
Winning the lottery with neural connectivity constraints: faster learning across cognitive tasks with spatially constrained sparse RNNs

Mikhail Khona*
Physics
MIT
mikail@mit.edu

Sarthak Chandra*
Brain and Cognitive Sciences
MIT
sarthakc@mit.edu

Joy J. Ma
Physics
MIT
joym@mit.edu

Ila Rani Fiete
Brain and Cognitive Sciences
MIT
fiete@mit.edu

Abstract

Recurrent neural networks (RNNs) are often used to model circuits in the brain, and can solve a variety of difficult computational problems requiring memory, error-correction, or selection [Hopfield, 1982, Maass et al., 2002, Maass, 2011]. However, fully-connected RNNs contrast structurally with their biological counterparts, which are extremely sparse ($\sim 0.1\%$). Motivated by the neocortex, where neural connectivity is constrained by physical distance along cortical sheets and other synaptic wiring costs, we introduce locality masked RNNs (LM-RNNs) that utilize task-agnostic predetermined graphs with sparsity as low as 4%. We study LM-RNNs in a multitask learning setting relevant to cognitive systems neuroscience with a commonly used set of tasks, 20-Cog-tasks [Yang et al., 2019]. We show through *reductio ad absurdum* that 20-Cog-tasks can be solved by a small pool of separated autapses that we can mechanistically analyze and understand. Thus, these tasks fall short of the goal of inducing complex recurrent dynamics and modular structure in RNNs. We next contribute a new cognitive multi-task battery, Mod-Cog, consisting of upto 132 tasks that expands by ~ 7 -fold the number of tasks and task-complexity of 20-Cog-tasks. Importantly, while autapses can solve the simple 20-Cog-tasks, the expanded task-set requires richer neural architectures and continuous attractor dynamics. On these tasks, we show that LM-RNNs with an optimal sparsity result in faster training and better data-efficiency than fully connected networks.

1 Introduction

Connectivity in biological neural networks are constrained by wiring-length costs, with a preference towards shorter synapses. The neocortex is effectively a two dimensional sheet which geometrically restricts physical distances between neurons — correspondingly, the spatial extent of connectivity in cortical circuits has been found to be significantly skewed towards shorter connections[Ercsey-Ravasz et al., 2013, Markov et al., 2013, Theodoni et al., 2022] (in particular connectivity extent appears to follow an exponential distribution).

Inspired by such biological constraints, we modify the architecture of a vanilla fully-connected RNN in one particular fashion — we construct a fixed sparse graph chosen by allowing local connections

among neurons laid on a two-dimensional sheet. We then use this sparse graph for the RNN, by only training weights between nodes that correspond to edges on the graph, and setting all other weights to zero. We refer to an RNN in this set up as a ‘Locality Masked RNN’ (LM-RNN). Another motivation for our work comes from continuous attractor network models of grid cells [Khona et al., 2022], where it was shown analytically that fixed and local topographic connectivity encouraged the formation of discrete modules.

Exploiting such simple locality constraints in LM-RNNs, when applied to multitask regimes relevant to cognitive systems neuroscience, results in distinct advantages: these networks require far fewer parameters to train, there is no cost to performance relative to an unconstrained network, and in fact learning is more rapid, sample-efficient and can achieve higher asymptotic performance. In contrast to other machine-learning approaches to network sparsification (see related work), we do not need any sophisticated pruning methods, any algorithms to construct and modify sparse skeletons nor any training data.

We focus our results on multitasking regimes for RNNs, by training them to simultaneously learn many cognitive tasks. We expect our results on the improvements conferred by LM-RNNs to primarily apply in this multitask learning setting, where the recognition of modular structure across tasks is important for generalization and effective learning.

Our main contributions are summarized by:

- We show that LM-RNNs can perform as well or better than dense networks, when accounting for the total number of nodes or the total number of synapses. Locality masking is thus an efficient prescription for choosing sparse subnetworks in a task agnostic and data independent fashion while still achieving high performance.
- We show that LM-RNNs reach this high performance faster and with lesser training than dense networks, indicating that sparse networks may be preferable to dense networks in memory and data-limited regimes for learning multiple tasks.
- We show that the tasks defined in [Yang et al., 2019] (an increasingly commonly used set of 20 tasks [Driscoll et al., 2022, Hummos, 2022, Flesch et al., 2022, Marton et al., 2021, Riveland and Pouget, 2022, da Costa et al., 2019, Duncker et al., 2020, Masse et al., 2022, Kao et al., 2021] to study representations in networks performing many cognitive tasks) can be solved by a small pool of unconnected autapses, which is essentially a feedforward structure.
- Despite not have any lateral recurrent connections, the pool of autapses shows the existence of cell clusters according to the nodal task variance metric used in [Yang et al., 2019]. Thus, our work reveals the limitations of using correlations or covariances between neurons to study recurrent mechanisms of computation. The shared inputs and input weights are a fundamental confound when using metrics such as nodal task variance used in [Yang et al., 2019]. We thus highlight the need for creation of better metrics.
- We mechanistically study how this pool of unconnected autapses solves 20-Cog-tasks.
- We then introduce Mod-Cog, a large battery consisting of upto 132 tasks inspired by cognitive science problems such as interval estimation, mental navigation and sequence generation which provides a useful setting to examine multitask learning and representation across tasks relevant to cognitive systems neuroscience.

1.1 Related work

Recent work has shown that effective pruning in recurrent networks can be done by biologically pruning plausible algorithms based on noise correlations between the presynaptic and postsynaptic neurons [Moore and Chaudhuri, 2020] that preserve the spectrum. In vision neuroscience, recent works have studied cortical topography by using a pretrained vision frontend with a readout layer with an additional spatial loss to encourage spatially nearby cells to have correlated receptive fields [Finzi et al., 2022, Lee et al., 2020a, Obeid and Konkle, 2021]. Other work has studied representations in RNNs trained to do simple tasks that have been embedded in 3-dimensional space [].

Several works in the realm of machine learning have tried to operationalize the idea of sparsity. These methods can be roughly categorized into 2 main classes:

Dense-to-sparse Ref. [Han et al., 2015] experimentally showed that training followed by pruning and retraining can give sparse networks with no loss of accuracy and ignited an interest in pruning methods. Following this work, the lottery ticket hypothesis states that dense, randomly-initialized, feed-forward networks contain subnetworks (“winning tickets”) that — when trained in isolation — reach test accuracy comparable to the original network in a similar number of iterations [Frankle and Carbin, 2018]. The best method to identify such winning tickets is Iterative Magnitude-based Pruning (IMP) [Frankle and Carbin, 2018, Frankle et al., 2019], which is computationally expensive and has to be run thoroughly for every different network. It has also been shown that parameters of the sparse initialization distribution and sign of weights at initialization are important factors [Zhou et al., 2019] which determine winning tickets. Overall, iterative pruning and retraining methods involve 3 steps: (1) pre-training a dense all-to-all model, (2) pruning synapses based on some criteria, and (3) re-training the pruned model to improve performance. This cycle needs to be done at least once and in many cases, multiple times, to get good performance. So this procedure requires at least the same training cost as training a dense model and often even more than that. In contrast, our proposed method for RNNs does not require seeing any data before pruning and trains over only the sparse remaining synapses; thus we do not require multiple cycles of pruning and training.

Other methods involving ways to encourage sparsity during the training process include L_1 (Lasso) regularization [Wen et al., 2018], L_0 regularization [Louizos et al., 2017, Savarese et al., 2020] and pruning using dynamically varying thresholds [Narang et al., 2017, Kusupati et al., 2020]. Unfortunately, all of the aforementioned methods require training the original dense network, in varying amounts, thus precluding the benefits that can be obtained by having a predetermined exact sparsity on the computation during training.

A class of methods which do not involve training data like SynFlow [Tanaka et al., 2020], GraSP [Wang et al., 2020], SNIP [Lee et al., 2019, 2020b] and FORCE [de Jorge et al., 2020] have been studied for only feedforward networks, while we study RNNs.

Sparse-to-sparse Another line of work concerning sparse-to-sparse training is most relevant to our study. This involves using a sparse interaction graph which is used to mask gradient updates. Older works maintained a static graph [Mocanu et al., 2016] and dealt only with feedforward networks but newer methods such as dynamic sparse training (DST) [Evci et al., 2020, Liu et al., 2021] have been proposed for both feedforward networks and RNNs which dynamically improve the sparse graph and provide better performance. These methods generally involve changing the topology of the sparse graph during training. Ref. [Liu et al., 2021] considers static Erdos-Renyi (ER) type sparse RNNs but for relatively denser values of sparsity (0.53 and 0.67) and more complex architectures like stacked LSTMs and Recurrent Highway Networks [Zilly et al., 2017]. Here we explore more extreme values of sparsity ($\sim 5\%$ and below) and show that they are optimal in the context of multitask learning regimes.

Lastly, static sparse networks have also found common usage in reservoir computing architectures, where large sparse networks are preferred to fully-connected networks to increase heterogeneity across nodes and allow for “richer” dynamics [Jaeger, 2001, Lukoševičius and Jaeger, 2009].

2 Results

2.1 Locality masked RNNs (LM-RNNs)

We restrict ourselves to simple RNNs for interpretability in the context of systems neuroscience. Our RNNs follow dynamics defined by:

$$\begin{aligned} \mathbf{h}_{t+1} &= \phi(\mathbf{W}\mathbf{h}_t + \mathbf{W}_{in}\mathbf{u}_t + \mathbf{b}^h), \\ \mathbf{o}_{t+1} &= \mathbf{W}_{out}\mathbf{h}_{t+1} + \mathbf{b}^o. \end{aligned}$$

Corresponding to the biological arrangement on neurons on a two-dimensional cortical sheet, we arrange the nodes of an RNN on the lattice points of a two-dimensional plane, as shown in Fig. 1a. Then, we constrain the weights for recurrent connections within the nodes of the RNN to be always zero for pairs nodes that lie at a Euclidean distance of larger than d . The training of the RNN then

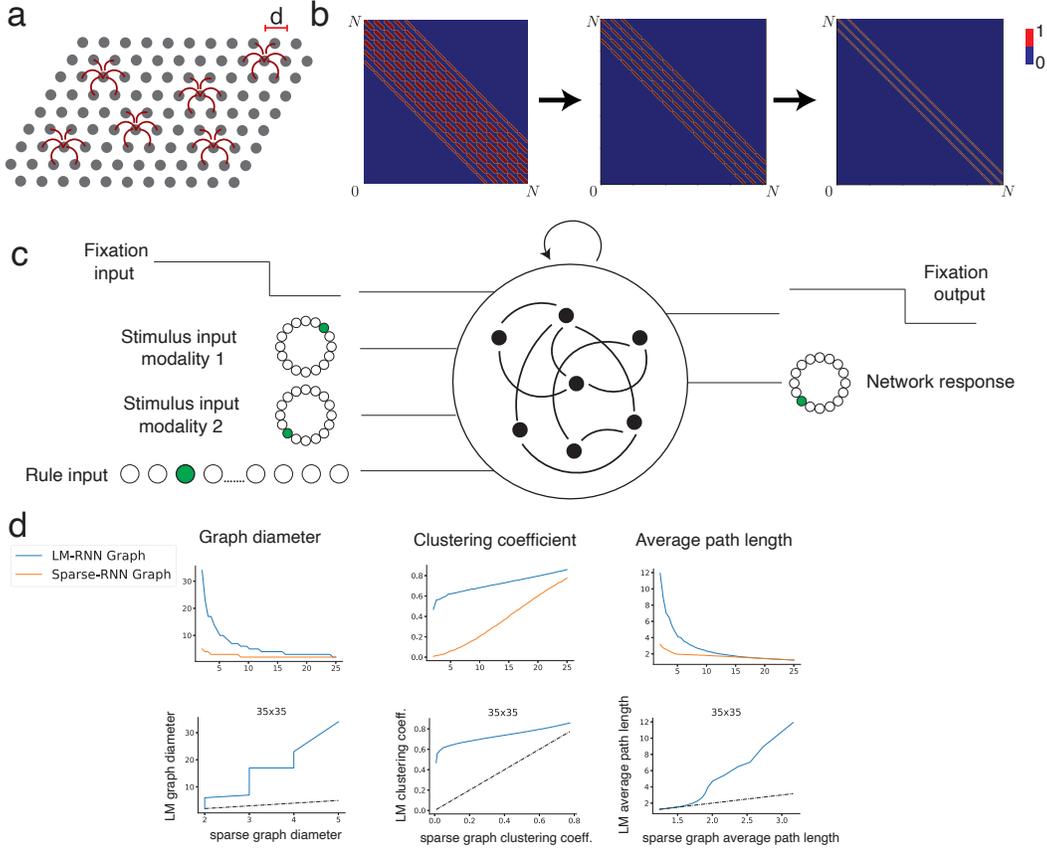


Figure 1: **Local connectivity constraint (locality masking) and schematic of cognitive multitasking RNN.**(a) In LM-RNNs neurons are arranged on a two-dimensional sheet, with nonzero weights permitted for nodes up to a distance $\leq d$ apart. (b) Locality masking on a two-dimensional sheet can be treated as a sparse mask on the hidden-to-hidden weights of an RNN (c) Schematic of the RNN setup in the context of 20-Cog-tasks and Mod-Cog: the network receives inputs encoding directions on two rings, a fixation signal and a rule input. The network is trained to output a fixation signal and a direction on a ring of output nodes. (d) A battery of graph theoretic metric to distinguish the connectivity graph of LM-RNNs from a random sparse (Erdos-Renyi) graph (left to right): Graph diameter, Clustering coefficient, Average path length.

proceeds in the usual fashion by using back propagation of the loss to update the unconstrained weights. We refer to such an RNN as a **Locality Masked RNN (LM-RNN)**.

This constraint on the weights of the RNN is implemented through a graph G , whose nodes are the units of the RNN, and edges correspond to pairs of units with unconstrained weights, determined by the spatial distance d . Operationally, the adjacency matrix of the graph, \mathbf{G} , is point-wise multiplied with the interaction matrix \mathbf{W} after each gradient step, effectively constraining which elements of the interaction matrix can be learnt by gradient descent:

$$\mathbf{W} \leftarrow \mathbf{G} \otimes \mathbf{W} \quad (1)$$

This graph adjacency matrix \mathbf{G} is static and unchanged throughout training.

In effect, the two-dimensional LM-RNN consists of a sparse subgraph G of the fully connected network, with each node connected to the $\sim \pi d^2$ nearest nodes to it. We posit that this sparse subgraph is like a “winning lottery ticket”, such that when trained in isolation the LM-RNN achieves comparable performance to a fully recurrently trained network. Moreover, we will demonstrate that these winning lottery tickets perform better in a more data-efficient manner than fully connected counterparts with a similar number of nodes or a similar number of synapses. This approach can be

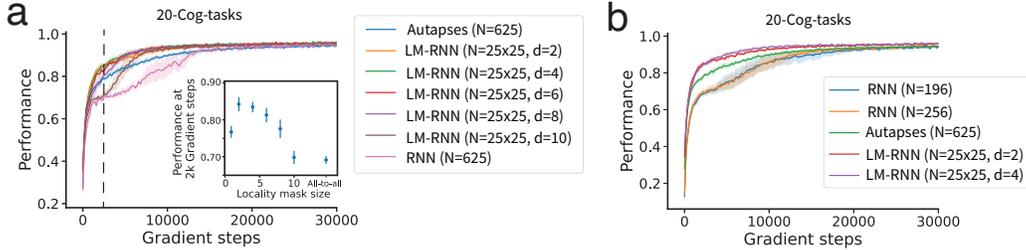


Figure 2: **LM-RNNs learn 20-Cog-tasks faster than fully connected networks with matched neuron or synapse number:** (a) LM-RNNs for all choices of locality mask sizes, d , perform better than fully-connected RNNs with the same number of nodes. This includes $d = 0$, i.e., a disconnected pool of autapses which also outperforms the fully-connected network. *Inset:* performance at 2×10^3 gradient steps across different values of d demonstrates a small optimal d that leads to the best performance. (b) LM-RNNs also outperform fully-connected RNNs constructed with the same number of synapses. RNN($N = 196$) and RNN($N = 256$) have 5.2×10^4 and 8.3×10^4 parameters respectively; $N = 625$ autapses, LM-RNN($N = 25 \times 25$, $d = 2$) and LM-RNN($N = 25 \times 25$, $d = 4$) have 4.4×10^4 , 5.2×10^4 and 7.4×10^4 parameters respectively

implemented easily in all training frameworks and is agnostic to the specific optimization algorithm being used.

Constructing a sparse graph in this particular fashion is distinct from a sparse random graph. To distinguish between our locality masks and sparse random graphs (Erdős-Rényi networks) with the same number of edges, we borrow several metrics from graph theory. In particular we examine three metrics — the average length of the shortest path between nodes, the graph diameter (i.e., the longest shortest path between two nodes) and correlation coefficient which (i.e., the density of mutually connected triplets of nodes)[Börner et al., 2007]. Notably, these metrics are properties fundamental to the structure of the connectivity in a graph, and are invariant to re-labelling and permutation of nodes in the graph. As can be seen in Fig.1d, there is a sharp contrast between the graph of an LM-RNN when compared with a random sparse RNN, with LM-RNNs having longer diameters and path lengths, but smaller clustering coefficients.

2.2 LM-RNNs learn the ‘20-Cog-tasks’ more rapidly than fully-connected RNNs

We apply LM-RNNs to a dataset used commonly in systems neuroscience, a set of 20 cognitive tasks introduced in [Yang et al., 2019], which we henceforth refer to as the 20-Cog-tasks. Each of these 20 tasks are constructed on the same input stimulus modalities — two rings of input units are used that each support a single activity bump, encoding a one-dimensional circular variable (which could represent direction of motion, for example), which we represent with vectors $\mathbf{u}_t^{\text{ring}1}(\theta_1)$ and $\mathbf{u}_t^{\text{ring}2}(\theta_2)$. Along with the two rings, the input into the RNN also comprises two additional inputs: first, a one-hot encoded rule input vector, indicating which task is to be performed, which we denote $\mathbf{u}^{\text{task-rule}}$ (“task name”); and second, a fixation input, a decrease of which is treated a ‘go’ signal for the RNN to provide the appropriate output, which we denote u_t^{fix} .

The expected output for each task is a response direction, which is again encoded in a ring of output units (which could represent, for example, a reach or saccade direction). These networks are trained using supervised learning with a cross entropy loss where the supervised target is a one-hot vector \mathbf{y} representing the required location of the bump on the ring, and the model’s output probabilities are constructed through a soft-max of the outputs, \mathbf{o}_t .

$$L_{CE}(\mathbf{o}_t, \mathbf{y}) = - \sum_i y_i \log \frac{e^{\mathbf{o}_{t,i}}}{\sum_j e^{\mathbf{o}_{t,j}}} \quad (2)$$

A schematic of the setup of the RNN is shown in Fig. 1c. For each trial of each task, the inputs are constructed as a gaussian bump on each of the two rings whose mean is drawn independently from a uniform distribution on the rings.

We compare LM-RNNs with different values of d against fully-connected RNNs with the same number of neurons (and hence many more parameters; cf. Fig. 2a) and fully-connected RNNs of a smaller size but with the same number of parameters (cf. Fig. 2b). In all cases, LM-RNNs for *any* value of d are far more sample-efficient and learn the tasks with same asymptotic performance as compared to the fully-connected counterparts. We also compare the performance of these models at the same fixed number of gradient steps early in training to show sample efficiency differences (Fig. 2a, *inset*). We hypothesize that this increased efficiency for LM-RNNs may be due to the ability to use a high-dimensional computational space, while having a significantly smaller number of parameters to be learned, resulting in the faster training observed at a similar level of performance.

2.3 ‘20-Cog-tasks’ are rapidly learned with simple autapse networks

While we demonstrated that LM-RNNs at all d perform better than fully-connected networks, we particularly note that $d = 0$, (i.e., a ‘network’ where each node is only connected to itself; in this case the network is simply a pool of disconnected autapses, and \mathbf{G} is an identity matrix) also performs better than a fully-connected network in terms of learning speed while reaching the *same asymptotic performance as a larger fully connected network*. Remarkably, as seen in Fig. 3, this pool of autapses continues to show ‘modularity’ in the network through the nodal-task-variance based metrics similar to the results of fully connected RNNs in [Yang et al., 2019] — however this apparent modularity clearly cannot be a result of any modular structures in the network due to the absence of any inter-node network connections. The autapse networks have no lateral connectivity and thus no way to share and reuse subtask structure across neurons. This suggests that such a task variance metric may simply be reflecting correlations between common inputs and similar input weights to hidden neurons. We verify this hypothesis in Fig. 3c where we plot the correlation between projection of the rule inputs to the hidden nodes, i.e., $\mathbf{b}^{\text{task-rule}}(\text{“task name”}) = \mathbf{W}^{\text{task-rule}} \mathbf{u}^{\text{task-rule}}(\text{“task name”})$ where $\mathbf{W}^{\text{task-rule}}$ is the submatrix of input weights formed by the rule-input appropriate columns of \mathbf{W}_{in} . Since the $\mathbf{u}^{\text{task-rule}}$ is presented as simply one-hot encoded vectors, the correlation of the rule inputs is simply $[\mathbf{W}^{\text{task-rule}}]^T \mathbf{W}^{\text{task-rule}}$. We observe that this correlation matrix of input projections in itself appears to cluster corresponding to the common subtask structure of 20-Cog-tasks.

We thus hope that our results motivate the study of better metrics and techniques to inspect functional modularity in RNNs, which we leave for future work.

While these results are in themselves indicative of the advantages conferred by LM-RNNs, we note that the 20-Cog-tasks are evidently too simplistic to make any strong claims, since they do not even require a network of connected neurons to accomplish the task. We present here first a simplified analysis demonstrating how a pool of disconnected autapses can solve 20-Cog-tasks, and thereafter present a more rigorous battery of cognitive tasks to more robustly demonstrate the utility of LM-RNNs.

2.4 Mechanistic analysis of the pool of autapses: A game of vector addition

For mechanistic interpretability, we first examine the dynamics of a single *linear* autapse in the presence of an input $b_t = \mathbf{W}_{in} \mathbf{u}_t$,

$$h_{t+1} = Wh_t + b_t. \tag{3}$$

This gives

$$h_t = W^t h_0 + \sum_{n=0}^{t-1} W^{t-n} b_n. \tag{4}$$

For a constant input, and for $0 \leq W < 1$ this can be simplified to give

$$h_t = \left[h_0 - \frac{b}{1-W} \right] W^t + \frac{b}{1-W}$$

Hence, we see that the autapse weight W defines an effective timescale for the autapse dynamics. Over the weight-dependent timescale $\tau = -\frac{1}{\log W}$ the autapse relaxes to a fixed point given by $h^* = \frac{b}{1-W}$. Thus, autapses with low weight ($W \approx 0$) very quickly converge to their associated bias dependent fixed point given by $h^* \approx b$ while autapses with high weight ($W = 1$) effectively function

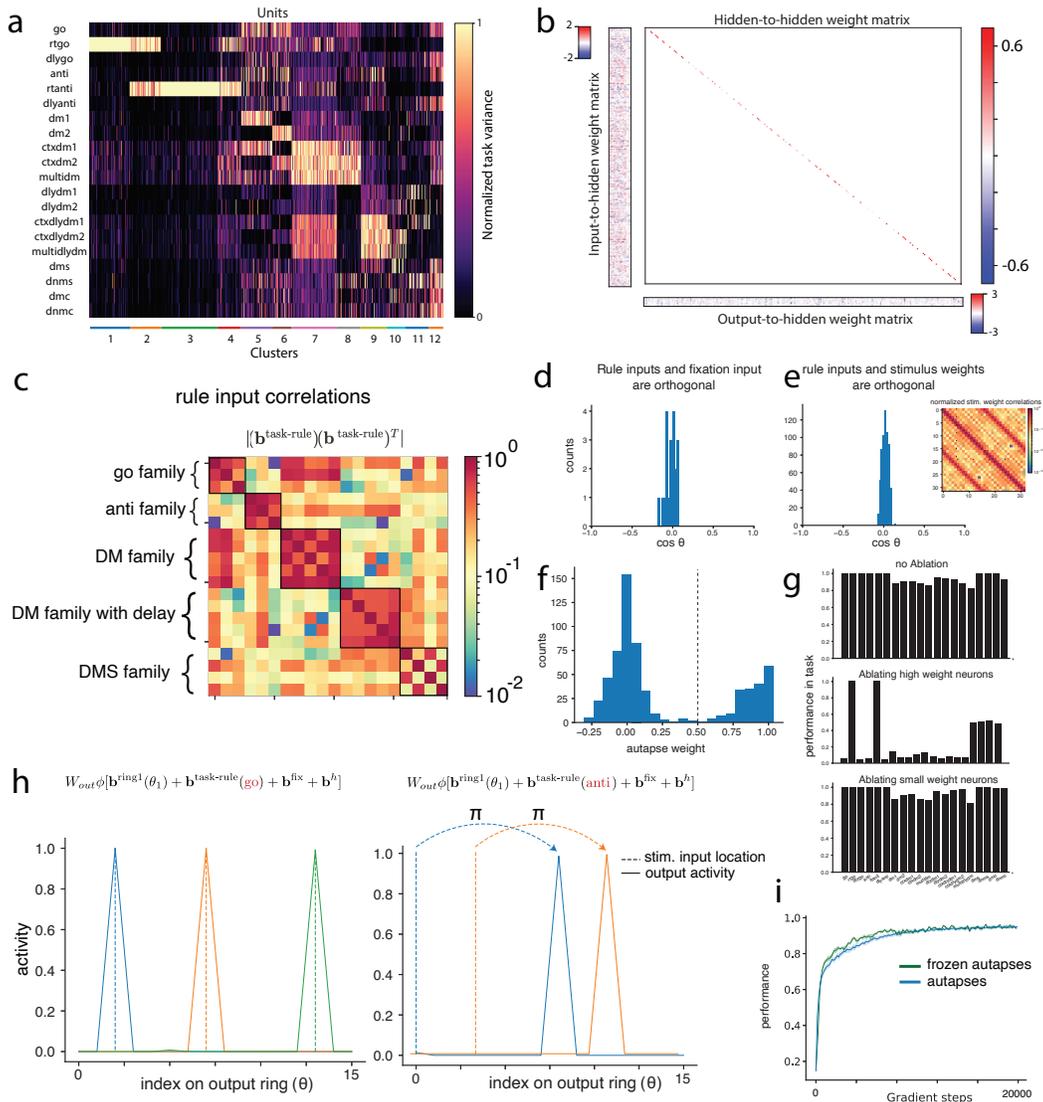


Figure 3: A suite of 20 cognitive tasks induces apparent modularity but is equally well-learned by a network of autapses (extreme spatial masking). (a) An autapse network from Fig.1 trained on the original 20-Cog-tasks showing the formation of 12 specialized clusters; (b) The weights of model including the diagonal hidden to hidden matrix which clearly show the existence of no modular structure. (c) Normalized rule-input weight correlations show block structure that is consistent with the 20-Cog-tasks family structure. (d) Rule input weights are orthogonal to fixation input weights (e) Rule input weights are orthogonal to stimulus input weights and (inset) stimulus input weight correlations show a circulant structure. (f) A histogram of autapse weights shows the existence of 2 clusters of autapses, those with autapse weight close to 1 and those with autapse weight close to 0. (g) Ablating the small weight autapses results in no loss of task performance while ablating the high weight autapses reduces all task performances to random chance, apart from the 2 reaction-time tasks *rt-go* and *rt-anti* which do not need memory. (h) Explicitly adding the appropriate input vectors is sufficient to solve the tasks without the need for dynamics: (left) for task *go* and (right) for task *anti*. (i) A pool of frozen autapses can learn the 20-Cog-tasks faster than a pool of autapses and reach the same asymptotic performance.

as perfect memory, retaining their initial state h_0 and perform addition of the network inputs as they are received:

$$h_t = h_0 + \sum_{n=0}^t b_n. \quad (5)$$

Examining the weights of autapses in the trained autapse pool reveals that the weights appear to be separable into two classes that correspond to the described above (Fig. 3f) — those with weights close to zero, that would be expected to rapidly approach a fixed point; and those with weights close to one, that would be expected to perform vector addition. We also observe that the rule inputs are approximately orthogonal to the fixation and ring inputs, Fig.3d,e, suggesting that the dynamics due to the rule inputs act in a subspace that is independent from the fixation and input stimuli.

We hypothesized that performing vector addition would be sufficient to be able to solve the 20-Cog-tasks. We verified this in two ways: firstly, we found that ablating autapses with lower weights resulted in no loss in performance, whereas ablating the larger weight synapses lead to significant performance deficits reducing the performance on almost all tasks to chance ¹, Fig. 3g; secondly, we found that explicit addition of the $\mathbf{b}^{\text{task-rule}}$ (“task rule”) vectors to the inputs from the other modalities was sufficient to generate outputs that performed the task, Fig. 3h.

2.4.1 A frozen pool of autapses

Motivated by these ablation studies, we trained an autapse pool with untrainable fixed autapse weights set to 1. We refer to this pool as the “frozen autapse pool”. This setup also learns all 20 tasks notably faster than a regular autapse pool, Fig. 3i. In this case the autapses are effectively “copy gates”, storing perfect memory of their previous input and adding them to the input vector currently being received (cf. equation 5).

2.5 Mod-Cog: An expanded battery of cognitive tasks

The reason the 20-Cog-tasks were trivially solved by a pool of (frozen) autapses is that the tasks were static, involving memory and fixed points but no dynamical computation like integration and sequence generation. These computations involve the manipulation of the information held in working memory. A circuit with integration properties is able in principle to powerfully generalize across tasks in a way that networks that exhibit the same set of states only as stable fixed points are not able to [Klukas et al., 2020, Khona and Fiete, 2022]. To perform a more robust demonstration of the utility of LM-RNNs, we construct a battery of new tasks based on the principle of integration, and related to cognitive science problems such as interval estimation, physical and mental navigation, and sequence generation. We build a set of modular and compositionally constructed tasks, using the neurogym framework [Molano-Mazon et al., 2022], which in turn was built on the AI opengym environment.

In particular, we build extensions, that incorporate additional complexity in two main forms: integration and sequence generation.

In the case of integration, we consider the set of ‘delay’ based tasks in the 20-Cog-tasks. In the 20-Cog-tasks, 12 out of 20 tasks involve a delay period in the presented input, wherein the network is expected to persistently hold the input presented in an internal memory before performing a task-relevant computation. To incorporate interval estimation in these tasks, we require the output to be displaced with respect to the originally expected output by a magnitude dependent on the length of the delay period. To this end, we choose the delay length randomly from a uniform distribution (as opposed to the fixed delay length considered in 20-Cog-tasks). As representative examples, in Fig. 4a,b we show the inputs and expected outputs for two different delay period lengths in the DlyGo_IntL task, constructed as an interval estimation extension to the DlyGo task from 20-Cog-tasks. For each of the 12 tasks, the interval-dependent-displacement may be either of clockwise or anti-clockwise, resulting in the introduction of 24 new tasks.

¹The chance performance for the delay match to sample family of tasks dms , dnms , dmc , dnmc is 50% due to the way samples are drawn: half of the trials are matching and the other half are non-matching, refer to methods of [Yang et al., 2019]. The performance of reaction-time tasks rt-go and rt-anti is not affected since they do not require memory.

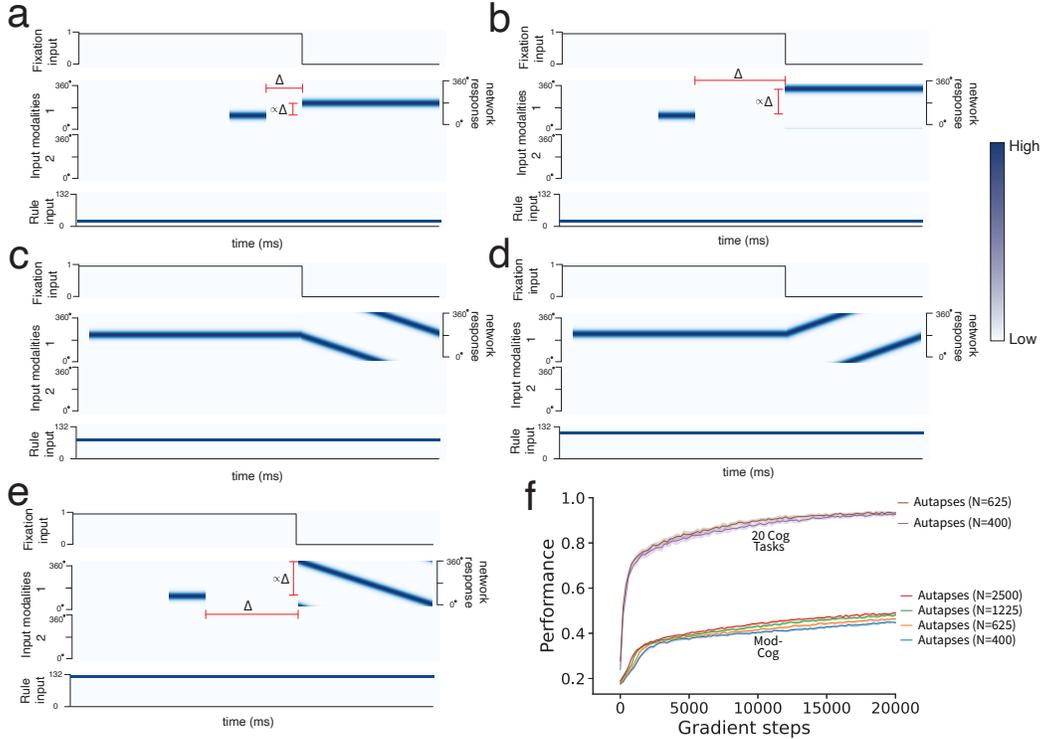


Figure 4: Mod-Cog: **A more-complex suite of upto 132 cognitive tasks.** Schematic of the tasks in Mod-Cog, that extend the 20-Cog-tasks. (a,b) Tasks DlyGo_IntL and DlyGo_IntL as interval estimation extensions to the DlyGo task from 20-Cog-tasks for two different lengths of delay periods. (c,d) Tasks Go_SeqR and Go_SeqL as sequence generation extensions to the Go task from 20-Cog-tasks. (e) The compositional extension DlyGo_IntL_SeqR that combines both interval estimation as well as sequence generation. (f) Mod-Cog introduces significantly more complex tasks as compared to 20-Cog-tasks, which can no longer be solved by a pool of autapses, regardless of their number.

In the case of sequence generation, the output of each task is modified to not be a single static direction, but instead a time-varying output corresponding to a drifting direction starting at a particular point (dependent on the particular task). The direction of drift can be changed dependent on the particular task. As representative examples, in Fig. 4c,d we show the inputs and expected outputs for the Go_SeqL and Go_SeqR tasks, constructed as an sequence generation extensions to the Go task from 20-Cog-tasks. This introduces 40 new tasks based on the earlier set of 20 tasks, with the output of each task drifting either clockwise or anti-clockwise.

This completes the construction of the 64 new tasks that we use in conjunction with the original 20 tasks as the tasks used for our main set of results hereafter in this paper. We refer to this set of 84 tasks as Mod-Cog. We note however that our modifications to the tasks are modular in nature (which is similar in spirit to the already existing modular subtask structure in the 20-Cog-tasks). This allows for an additional extension of 48 more tasks that may be generated by a composition of the sequence generation and interval estimation extensions (such as the DlyGo_IntL_SeqR task shown in Fig. 4e). For simplicity, we do not use these additional 48 tasks in our main results; however they are included in the repository of tasks that we provide at github link (will be inserted upon acceptance; provided as .zip file in supplementary material).

The rule input used for Mod-Cog is encoded as a one-hot vector, similar to the setup used in [Yang et al., 2019]. This ensures that, by construction, the rule input cannot be directly used as a signal to help decompose tasks into having common subtasks.

To demonstrate that Mod-Cog is significantly harder than 20-Cog-tasks, we demonstrate in Fig. 4f that a pool of autapses is incapable of achieving significant performance levels, in sharp contrast

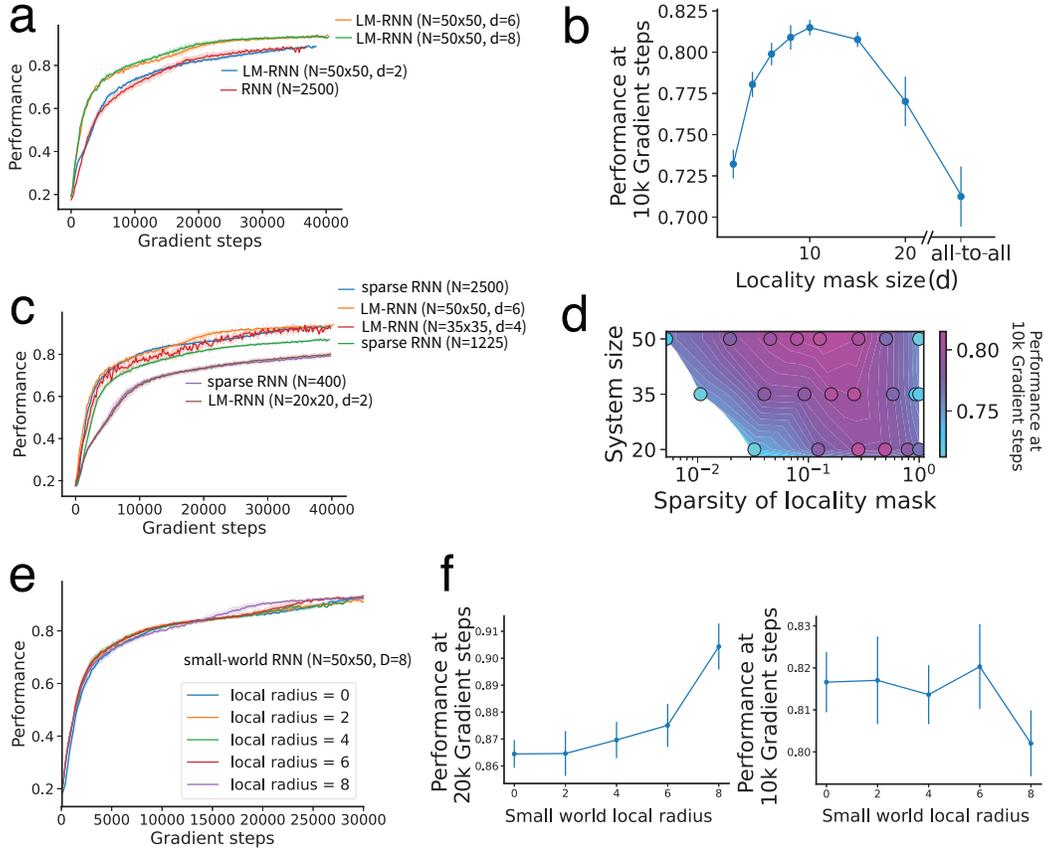


Figure 5: **Faster learning and better asymptotic performance of LM-RNNs and sparse RNNs on Mod-Cog** (a) LM-RNNs for all values of d perform equally well or better than a fully-connected network with the same number of nodes. This improvement takes the form of faster learning as well as better asymptotic performance. (b) Network performance at 10^4 gradient steps for different values of d and networks of size 50×50 , demonstrating an optimal locality size that leads to the fastest learning. (c) LM-RNN performance for varying system size for a fixed network sparsity of $\sim 4\%$, as compared with fixed random sparse networks with the same sparsity. Performance for sparse RNNs is similar to or only slightly worse than LM-RNNs with the same sparsity (d) Network performance as a function of system size ($N^2 = \text{System Size}$) and sparsity of the locality mask (i.e., the fraction of nonzero entries in the hidden-to-hidden weight matrix). At larger system sizes, the optimal sparsity for best performance is lower. (d) (e)

with the 20-Cog-tasks. Moreover, this is independent of the number of autapses — even pools that are ~ 6 times larger than those necessary for solving 20-Cog-tasks are unable to produce larger than 50% performance accuracy.

2.6 RNN and LM-RNN performance on Mod-Cog

As earlier, we examine the performance of LM-RNNs with varying d on Mod-Cog, and compare them with fully-connected RNNs as a function of the number of gradient steps in training (cf. Fig. 5a). Here again we see that LM-RNNs, for appropriately chosen values of d train significantly more rapidly as compared to fully connected networks. Due to the increased task complexity (as evidenced by Fig. 4f), locality masks corresponding to very small values of d result in suboptimal performance (which is nonetheless similar to the $d \rightarrow \infty$ corresponding to the fully connected network). Instead, intermediate small values d outperform all other values of d , as shown in Fig. 5b, indicating an optimal nontrivial locality mask d . For smaller networks, the optimal value of d corresponds to increasingly larger fractions of all edges in the network (cf. Fig. 5d), indicating

	Performance at 10k gradient steps	No. of nonzero weights	No. of nonzero hidden-to-hidden weights
LM-RNN($N = 50 \times 50, d = 6$)	0.799 ± 0.007	6.17×10^5	2.82×10^5
LM-RNN($N = 50 \times 50, d = 10$)	0.815 ± 0.005	1.13×10^6	7.92×10^5
RNN($N = 676$)	0.69 ± 0.01	5.47×10^5	4.57×10^5
RNN($N = 2500$)	0.71 ± 0.02	6.58×10^6	6.25×10^6
L_1 RNN($N = 2500$), $\lambda = 10^{-4}$	0.59 ± 0.02	4.00×10^5	6.51×10^4
L_1 RNN($N = 2500$), $\lambda = 10^{-5}$	0.72 ± 0.01	5.16×10^5	1.81×10^5
L_1 RNN($N = 2500$), $\lambda = 10^{-6}$	0.73 ± 0.02	1.02×10^6	6.87×10^5
L_1 RNN($N = 2500$), $\lambda = 10^{-7}$	0.71 ± 0.01	1.86×10^6	1.52×10^6

Table 1: **Comparison across networks with similar number of trainable parameters and nonzero weights.** L_1 RNNs are fully-connected RNNs with an L_1 regularization to promote sparse solutions. In this case, the number of nonzero weights is counted as the number of weights above a small threshold in the trained RNN. See Appendix Fig.6 for learning curves

that the optimal may depend more directly on the complexity of the tasks to be solved, rather than scaling with the number of nodes in the network.

For most LM-RNNs, a *random* sparse graph with the same number of incoming edges from each node performs almost as well as the graph chosen through locality masking, as shown in Fig. 5c. Nevertheless, in each case the locality mask performs as well as or slightly better than the random sparse graph, while performing significantly better than dense fully-connected networks. For the case of small values of d , we hypothesize that the slight improvement of locality masks may arise from the presence of disconnected components in random sparse graphs in networks with low degrees. For clustering visualizations based on nodal task variance [Yang et al., 2019], see Appendix Figs. 7,9,8,10 and 11.

We also examined small-world networks[Watts and Strogatz, 1998] as graphs that interpolate between random sparse graphs and the locality masks of LM-RNNs. In particular, we considered graphs with a fixed number of total synapses corresponding to an LM-RNN with $d = 8$. Small-world graphs were constructed with this number of synapses: connections were formed between nodes within a ‘local radius’ distance of each other, and random sparse connections were added until the target total number of synapses. A local radius of zero then corresponds to a completely random sparse graph, and a local radius of 8 corresponds to the LM-RNN that we have been considering thus far. We observed that there was no significant variation of performance across different local radii at 10^3 gradient steps, however, at 2×10^3 gradient steps the LM-RNN appears to perform slightly better than any small world network, Fig. 5f (left).

2.7 Controls – comparisons with various baselines for sparsity

We perform comparisons across 3 baselines. First, fully-connected networks with the same number of nodes; second, fully-connected networks with the same number of synapses; and third, fully-connected networks with the same number of nodes, but with an additional L_1 regularization term in the loss function to promote the discovery of sparse solutions through training. As we demonstrate in Table 1, LM-RNNs reach a higher performance earlier than all other comparable models. While sparse networks achieved through L_1 regularized models perform better than fully-connected models at some choices of regularization strengths, they achieve this better performance slower than LM-RNNs with static sparsity. Thus while sparsity is clearly beneficial to improved performance in the multitask setting of Mod-Cog, choosing this sparsity in a fixed, task-independent fashion at the start of training results in faster training.

3 Discussion

Through our results, we have demonstrated that a simple fixed sparse graph can provide a large increase in the sample-efficiency of the training process in single-layer recurrent networks. These results could in principle be extended to other architectures and tasks but we restricted our experiments to simple RNNs and cognitive tasks for their relevance to systems and computational neuroscience.

As RNNs trained on multitask learning problems become more popular as models of PFC and other associated brain regions [Driscoll et al., 2022, Hummos, 2022, Riveland and Pouget, 2022, da Costa et al., 2019, Duncker et al., 2020, Masse et al., 2022, Kao et al., 2021], we need to be more cognizant about how the computational complexity of the task set used affects learnt representations. By demonstrating that a pool of unconnected autapses can solve the 20-Cog-tasks [Yang et al., 2019], we have shown that this apparently complex set of tasks can be solved without any communication between neurons. This raises the hypothesis that appropriately tuned input weights corresponding to the “rule inputs” along with static memory are enough to solve many tasks. Here we constructed Mod-Cog as a set of tasks with increased complexity by adding more modular sub-components that require the RNN to perform non-trivial computations. Thus we provide a more reasonable multitask setting that leads to richer solutions and may be a better testbed for exploring shared motifs and representations across modular subtasks. Another possible direction for future study to increase task complexity for 20-Cog-tasks has been to eliminate rule inputs and force the network to infer the task needed to be solved. This would likely involve some flavor of predictive coding. A recent work Hummos [2022] makes progress in this direction.

We have found that given a certain amount of task complexity and network size, there is an optimal amount of locality masking that provides the most benefit to learning, both in terms of sample efficiency and asymptotic performance. While we have shown this result for recurrent networks, qualitatively similar results on an optimal value of sparsity have been derived analytically for cerebellum-like feedforward architectures [Litwin-Kumar et al., 2017]. It remains an open question to theoretically investigate the relationship between task complexity and the amount of sparsity that is needed: if the network is too sparse, it will not have enough expressivity to perform well on the dataset, while if the network is too dense, the benefits provided by sparsity will not be exploited.

Although our matrices are extremely sparse and would be well suited to sparse matrix representations, we still maintain dense matrix datatypes for all of the training and evaluation processes since standard libraries like PyTorch do not have native support for these data structures. As sparse matrices get more common and their usefulness more apparent, as has been pointed out before [Liu et al., 2021], it will be very useful for native deep learning software and hardware implementations on GPUs to exploit the potential efficiencies of very sparse matrix structures.

4 Acknowledgements

We thank Guangyu Robert Yang for open-sourcing the 20 cognitive tasks set and the NeuroGym framework [Molano-Mazon et al., 2022] released under an MIT License. This work was supported by the MathWorks Science Fellowship to MK, the Simons Foundation through the Simons Collaboration on the Global Brain, the ONR, and the Howard Hughes Medical Institute through the Faculty Scholars Program to IRF.

References

- Katy Börner, Soma Sanyal, Alessandro Vespignani, et al. Network science. *Annu. rev. inf. sci. technol.*, 41(1):537–607, 2007.
- Pedro F. da Costa, Sebastian Popescu, Robert Leech, and Romy Lorenz. Elucidating cognitive processes using lstms. In *Conference on Cognitive Computational Neuroscience*, 2019. URL <https://ccneuro.org/2019/proceedings/0000272.pdf>.
- Pau de Jorge, Amartya Sanyal, Harkirat S Behl, Philip HS Torr, Gregory Rogez, and Puneet K Dokania. Progressive skeletonization: Trimming more fat from a network at initialization. *arXiv preprint arXiv:2006.09081*, 2020.
- Laura Driscoll, Krishna Shenoy, and David Sussillo. Flexible multitask computation in recurrent networks utilizes shared dynamical motifs. *bioRxiv*, 2022.
- Lea Duncker, Laura Driscoll, Krishna V Shenoy, Maneesh Sahani, and David Sussillo. Organizing recurrent network dynamics by task-computation to enable continual learning. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 14387–14397. Curran Associates, Inc., 2020. URL <https://proceedings.neurips.cc/paper/2020/file/a576eafbce762079f7d1f77fca1c5cc2-Paper.pdf>.
- Mária Ercsey-Ravasz, Nikola T Markov, Camille Lamy, David C Van Essen, Kenneth Knoblauch, Zoltán Toroczkai, and Henry Kennedy. A predictive network model of cerebral cortical connectivity based on a distance rule. *Neuron*, 80(1):184–197, 2013.
- Utku Evci, Trevor Gale, Jacob Menick, Pablo Samuel Castro, and Erich Elsen. Rigging the lottery: Making all tickets winners. In *International Conference on Machine Learning*, pages 2943–2952. PMLR, 2020.
- Dawn Finzi, Eshed Margalit, Kendrick Kay, Daniel LK Yamins, and Kalanit Grill-Spector. Topographic DCNNs trained on a single self-supervised task capture the functional organization of cortex into visual processing streams. In *SVRHM 2022 Workshop @ NeurIPS*, 2022. URL <https://openreview.net/forum?id=E1iY-d13smd>.
- Timo Flesch, David G Nagy, Andrew Saxe, and Christopher Summerfield. Modelling continual learning in humans with hebbian context gating and exponentially decaying task signals. *arXiv preprint arXiv:2203.11560*, 2022.
- Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural networks. *arXiv preprint arXiv:1803.03635*, 2018.
- Jonathan Frankle, Gintare Karolina Dziugaite, Daniel M Roy, and Michael Carbin. Stabilizing the lottery ticket hypothesis. *arXiv preprint arXiv:1903.01611*, 2019.
- Song Han, Jeff Pool, John Tran, and William Dally. Learning both weights and connections for efficient neural network. *Advances in neural information processing systems*, 28, 2015.
- John J Hopfield. Neural networks and physical systems with emergent collective computational abilities. *Proceedings of the national academy of sciences*, 79(8):2554–2558, 1982.
- Ali Hummos. Thalamus: a brain-inspired algorithm for biologically-plausible continual learning and disentangled representations. *arXiv preprint arXiv:2205.11713*, 2022.
- Herbert Jaeger. The “echo state” approach to analysing and training recurrent neural networks-with an erratum note. *Bonn, Germany: German National Research Center for Information Technology GMD Technical Report*, 148(34):13, 2001.
- Ta-Chu Kao, Kristopher Jensen, Gido van de Ven, Alberto Bernacchia, and Guillaume Hennequin. Natural continual learning: success is a journey, not (just) a destination. *Advances in Neural Information Processing Systems*, 34, 2021.

- Mikhail Khona and Ila R Fiete. Attractor and integrator networks in the brain. *Nature Reviews Neuroscience*, pages 1–23, 2022.
- Mikhail Khona, Sarthak Chandra, and Ila R Fiete. From smooth cortical gradients to discrete modules: spontaneous and topologically robust emergence of modularity in grid cells. *bioRxiv*, pages 2021–10, 2022.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Mirko Klukas, Marcus Lewis, and Ila Fiete. Efficient and flexible representation of higher-dimensional cognitive variables with grid cells. *PLoS computational biology*, 16(4):e1007796, 2020.
- Aditya Kusupati, Vivek Ramanujan, Raghav Somani, Mitchell Wortsman, Prateek Jain, Sham Kakade, and Ali Farhadi. Soft threshold weight reparameterization for learnable sparsity. In *International Conference on Machine Learning*, pages 5544–5555. PMLR, 2020.
- Hyodong Lee, Eshed Margalit, Kamila M Jozwik, Michael A Cohen, Nancy Kanwisher, Daniel LK Yamins, and James J DiCarlo. Topographic deep artificial neural networks reproduce the hallmarks of the primate inferior temporal cortex face processing network. *bioRxiv*, 2020a.
- Namhoon Lee, Thalaiyasingam Ajanthan, and Philip Torr. SNIP: SINGLE-SHOT NETWORK PRUNING BASED ON CONNECTION SENSITIVITY. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=B1VZqjAcYX>.
- Namhoon Lee, Thalaiyasingam Ajanthan, Stephen Gould, and Philip H. S. Torr. A signal propagation perspective for pruning neural networks at initialization. In *International Conference on Learning Representations*, 2020b. URL <https://openreview.net/forum?id=HJeTo2VFwH>.
- Ashok Litwin-Kumar, Kameron Decker Harris, Richard Axel, Haim Sompolinsky, and LF Abbott. Optimal degrees of synaptic connectivity. *Neuron*, 93(5):1153–1164, 2017.
- Shiwei Liu, Decebal Constantin Mocanu, Yulong Pei, and Mykola Pechenizkiy. Selfish sparse rnn training. In *International Conference on Machine Learning*, pages 6893–6904. PMLR, 2021.
- Christos Louizos, Max Welling, and Diederik P Kingma. Learning sparse neural networks through l_0 regularization. *arXiv preprint arXiv:1712.01312*, 2017.
- Mantas Lukoševičius and Herbert Jaeger. Reservoir computing approaches to recurrent neural network training. *Computer Science Review*, 3(3):127–149, 2009.
- Wolfgang Maass. Liquid state machines: motivation, theory, and applications. *Computability in context: computation and logic in the real world*, pages 275–296, 2011.
- Wolfgang Maass, Thomas Natschläger, and Henry Markram. Real-time computing without stable states: A new framework for neural computation based on perturbations. *Neural computation*, 14(11):2531–2560, 2002.
- Nikola T Markov, Mária Ercsey-Ravasz, David C Van Essen, Kenneth Knoblauch, Zoltán Toroczkai, and Henry Kennedy. Cortical high-density counterstream architectures. *Science*, 342(6158):1238406, 2013.
- Christian David Marton, Guillaume Lajoie, and Kanaka Rajan. Efficient and robust multi-task learning in the brain with modular task primitives. *arXiv preprint arXiv:2105.14108*, 2021.
- Nicolas Y Masse, Matthew C Rosen, Doris Y Tsao, and David J Freedman. Rapid learning with highly localized synaptic plasticity. *bioRxiv*, 2022.

- Decebal Constantin Mocanu, Elena Mocanu, Phuong H Nguyen, Madeleine Gibescu, and Antonio Liotta. A topological insight into restricted boltzmann machines. *Machine Learning*, 104(2): 243–270, 2016.
- Manuel Molano-Mazon, Joao Barbosa, Jordi Pastor-Ciurana, Marta Fradera, Ru-Yuan Zhang, Jeremy Forest, Jorge del Pozo Lerida, Li Ji-An, Christopher J Cueva, Jaime de la Rocha, et al. Neurogym: An open resource for developing and sharing neuroscience tasks. *PsyArXiv*, 2022.
- Eli Moore and Rishidev Chaudhuri. Using noise to probe recurrent neural network structure and prune synapses. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 14046–14057. Curran Associates, Inc., 2020. URL <https://proceedings.neurips.cc/paper/2020/file/a1ada9947e0d683b4625f94c74104d73-Paper.pdf>.
- Sharan Narang, Erich Elsen, Gregory Diamos, and Shubho Sengupta. Exploring sparsity in recurrent neural networks. *arXiv preprint arXiv:1704.05119*, 2017.
- Dina Obeid and Talia Konkle. Wiring minimization of deep neural networks reveal conditions in which multiple visuotopic areas emerge. *Journal of Vision*, 21(9):2135–2135, 2021.
- Reidar Riveland and Alexandre Pouget. A neural model of task compositionality with natural language instructions. *bioRxiv*, 2022.
- Peter J Rousseeuw. Silhouettes: a graphical aid to the interpretation and validation of cluster analysis. *Journal of computational and applied mathematics*, 20:53–65, 1987.
- Pedro Savarese, Hugo Silva, and Michael Maire. Winning the lottery with continuous sparsification. *Advances in Neural Information Processing Systems*, 33:11380–11390, 2020.
- Hidenori Tanaka, Daniel Kunin, Daniel L Yamins, and Surya Ganguli. Pruning neural networks without any data by iteratively conserving synaptic flow. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 6377–6389. Curran Associates, Inc., 2020. URL <https://proceedings.neurips.cc/paper/2020/file/46a4378f835dc8040c8057beb6a2da52-Paper.pdf>.
- Panagiota Theodoni, Piotr Majka, David H Reser, Daniel K Wójcik, Marcello GP Rosa, and Xiao-Jing Wang. Structural attributes and principles of the neocortical connectome in the marmoset monkey. *Cerebral Cortex*, 32(1):15–28, 2022.
- Chaoqi Wang, Guodong Zhang, and Roger Grosse. Picking winning tickets before training by preserving gradient flow. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=SkgSACVKPH>.
- Duncan J Watts and Steven H Strogatz. Collective dynamics of ‘small-world’ networks. *nature*, 393(6684):440–442, 1998.
- Wei Wen, Yuxiong He, Samyam Rajbhandari, Minjia Zhang, Wenhan Wang, Fang Liu, Bin Hu, Yiran Chen, and Hai Li. Learning intrinsic sparse structures within long short-term memory. In *International Conference on Learning Representations*, 2018. URL <https://openreview.net/forum?id=rk6cfpRjZ>.
- Guangyu Robert Yang, Madhura R Joglekar, H Francis Song, William T Newsome, and Xiao-Jing Wang. Task representations in neural networks trained to perform many cognitive tasks. *Nature neuroscience*, 22(2):297–306, 2019.
- Hattie Zhou, Janice Lan, Rosanne Liu, and Jason Yosinski. Deconstructing lottery tickets: Zeros, signs, and the supermask. *Advances in neural information processing systems*, 32, 2019.
- Julian Georg Zilly, Rupesh Kumar Srivastava, Jan Koutník, and Jürgen Schmidhuber. Recurrent highway networks. In *International conference on machine learning*, pages 4189–4198. PMLR, 2017.

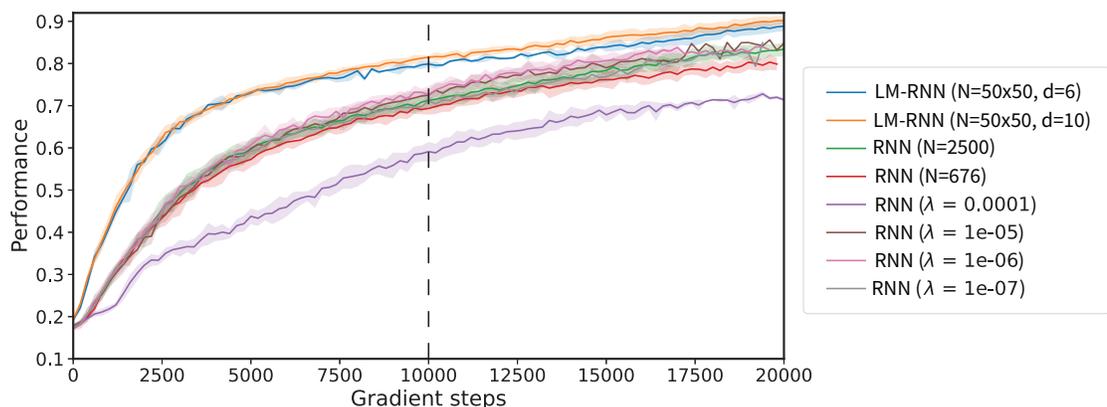


Figure 6: **Performance curves** of models trained with varying values of sparsity (L1) regularization showing that a fixed local mask (of similar sparsity, cf. Table 1) is better.

A Methods and Hyperparameters

PyTorch was used for all simulations. All networks were trained with supervised learning using a cross entropy loss. The optimizer used was Adam [Kingma and Ba, 2014] with a learning rate of 10^{-3} . Each trial was drawn randomly and independently and the tasks were randomly interleaved. The Neurogym [Molano-Mazon et al., 2022] environment was used to create tasks and the training data for the model.

For every performance curve, we trained 15 RNNs with different random seeds and used the averaged curve for plotting and computing the optimal locality mask sizes.

An expanded description of the Mod-Cog tasks and how to create them will be made available at the Github repository after acceptance.

To measure the number of clusters that the hidden nodes of the RNN partition into to solve a task, we use first use agglomerative hierarchical clustering² to obtain a linkage tree for the variance of each node across different inputs for a given task. Then, the silhouette score [Rousseeuw, 1987] is evaluated for each possible linkage-based cluster partition and the partition with the highest score is selected to represent the clustering of the RNN.

Graph theoretic measures shown in Fig. 1d were computed using the Python package ‘networkx’ version 2.8.4.

²<https://docs.scipy.org/doc/scipy/reference/generated/scipy.cluster.hierarchy.linkage.html>

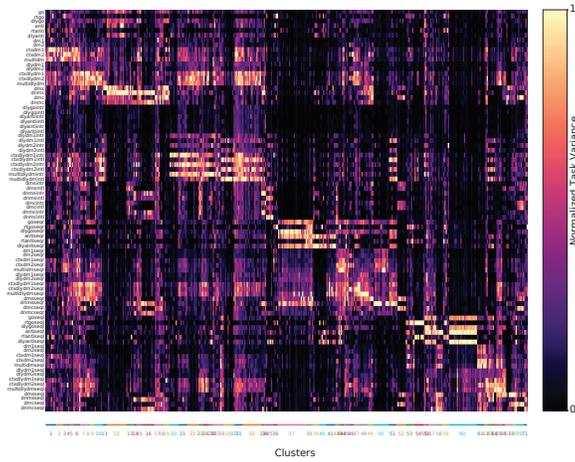


Figure 7: An LM-RNN(N=20x20,d=2) trained on 84 Mod-Cog tasks showing the formation of clusters

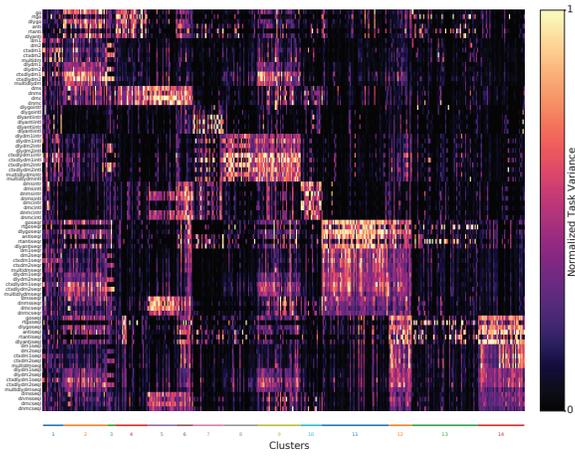


Figure 8: An LM-RNN(N=20x20,d=6) trained on 84 Mod-Cog tasks showing the formation of clusters

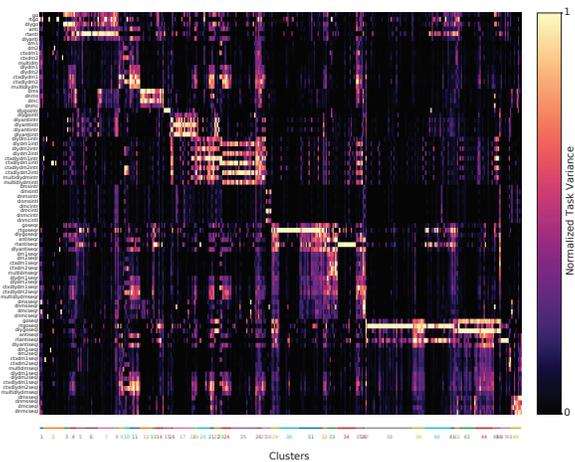


Figure 9: An RNN(N=400) trained on 84 Mod-Cog tasks showing the formation of clusters

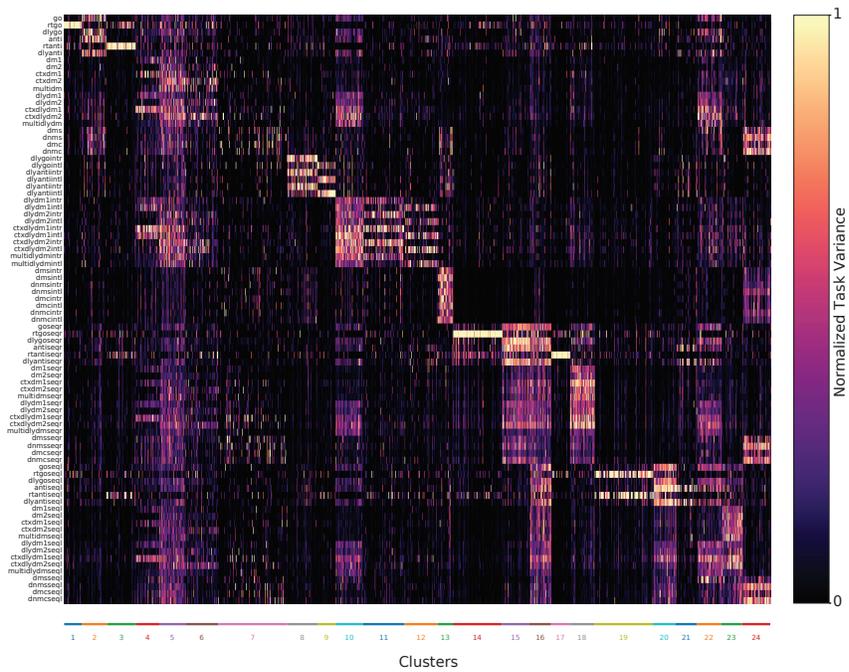


Figure 10: 24 Clusters formed in an LM-RNN($N = 50 \times 50$, $d = 10$) trained on 84 Mod-Cog tasks.

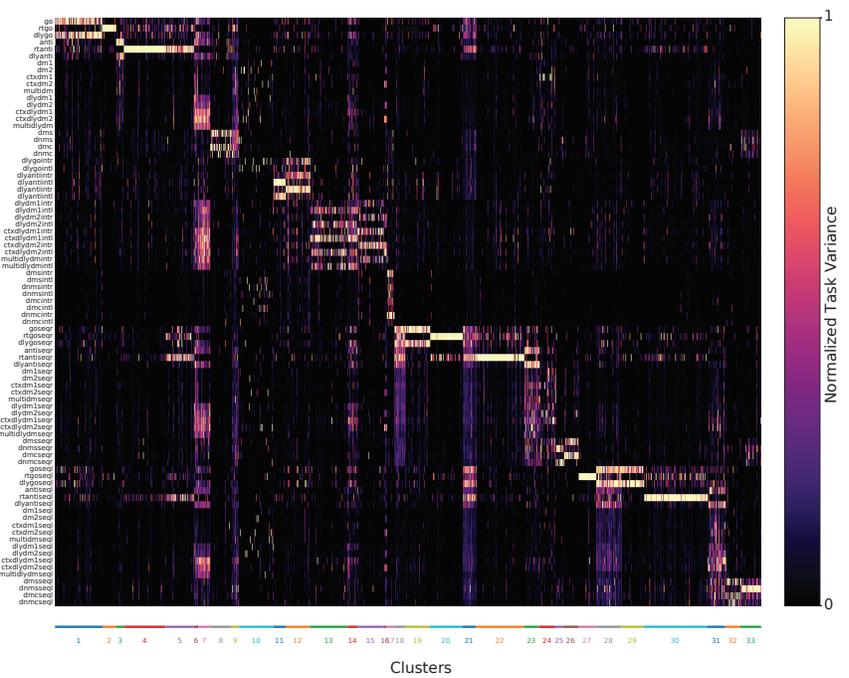


Figure 11: 33 Clusters formed in an RNN($N=2500$) trained on 84 Mod-Cog tasks, contrast with previous figure where the number of clusters is smaller.