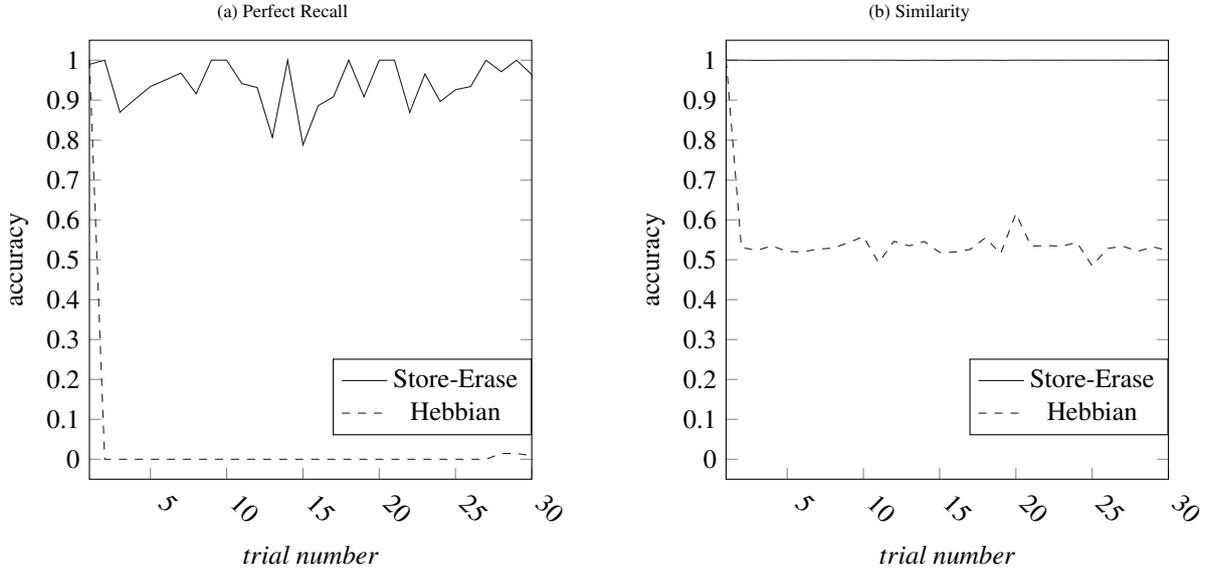


Figure 11: Learning parse trees without weight resets. The plot shows the performance an AGN with $N = 1024$ neurons trained with either the fast store-erase learning rule (solid lines) or traditional Hebbian learning (dashed lines), reported as the average similarity for *lex* activity patterns representing symbols stored in tree nodes. Each trial involved learning a single parse tree, and evaluating symbol recall. Attractor states were recycled between trees, but each tree was represented by a unique set of attractor transitions and inter-regional associations (from *mem* to *ctx* and *lex*). In between trials, the weights of the model were not reset. The size of randomly selected parse trees was limited to 64 nodes to prevent errors in attractor convergence.

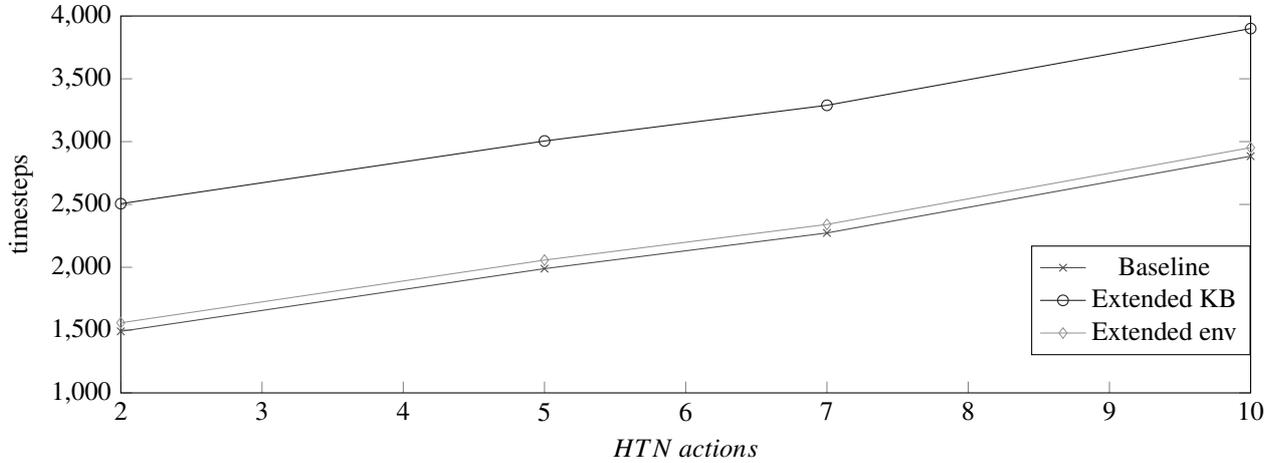


Figure 12: HTN planning runtime. The full programmable neural network was pre-trained to perform the HTN planning task, and evaluated on inputs representing a simple repair task domain (Section 2.3.6). Each line shows the number of timesteps taken to parse the domain knowledge-base and environment bindings, and perform decomposition of a sequence of high-level actions (y-axis). The x-axis indicates the number of actions in the target HTN (including internal and leaf nodes). Three conditions were evaluated. The baseline condition involved a small knowledge-base and environment. The “extended KB” condition involved a knowledge-base that was roughly twice the size of the baseline knowledge-base. The “extended env” condition involved twice the environmental bindings as the baseline condition. The results show that the computational complexity of the planning stage of the task is independent of the size of the knowledge-base and environmental bindings, and depends only on the resulting HTN tree.

706 sitions are selected during model execution by multiplicative contextual gating signals that control memory retrieval
707 and iteration through learned data structures. We refer to this as *functional branching*, as the branches in attractor
708 sequences are determined by patterns of activity and are not stored in distinct connectivity matrices.

709 Empirical results demonstrate that our model can reliably store and retrieve attractor graphs representing composi-
710 tional data structures such as associative arrays, linked lists, and trees. While the number of learned attractor states is
711 limited, our model can learn attractor graphs with large numbers of transitions (edges), and with very high branching
712 factors (vertex degrees), and individual attractors may be used as components of several data structures. We evaluated
713 two forms of one-step associative learning: traditional Hebbian learning and the fast store-erase learning rule (Katz
714 et al., 2019). While the two learning rules yielded similar memory capacities, the store-erase rule significantly outper-
715 formed Hebbian learning on attractor graphs with very high branching factors. This reflects a reduction of interference
716 across contexts that is likely due to the anti-Hebbian component of the store-erase rule, which also enables controlled
717 erasure of learned associations. The ability to erase and overwrite transitions permits rapid reorganization of attractor
718 graphs, making the network an effective model of reusable working memory.

719 We also showed that compositional data structures can be efficiently manipulated via procedural gating control in a
720 programmable neural network. The network successfully performed a hierarchical planning task involving rule-based
721 decomposition of action sequences. A significant limitation of this model is that it does not leverage compositional
722 memory to store programmatic procedures, and instead relies upon a comparably simple assembly-like language
723 with linear program sequences (Katz et al., 2019). Because attractor graphs can represent tree data structures, these
724 processes may instead be represented as abstract syntax trees for programs written in a high level programming
725 language. In addition, a unified program/data memory would allow implementation of homoiconic programming
726 languages such as Lisp and Scheme, making it possible for the model to modify learned programs and synthesize new
727 ones.

728 Another significant limitation of our model is that it does not identify opportunities to reuse existing memory
729 structures, and does not have a stable long-term memory. As mentioned in the Introduction, the limited capacity of
730 human working memory is offset by the ability to organize and store structured representations that afford access to a
731 broad range of information. We have shown how representations in working memory can be structured or “chunked”
732 according to learned programmatic procedures. The capacity of working memory would be greatly enhanced by a
733 long-term memory containing structures that can be integrated with the contents of working memory. Structures in
734 working memory could then be replaced by pointers to existing long-term memory structures, effectively “chunking”
735 them into compressed units to reduce working memory load.

736 Our model differs from contemporary machine learning approaches to compositionality in several ways. Most no-
737 tably, individual representations are fixed-point dynamical attractors learned with one-step associative learning rather
738 than the error-based gradient descent learning methods common in deep learning. These attractor states are composed
739 into complex structures with context-dependent transitions that represent relations between discrete elements in mem-
740 ory. Because these relations are stored in connectivity weights, our model does not rely on persistent maintenance of
741 multiple activity patterns. Instead, memories are retrieved as needed via top-down control of attractor transitions. This
742 “activity-silent” form of working memory has a strong basis in neuroscientific theory (Manohar et al., 2019; Mongillo
743 et al., 2008; Rose et al., 2016; Stokes, 2015; Barbosa et al., 2019), and has not previously been used for compositional
744 learning in artificial neural networks.

745 AGNs do not require specialized operations for compressing elements into structured representations, such as cir-
746 cular convolutions. Instead of creating summary vectors, AGNs learn direct relations between elements with arbitrary
747 encodings (activity patterns) using one-step associative learning. Because structure is learned in connectivity weights,
748 compositional data structures can be modified without changing any activity state encodings. This is particularly
749 advantageous for nested structures: modifications to a sub-structure do not require modifications to encapsulating
750 structures. For example, a leaf may be added to a tree without modifying the encodings of the leaf’s ancestor nodes.
751 In contrast, semantic pointers are semantically related to the content they represent, and cannot be modified without
752 creation of new semantic pointers (Blouw et al., 2016).

753 Attractor graphs are capable of representing any labeled directed multigraph that does not contain two edges
754 with a shared source node and edge label. This represents a very general class of possible data structures, including
755 associative arrays, linked lists, and trees, but also graphs with cycles that cannot be represented as semantic pointers
756 due to recursive dependencies. This expressive capability exceeds that of other attractor-based models that focus on
757 sequence learning (Rajan et al., 2016; Yamashita and Tani, 2008; Jensen, 2006), or that have architecturally separated

758 representations of each hierarchical level, as in Dynamic Field Theory (Durán et al., 2012). In our model, attractor
 759 graphs in the *mem* region can represent arbitrarily nested hierarchical structures without distinct regions for each level
 760 of the hierarchy.

761 Acknowledgements

762 This work was supported by ONR award N00014-19-1-2044.

763 Appendix A. Planning Task Domain

764 The hierarchical planning task described in Section 2.3.6 involves rule-based decomposition of sequential behav-
 765 iors according to environmental conditions. Here we provide the decomposition rules, environmental bindings, and
 766 sequences of top-level actions used for the experiments in Section 3.6.

767 Rules in the knowledge-base specify how compound actions can be decomposed into sequences of sub-actions
 768 according to properties of the environment. For example, opening a door involves different actions that depend on the
 769 type of door being opened (e.g., a sliding door is opened by grasping the handle, sliding the door open, and releasing
 770 the handle). We use the following notation to express these rules:

$$compound_action(env_value) = \begin{cases} (sub_action_1^a, sub_action_2^a, \dots) & \text{if } env_val = val^a \\ (sub_action_1^b, sub_action_2^b, \dots) & \text{if } env_val = val^b \\ \dots & \dots \end{cases}$$

771 where *compound_action* is the action to be decomposed, *env_value* is the value of the environmental binding that
 772 determines the applicable decomposition rule, $(sub_action_1^a, sub_action_2^a, \dots)$ is a sequence of sub-actions for rule *a*,
 773 and val^a is the required environmental binding value to apply rule *a*. The full set of rules is enumerated below. Rules
 774 marked with a ★ are only learned in the “extended knowledge-base” condition in Section 3.6, and are not used during
 775 decomposition of the top-level sequences that were tested.

$$\begin{aligned} open_door(door_type) &= \begin{cases} (grasp, slide_open, release) & \text{if } door_type = sliding \\ (enter_passcode, enter_open) & \text{if } door_type = electronic \\ \star (unlatch, grasp, pull_open, release) & \text{if } door_type = latch \end{cases} \\ close_door(door_type) &= \begin{cases} (grasp, slide_closed, release) & \text{if } door_type = sliding \\ (enter_close) & \text{if } door_type = electronic \\ \star (grasp, push_closed, release, latch) & \text{if } door_type = latch \end{cases} \\ check_component(indicator) &= \begin{cases} (check_led) & \text{if } indicator = led \\ \star (check_pressure_gauge) & \text{if } indicator = pressure_gauge \\ \star (check_display) & \text{if } indicator = display \end{cases} \\ check_led(led_color) &= \begin{cases} (report_working) & \text{if } led_color = green \\ (repair_component) & \text{if } led_color = yellow \\ \star (report_broken) & \text{if } led_color = red \end{cases} \\ check_display(display_reading) &= \begin{cases} \star (report_working) & \text{if } display_reading = ok \\ \star (repair_component) & \text{if } display_reading = warning \\ \star (report_broken) & \text{if } display_reading = error \end{cases} \\ check_pressure_gauge(pressure_reading) &= \begin{cases} \star (report_working) & \text{if } pressure_reading = high \\ \star (repair_component) & \text{if } pressure_reading = medium \\ \star (report_broken) & \text{if } pressure_reading = low \end{cases} \end{aligned}$$

$$\text{repair_component}(\text{interaction_point}) = \begin{cases} (\text{flip_switch}) & \text{if } \text{interaction_point} = \text{switch} \\ (\text{press_button}) & \text{if } \text{interaction_point} = \text{button} \\ \star (\text{screw_valve}) & \text{if } \text{interaction_point} = \text{valve} \end{cases}$$

776 Environmental bindings are stored in an associative array, and each binding takes the form of a simple key-value
 777 pair. The full environment is listed below. Items marked with a \star are only learned in the “extended environment”
 778 condition in Section 3.6, and are not accessed during decomposition of the test top-level sequences.

door_type : electronic
indicator : led
led_color : yellow
interaction_point : button
 \star *weather : cloudy*
 \star *pants : jeans*
 \star *time : evening*
 \star *mood : tired*

779 Four different top-level action sequences were used for testing. These sequences are listed below, along with the
 780 total number of actions contained in the resulting HTN (including internal and leaf nodes):

2 actions : (repair_component)
5 actions : (open_door, close_door)
7 actions : (open_door, press_button, report_repaired, close_door)
10 actions : (open_door, check_component, close_door)

781 References

- 782 Aizawa, K., 2003. The productivity of thought, in: The Systematicity Arguments. Springer, pp. 43–55.
 783 Andreas, J., Rohrbach, M., Darrell, T., Klein, D., 2016. Neural module networks, in: Proceedings of the IEEE Conference on Computer Vision
 784 and Pattern Recognition, pp. 39–48.
 785 Ba, J., Hinton, G.E., Mnih, V., Leibo, J.Z., Ionescu, C., 2016. Using fast weights to attend to the recent past, in: Advances in Neural Information
 786 Processing Systems, pp. 4331–4339.
 787 Baan, J., Leible, J., Nikolaus, M., Rau, D., Ulmer, D., Baumgärtner, T., Hupkes, D., Bruni, E., 2019. On the realization of compositionality in neural
 788 networks, in: Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP, pp. 127–137.
 789 Baars, B.J., 2007. The global workspace theory of consciousness. The Blackwell Companion to Consciousness , 236–246.
 790 Baars, B.J., Franklin, S., 2003. How conscious experience and working memory interact. Trends in Cognitive Sciences 7, 166–172.
 791 Baayen, R.H., 1994. Productivity in language production. Language and Cognitive Processes 9, 447–469.
 792 Baddeley, A., 1993. Working memory and conscious awareness. Theories of Memory , 11–28.
 793 Bahdanau, D., Cho, K., Bengio, Y., 2014. Neural machine translation by jointly learning to align and translate. arXiv preprint arXiv:1409.0473 .
 794 Barbosa, J., Stein, H., Martinez, R., Galan, A., Adam, K., Li, S., Valls-Solé, J., Constantinidis, C., Compte, A., 2019. Interplay between persistent
 795 activity and activity-silent dynamics in prefrontal cortex during working memory. bioRxiv , 763938.
 796 Barsalou, L.W., 1993. Flexibility, structure, and linguistic vagory in concepts: Manifestations of compositional system of perceptual symbols.
 797 Theories of Memory , 29.
 798 Besold, T.R., Garcez, A.d., Bader, S., Bowman, H., Domingos, P., Hitzler, P., Kühnberger, K.U., Lamb, L.C., Lowd, D., Lima, P.M.V., et al., 2015.
 799 Neural-symbolic learning and reasoning: A survey and interpretation, in: Knowledge Representation and Reasoning: Integrating Symbolic and
 800 Neural Approaches: Papers from the 2015 AAAI Spring Symposium.
 801 Bienenstock, E., Geman, S., Potter, D., 1997. Compositionality, MDL priors, and object recognition, in: Advances in Neural Information Processing
 802 Systems, pp. 838–844.
 803 Blouw, P., Solodkin, E., Thagard, P., Eliasmith, C., 2016. Concepts as semantic pointers: A framework and computational model. Cognitive
 804 Science 40, 1128–1162.
 805 Borisyuk, R., Chik, D., Kazanovich, Y., da Silva Gomes, J., 2013. Spiking neural network model for memorizing sequences with forward and
 806 backward recall. Biosystems 112, 214–223.
 807 Botvinick, M.M., 2008. Hierarchical models of behavior and prefrontal function. Trends in Cognitive Sciences 12, 201–208.

808 Bunel, R., Hausknecht, M., Devlin, J., Singh, R., Kohli, P., 2018. Leveraging grammar and reinforcement learning for neural program synthesis.
809 arXiv preprint arXiv:1805.04276 .

810 Burke, M., Penkov, S., Ramamoorthy, S., 2019. From explanation to synthesis: Compositional program induction for learning from demonstration.
811 Robotics: Science and Systems XV doi:10.15607/RSS.2019.XV.015.

812 Campitelli, G., Gobet, F., Head, K., Buckley, M., Parker, A., 2007. Brain localization of memory chunks in chessplayers. *International Journal of*
813 *Neuroscience* 117, 1641–1659.

814 Chabuk, T., Reggia, J.A., 2013. The added value of gating in evolved neurocontrollers, in: *The 2013 International Joint Conference on Neural*
815 *Networks (IJCNN)*, pp. 1–8. doi:10.1109/IJCNN.2013.6706895.

816 Chung, J., Gulcehre, C., Cho, K., Bengio, Y., 2014. Empirical evaluation of gated recurrent neural networks on sequence modeling. arXiv preprint
817 arXiv:1412.3555 .

818 Colom, R., Rebollo, I., Palacios, A., Juan-Espinoso, M., Kyllonen, P.C., 2004. Working memory is (almost) perfectly predicted by g. *Intelligence*
819 32, 277–296.

820 Conway, A.R., Kane, M.J., Engle, R.W., 2003. Working memory capacity and its relation to general intelligence. *Trends in Cognitive Sciences* 7,
821 547–552.

822 Cowan, N., 2010. The magical mystery four: How is working memory capacity limited, and why? *Current Directions in Psychological Science*
823 19, 51–57.

824 Danihelka, I., Wayne, G., Uria, B., Kalchbrenner, N., Graves, A., 2016. Associative long short-term memory, in: *International Conference on*
825 *Machine Learning*, pp. 1986–1994.

826 Dehaene, S., Changeux, J.P., 1997. A hierarchical neuronal network for planning behavior. *Proceedings of the National Academy of Sciences* 94,
827 13293–13298.

828 Durán, B., Sandamirskaya, Y., Schöner, G., 2012. A dynamic field architecture for the generation of hierarchically organized sequences, in:
829 *International Conference on Artificial Neural Networks*, Springer. pp. 25–32.

830 Eliasmith, C., Stewart, T.C., Choo, X., Bekolay, T., DeWolf, T., Tang, Y., Rasmussen, D., 2012. A large-scale model of the functioning brain.
831 *Science* 338, 1202–1205.

832 Erilhagen, W., Schöner, G., 2002. Dynamic field theory of movement preparation. *Psychological Review* 109, 545.

833 Erol, K., 1996. Hierarchical task network planning: formalization, analysis, and implementation. Ph.D. thesis.

834 Fodor, J.A., Pylyshyn, Z.W., 1988. Connectionism and cognitive architecture: A critical analysis. *Cognition* 28, 3–71. doi:10.1016/
835 0010-0277(88)90031-5.

836 Gayler, R.W., 2003. Vector symbolic architectures answer jackendoff’s challenges for cognitive neuroscience, in: *International Conference on*
837 *Cognitive Science*, Citeseer.

838 Ghallab, M., Nau, D., Traverso, P., 2004. *Automated Planning: Theory and Practice*. Elsevier.

839 Graves, A., Wayne, G., Danihelka, I., 2014. Neural turing machines. arXiv preprint arXiv:1410.5401 .

840 Graves, A., Wayne, G., Reynolds, M., Harley, T., Danihelka, I., Grabska-Barwińska, A., Colmenarejo, S.G., Grefenstette, E., Ramalho, T., Agapiou,
841 J., et al., 2016. Hybrid computing using a neural network with dynamic external memory. *Nature* 538, 471–476.

842 Hauge, T.C., Katz, G.E., Davis, G.P., Huang, D.W., Reggia, J.A., Gentili, R.J., 2020. High-level motor planning assessment during performance of
843 complex action sequences in humans and a humanoid robot. *International Journal of Social Robotics* , 1–18.

844 Hauge, T.C., Katz, G.E., Davis, G.P., Jaquess, K.J., Reinhard, M.J., Costanzo, M.E., Reggia, J.A., Gentili, R.J., 2019. A novel application of
845 levenshtein distance for assessment of high-level motor planning underlying performance during learning of complex motor sequences. *Journal*
846 *of Motor Learning and Development* 1, 1–20.

847 Hochreiter, S., Schmidhuber, J., 1997. Long short-term memory. *Neural Computation* 9, 1735–1780.

848 Hoshino, O., Usuba, N., Kashimori, Y., Kambara, T., 1997. Role of itinerancy among attractors as dynamical map in distributed coding scheme.
849 *Neural Networks* 10, 1375–1390.

850 Hupkes, D., Dankers, V., Mul, M., Bruni, E., 2020. Compositionality decomposed: How do neural networks generalise? *Journal of Artificial*
851 *Intelligence Research* 67, 757–795.

852 Hupkes, D., Singh, A., Korrel, K., Kruszewski, G., Bruni, E., 2018. Learning compositionally through attentive guidance. arXiv preprint
853 arXiv:1805.09657 .

854 Jaeggi, S.M., Buschkuhl, M., Jonides, J., Perrig, W.J., 2008. Improving fluid intelligence with training on working memory. *Proceedings of the*
855 *National Academy of Sciences* 105, 6829–6833.

856 Jensen, O., 2006. Maintenance of multiple working memory items by temporal segmentation. *Neuroscience* 139, 237–249.

857 Kalyan, A., Mohta, A., Polozov, O., Batra, D., Jain, P., Gulwani, S., 2018. Neural-guided deductive search for real-time program synthesis from
858 examples. arXiv preprint arXiv:1804.01186 .

859 Kamp, H., Partee, B., 1995. Prototype theory and compositionality. *Cognition* 57, 129–191.

860 Katz, G.E., Davis, G.P., Gentili, R.J., Reggia, J.A., 2019. A programmable neural virtual machine based on a fast store-erase learning rule. *Neural*
861 *Networks* 119, 10–30.

862 Kipf, T., Li, Y., Dai, H., Zambaldi, V., Sanchez-Gonzalez, A., Grefenstette, E., Kohli, P., Battaglia, P., 2019. Compile: Compositional imitation
863 learning and execution, in: *International Conference on Machine Learning (ICML)*, pp. 3418–3428.

864 Lake, B., Baroni, M., 2018. Generalization without systematicity: On the compositional skills of sequence-to-sequence recurrent networks, in:
865 *International Conference on Machine Learning*, pp. 2873–2882.

866 Lake, B.M., 2019. Compositional generalization through meta sequence-to-sequence learning, in: *Advances in Neural Information Processing*
867 *Systems*, pp. 9788–9798.

868 Lake, B.M., Ullman, T.D., Tenenbaum, J.B., Gershman, S.J., 2017. Building machines that learn and think like people. *Behavioral and Brain*
869 *Sciences* 40.

870 Loula, J., Baroni, M., Lake, B.M., 2018. Rearranging the familiar: Testing compositional generalization in recurrent networks. arXiv preprint
871 arXiv:1807.07545 .

872 Manohar, S.G., Zokaei, N., Fallon, S.J., Vogels, T., Husain, M., 2019. Neural mechanisms of attending to items in working memory. *Neuroscience*

873 & Biobehavioral Reviews .

874 Marcus, G., 2018. Deep Learning: A Critical Appraisal. arXiv:1801.00631 .

875 Marcus, G., 2020. The next decade in ai: four steps towards robust artificial intelligence. arXiv preprint arXiv:2002.06177 .

876 Masse, N.Y., Grant, G.D., Freedman, D.J., 2018. Alleviating catastrophic forgetting using context-dependent gating and synaptic stabilization.

877 Proceedings of the National Academy of Sciences 115, E10467–E10475.

878 Miconi, T., Stanley, K., Clune, J., 2018. Differentiable plasticity: training plastic neural networks with backpropagation, in: International Confer-

879 ence on Machine Learning, pp. 3559–3568.

880 Miller, P., 2016. Itinerancy between attractor states in neural systems. Current Opinion in Neurobiology 40, 14–22.

881 Mongillo, G., Barak, O., Tsodyks, M., 2008. Synaptic theory of working memory. Science 319, 1543–1546.

882 Nefdt, R.M., 2020. A puzzle concerning compositionality in machines. Minds and Machines , 1–29.

883 Oberauer, K., 2009. Design for a working memory. Psychology of Learning and Motivation 51, 45–100.

884 Pelletier, F.J., 1994. The principle of semantic compositionality. Topoi 13, 11–24.

885 Pham, T., Tran, T., Venkatesh, S., 2018. Graph memory networks for molecular activity prediction, in: 2018 24th International Conference on

886 Pattern Recognition (ICPR), IEEE, pp. 639–644.

887 Piantadosi, S.T., Tenenbaum, J.B., Goodman, N.D., 2016. The logical primitives of thought: Empirical foundations for compositional cognitive

888 models. Psychological Review 123, 392.

889 Plate, T.A., 1995. Holographic reduced representations. IEEE Transactions on Neural networks 6, 623–641.

890 Rajan, K., Harvey, C.D., Tank, D.W., 2016. Recurrent network models of sequence generation and memory. Neuron 90, 128–142.

891 Reggia, J.A., Katz, G.E., Davis, G.P., 2019. Modeling working memory to identify computational correlates of consciousness. Open Philosophy

892 2, 252–269.

893 Reggia, J.A., Katz, G.E., Davis, G.P., 2020. Artificial conscious intelligence. Journal of Artificial Intelligence and Consciousness 7, 95–107.

894 Reverberi, C., Görden, K., Haynes, J.D., 2012. Compositionality of rule representations in human prefrontal cortex. Cerebral Cortex 22, 1237–

895 1246.

896 Riesenhuber, M., Poggio, T., 1999. Hierarchical models of object recognition in cortex. Nature Neuroscience 2, 1019–1025.

897 Rikhye, R.V., Gilra, A., Halassa, M.M., 2018. Thalamic regulation of switching between cortical representations enables cognitive flexibility.

898 Nature Neuroscience 21, 1753–1763.

899 Rose, N.S., LaRocque, J.J., Riggall, A.C., Gosseries, O., Starrett, M.J., Meyering, E.E., Postle, B.R., 2016. Reactivation of latent working memories

900 with transcranial magnetic stimulation. Science 354, 1136–1139.

901 Sandamirskaya, Y., Zibner, S.K., Schneegans, S., Schöner, G., 2013. Using dynamic field theory to extend the embodiment stance toward higher

902 cognition. New Ideas in Psychology 31, 322–339.

903 Silver, D., Huang, A., Maddison, C.J., Guez, A., Sifre, L., Van Den Driessche, G., Schrittwieser, J., Antonoglou, I., Panneershelvam, V., Lanctot,

904 M., et al., 2016. Mastering the game of go with deep neural networks and tree search. Nature 529, 484–489.

905 Stewart, T.C., Bekolay, T., Eliasmith, C., 2011. Neural representations of compositional structures: Representing and manipulating vector spaces

906 with spiking neurons. Connection Science 23, 145–153.

907 Stokes, M.G., 2015. ‘activity-silent’ working memory in prefrontal cortex: a dynamic coding framework. Trends in Cognitive Sciences 19, 394–405.

908 Sukhbaatar, S., Weston, J., Fergus, R., et al., 2015. End-to-end memory networks, in: Advances in Neural Information Processing Systems, pp.

909 2440–2448.

910 Sylvester, J., Reggia, J., 2016. Engineering neural systems for high-level problem solving. Neural Networks 79, 37–52. doi:10.1016/j.neunet.

911 2016.03.006.

912 Sylvester, J., Reggia, J., Weems, S., Bunting, M., 2013. Controlling working memory with learned instructions. Neural Networks 41, 23–38.

913 Szabó, Z., 2012. The case for compositionality. The Oxford Handbook of Compositionality 64, 80.

914 Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I., 2017. Attention is all you need, in: Advances

915 in Neural Information Processing Systems, pp. 5998–6008.

916 Vecoven, N., Ernst, D., Wehenkel, A., Drion, G., 2020. Introducing neuromodulation in deep neural networks to learn adaptive behaviours. PloS

917 One 15, e0227922.

918 Van der Velde, F., Kamps, M.d., 2006. Neural blackboard architectures of combinatorial structures in cognition. Behavioral and Brain Sciences

919 29, 37–70.

920 Witkin, A.P., Tenenbaum, J.M., 1983. On the role of structure in vision, in: Human and Machine Vision. Elsevier, pp. 481–543.

921 Yamashita, Y., Tani, J., 2008. Emergence of functional hierarchy in a multiple timescale neural network model: a humanoid robot experiment.

922 PLoS Comput Biol 4, e1000220.