

1 Top-down signaling dynamically mediates information processing 2 in biologically inspired RNNs

3 Tomas Gallo Aquino¹ *, Robert Kim ² *, Nuttida Rungratsameetaweemana¹

4 ¹ Department of Biomedical Engineering, Columbia University, New York, NY 10027, USA

5 ² Department of Neurology, Cedars-Sinai Medical Center, Los Angeles, CA 90048, USA

6 * Equal contribution

7 Correspondence: nr2869@columbia.edu (N.R.)

8 Abstract

9 Recent studies have proposed employing biologically plausible recurrent neural networks (RNNs)
10 to explore flexible decision making processes in the brain. However, the mechanisms underlying
11 the integration of bottom-up factors (such as incoming sensory signals) and top-down factors (such
12 as task instructions and selective attention) remain poorly understood, both within the context of
13 these models and the brain. To address this question, we trained biologically inspired RNNs on
14 complex cognitive tasks that require adaptive integration of these factors. By performing extensive
15 dynamical systems analyses, we show that our RNN model is capable of seamlessly incorporat-
16 ing top-down signals with sensory signals to perform the complex tasks. Furthermore, through
17 comprehensive local connectivity analyses, we identified important inhibitory feedback signals that
18 efficiently modulate the bottom-up sensory coding in a task-driven manner. Finally, we introduced
19 an anatomical constraint where a specific subgroup of neurons receives the sensory input signal,
20 effectively creating a designated sensory area within the RNN. Through this constraint, we show
21 that these “sensory” neurons possess the remarkable ability to multiplex and dynamically combine
22 both bottom-up and top-down information. These findings are consistent with recent experimen-
23 tal results highlighting that such integration is a key factor in facilitating flexible decision making.
24 Overall, our work provides a framework for generating testable hypotheses for the hierarchical
25 encoding of task-relevant information.

26 **Introduction**

27 Models based on recurrent neural networks (RNNs) of continuous-variable firing rate units have
28 been widely used to reproduce previously observed experimental findings and to explore neural dy-
29 namics associated with cognitive functions such as attention and working memory [1–4]. However,
30 these modeling studies primarily focus on tasks with fixed task and context, thereby forcing the
31 model to predominantly rely on bottom-up processing (i.e., encoding sensory stimuli).

32 Such an approach does not fully capture the dynamics of the natural environment, where in-
33 formative top-down information is consistently available and plays a critical role in optimizing the
34 translation of incoming sensory signals into appropriate motor responses. In the real world, there is
35 an abundance of contextual cues, task instructions, and prior knowledge that guide and shape our
36 perception and decision-making processes. Neglecting to take these important top-down factors
37 into models limits our understanding of how neural systems combine bottom-up and top-down
38 information to generate optimal behavior.

39 To address this knowledge gap, we carefully designed cognitive tasks that require seamless
40 integration of bottom-up sensory inputs with top-down factors. By training biologically inspired
41 RNNs on these tasks and studying the trained networks, we show that disinhibitory circuitry
42 governed by inhibitory-to-inhibitory connections are critical for the intricate interplay between the
43 two essential components of cognitive processing. More specifically, our contributions in this paper
44 are:

- 45 1. We develop a biologically plausible RNN model with important biological properties in-
46 cluding heterogeneous neuronal timescales, uniform neurotransmitter release characteristics
47 within individual neurons, and sensory input constraints.
- 48 2. We introduce a set of complex decision-making tasks that mimic empirical neuroscience
49 studies, featuring: a) orthogonal manipulation of bottom-up and top-down signals and b)
50 diverse top-down signaling with varying temporal dynamics.
- 51 3. Our results unveil mechanistic insights into the integration of top-down and bottom-up
52 signals during information processing in a biological context.
- 53 4. We present a comprehensive theoretical framework delineating the computational principles
54 employed by our RNN model to flexibly harness distinct top-down signals to guide the
55 analysis of sensory input, thereby optimizing task performance.

56 **Biologically plausible recurrent neural network (RNN) model**

In this study, we employed the following continuous-variable recurrent neural network (RNN) model:

$$\tau \frac{dx}{dt} = -x(t) + wr(t) + w_{in}u(t) \quad (1)$$

$$r(t) = \sigma(x(t)) = \frac{1}{1 + \exp(-x(t))} \quad (2)$$

$$o(t) = w_{out}r(t) + b \quad (3)$$

57 where $\tau \in \mathbb{R}^{1 \times N}$ refers to the synaptic time constants, $x \in \mathbb{R}^{N \times T}$ denotes the synaptic current
 58 variable from from N units across T time-points. By applying a sigmoid nonlinearity, we estimated
 59 the firing rates $r \in \mathbb{R}^{N \times T}$ based on the synaptic current values (x). The connection weights from
 60 the time-varying inputs ($u \in \mathbb{R}^{N_{in} \times T}$; N_{in} , the number of stimuli) to the network were represented
 61 by the weight matrix $w_{in} \in \mathbb{R}^{N \times N_{in}}$. Additionally, $w \in \mathbb{R}^{N \times N}$ contains connection weights between
 62 the N units. To compute the network output ($o \in \mathbb{R}^{1 \times T}$), we linearly combined all the firing rates
 63 specified by the output connection weight matrix, $w_{out} \in \mathbb{R}^{1 \times N}$, along with the constant bias
 64 term, b . The network size (N) was set to 1000, and the input signals were given to all the units
 65 except for the RNNs used in **Hierarchical organization of inhibitory-to-inhibitory feedback**
 66 **connections** where only the first 200 units received the input.

67 The dynamics were discretized using the first-order Euler approximation method and using the
 68 step size (Δt) of 5 ms:

$$x_t = \left(1 - \frac{\Delta t}{\tau}\right)x_{t-1} + \frac{\Delta t}{\tau}(wr_{t-1} + w_{in}u_{t-1}) \quad (4)$$

69 where $x_t = x(t)$. Given that the units in this network model communicate through differentiable
 70 and continuous signals, a gradient-descent supervised method, known as backpropagation through
 71 time (BPTT; [5]), was employed to train our RNNs to perform cognitive tasks. Specifically, the
 72 trainable parameters of the model included w , τ , w_{out} , and b . We used Adam (adaptive moment
 73 estimation) optimization algorithm to update these parameters. The learning rate was set to 0.01,
 74 and the TensorFlow default values were used for the rest of the parameters including the first
 75 and second moment decay rates. To further impose biological constraints, we enforced Dale's
 76 law (uniform neurotransmitter release characteristics within separate excitatory and inhibitory
 77 neurons) using methods similar to those implemented in previous studies ([2, 6, 7]). To adhere to
 78 empirical findings regarding the ratio of excitatory and inhibitory neurons observed in the brain,
 79 each RNN consists of 80% excitatory and 20% inhibitory units (i.e., E-I ratio of 80/20; [8–10]).

Importantly, instead of fixing the synaptic decay constant (τ) to a fixed value for all the units, we optimized the parameter for each unit. The parameter was trained to range from 20 ms to 125 ms to model heterogeneous synaptic dynamics of different receptors in the cortex [11, 12]. We initialized the synaptic decay time constant parameter (τ) using

$$\tau = \sigma(\mathcal{N}(0, 1))\tau_{step} + \tau_{min}$$

80 where $\sigma(\cdot)$ is the sigmoid function and $\mathcal{N}(0, 1)$ refers to the standard normal Gaussian distribution.
81 $\tau_{min} = 20$ ms and $\tau_{step} = 105$ ms were used to constrain the parameter to range from 20 ms to
82 125 ms.

83 The schematic diagram of the model is shown in fig. 1a. All the models were implemented with
84 TensorFlow 1.10.0 and trained on NVIDIA GPUs (RTX A4000).

85 Experiments

86 **Dynamic contextual decision-making task.** With the goal of elucidating the computational
87 principles governing the seamless integration of diverse information streams in dynamic information
88 processing, we developed a novel dynamic context decision-making task for our RNN models. In this
89 task, we orthogonally manipulated bottom-up sensory signals as well as two distinct time-varying
90 top-down instruction signals: task context and attention. This experimental design enabled us to
91 independently investigate the temporal dynamics and circuit configurations involved in processing
92 each signal and their subsequent integration. Specifically, we trained RNNs to perform variants
93 of a delayed match-to-sample (DMS) task paired with an opposite anti-DMS task (fig. 1a,b). On
94 each trial, two stimuli were presented sequentially, containing information about either one (fig. 1a)
95 or two sensory modalities (fig. 1b; only one of which was contextually relevant in each trial). In
96 the pro-DMS task, the agent’s goal was to produce a +1 output when the two stimuli of the
97 relevant modality matched, and a -1 when they did not. Conversely, the opposite responses were
98 expected when an anti-DMS context was instructed. In the one-modality DMS task, a task cue
99 signal representing which task to perform (pro-DMS vs. anti-DMS) was presented during the
100 instruction period. In the two-modality task, the signal for which modality to attend to (attention
101 cue) and which task to perform (task cue) was conveyed during the instruction period. Crucially,
102 we manipulated the temporal dynamics of instructions, by presenting instructions early (before
103 stimuli) or late (during the delay period in between two stimuli). Stimuli and instruction signals
104 were presented for 250 ms each.

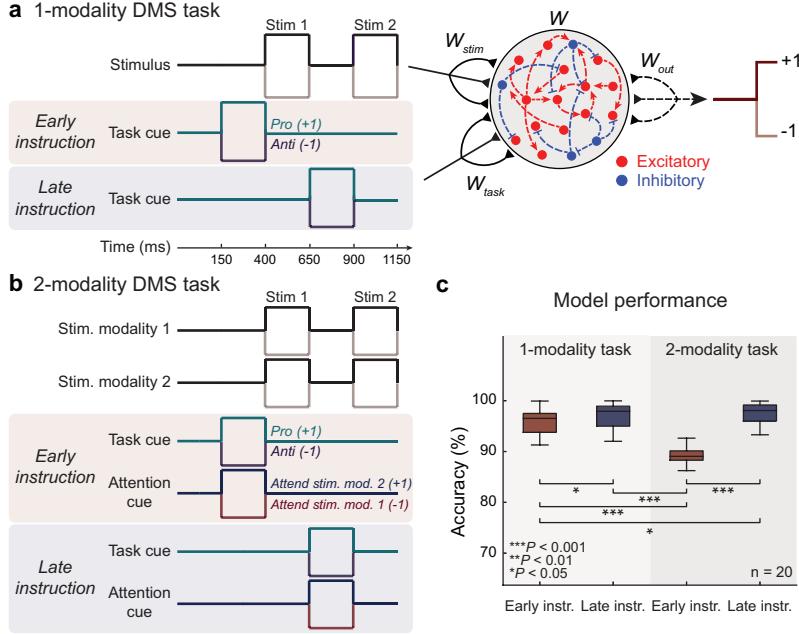


Fig. 1 | Delayed match-to-sample (DMS) tasks and model schematic. **a)** Schematic of a 1-modality DMS task with sequential stimuli separated by a delay interval. A task cue, indicating the stimulus-response mapping, is presented either before (*Early instruction*) or after stimulus 1 (*Late instruction*). The RNNs generate an output of +1 or -1 for matched and mismatched trials, respectively (*Pro*; and vice versa for *Anti*). **b)** Schematic of a 2-modality DMS task where each stimulus comprises two modalities. In addition to the task cue, an attention cue is presented on each trial, indicating the modality on which the RNNs need to perform the DMS task. The trial timing for both 1- and 2-modality DMS tasks is depicted in *a*). The task and attention cues are presented either before (*Early instruction*) or after stimulus 1 (*Late instruction*). **c)** Testing performance of the RNNs on the DMS tasks. The average accuracy of the trained RNNs is presented, with statistical significance computed using two-sided Wilcoxon rank-sum tests. Boxplot: central lines, median; bottom and top edges, lower and upper quartiles; whiskers, $1.5 \times$ interquartile range; outliers are not plotted.

105 **Training & testing protocol.** Our model training was deemed successful if the following two
 106 criteria were satisfied within the first 70,000 epochs: 1) loss value (defined as the root mean squared
 107 error between the network output and target signals) < 7 ; and 2) task performance $> 95\%$ for
 108 the one-modality task and $> 85\%$ for the two-modality task. If the network did not meet the
 109 criteria within the first 70,000 epochs, the training was terminated. For each successfully trained
 110 RNN, we simulated 400 trials for the one-modality task and 1000 trials for the two-modality task,
 111 randomizing each trial to contain one possible combination of instruction order, task instruction,
 112 and stimulus identities. Using these criteria, we trained 20 RNNs for each task. Investigating the
 113 performance based on the timing of the top-down instruction signals revealed that the average task

114 accuracy was lower when the instruction was given before the first stimulus presentation for both
115 tasks (fig. 1c).

116 Characterization of network motifs and circuit configurations

117 **Decoding analysis.** To characterize the stability of information representation in our RNN mod-
118 els, we performed a decoding analysis utilizing support vector machine (SVM) models. For each
119 trained RNN, we obtained firing rate estimates of all the neurons across multiple trials. We then
120 used the first half of the trials to train SVM models to decode task-related information. Subse-
121 quently, we tested each SVM on the left-out trials and computed the decoding accuracy at each
122 time step.

123 Additionally, to determine whether RNNs preserved the geometry of stimulus coding across
124 task conditions, we performed cross-condition generalization analyses, training and testing across
125 time bins and across task conditions of interest[13]. We determined that top-down task signals,
126 indicating pro-DMS vs. anti-DMS trials, were easily decodable from both the one-modality and
127 two-modality models, immediately following the instruction presentation (fig. 2a-d). For both
128 models, at the time of the second stimulus presentation, task signal decoding was stronger along
129 the diagonal (i.e., training and testing at the time of second stimulus) than when training at
130 the time of the first stimulus and testing cross-temporally at the second stimulus. This indicates
131 that while task demands are represented throughout the trial, there is a dynamic shift in their
132 representations in between stimulus presentation windows (as captured by the linear SVMs).

133 Furthermore, we set out to decode the top-down attention signal indicating which modality is
134 relevant in each trial for the two-modality model. Once again, we were able to successfully decode
135 these signals from neural populations immediately following the instruction period (fig. 2e-f). More
136 importantly, we could also decode the identity of the attention signal during the first stimulus and
137 delay periods for both attended (fig. 2g) and unattended (fig. 2h) modalities (for trials where the
138 attention signal was given before the first stimulus). This indicates that both top-down attention
139 and bottom-up stimulus signals are encoded by the network, and that the model is able to maintain
140 bottom-up modality information even in the presence of top-down signals. We also measured SVM
141 decoding performance for stimulus identity when training in early instruction trials and testing
142 in late instruction trials, and vice versa. We found stable representation for stimulus identities in
143 the one-modality models, as well as both sensory modalities in the two-modality model variants
144 (Supplementary Figure 1).

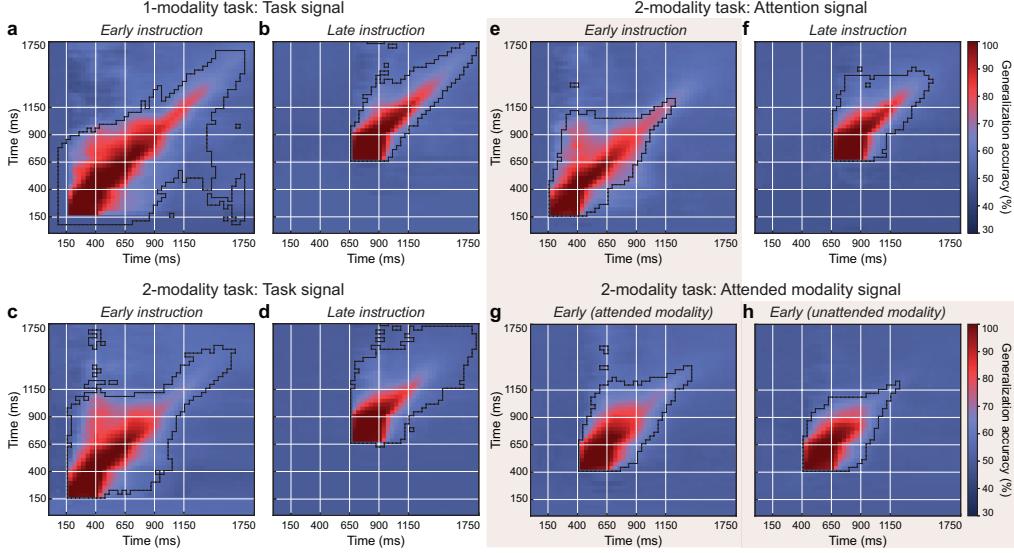


Fig. 2 | Temporal encoding of top-down signals during DMS tasks. Cross-temporal decodability of task and attention signals in the 1- and 2-modality DMS tasks is shown as the average generalization accuracy for each pair of training/testing time windows. Highlighted boundaries indicate significant decoding accuracies, $p < 1e - 20$, binomial test, Bonferroni corrected. **a)** and **b)** Dynamic coding of the task signal (*Pro/Anti*) extending to the onset of the response period (1150 ms) is evident in both the *Early* and *Late instruction* conditions of the 1-modality task. **c)** and **d)** Similar patterns of dynamic codes for the task signal are observed for the 2-modality task, illustrating the robust encoding of the task identity. **e)** and **f)** Reliable decoding of the attention signal, indicating whether to attend to stimulus modality 1 or 2, can be observed over time in both the *Early* and *Late instruction* conditions. **g)** and **h)** In the *Early instruction* condition, the representation of the attention signal can be found during the stimulus 1 and delay interval, in the RNN activity corresponding to both the attended (**g**) and unattended stimulus modalities (**h**). White lines denote the reference event of each trial period (task/instruction cue, stimulus onsets, response).

145 **State space analysis.** To characterize how RNNs flexibly adapted to different task and attention
 146 demands, we performed a state space analysis, treating RNNs as a high-dimensional dynamical
 147 system. The goal of this analysis was to map how the energy landscape of the network changed as
 148 a function of inputs, providing a summarized account for the effect of task and attention demands
 149 over network dynamics. For this, we defined and minimized an auxiliary kinetic energy function
 150 $q(x, u) = \frac{1}{2}|F(x, u)|^2$, where $F(x, u) = \frac{dx}{dt}|_u$, utilizing the previously defined system dynamics from
 151 Equation 3, where x is the population state vector and u is the input to the network, containing
 152 stimulus, task, and modality cue information. Notably, $q(x, u) = 0$ if and only if the neural
 153 population state x is a fixed point, a point in which the velocity of the dynamic system is zero
 154 along every dimension [14, 15]. Along with fixed points, finding non-zero local minima of the

155 energy landscape, defined as slow points, is also useful, since both fixed points and slow points may
 156 be used to locally linearize the dynamic system and summarize its activity [14].

157 To identify slow points and visualize networks' energy landscape, we minimized $q(x, u)$ for
 158 all the possible values of u during the early instruction period, and demonstrated how top-down
 159 instructions reshape the activity of one example network prior to receiving bottom-up stimulus
 160 information. For the one-modality model, we tested the effect of instruction inputs in the pro-DMS,
 161 anti-DMS, or late instruction conditions. For the two-modality model, we tested u in each of the
 162 4 possible combinations between pro/anti-DMS inputs and the two cued attentional modalities.

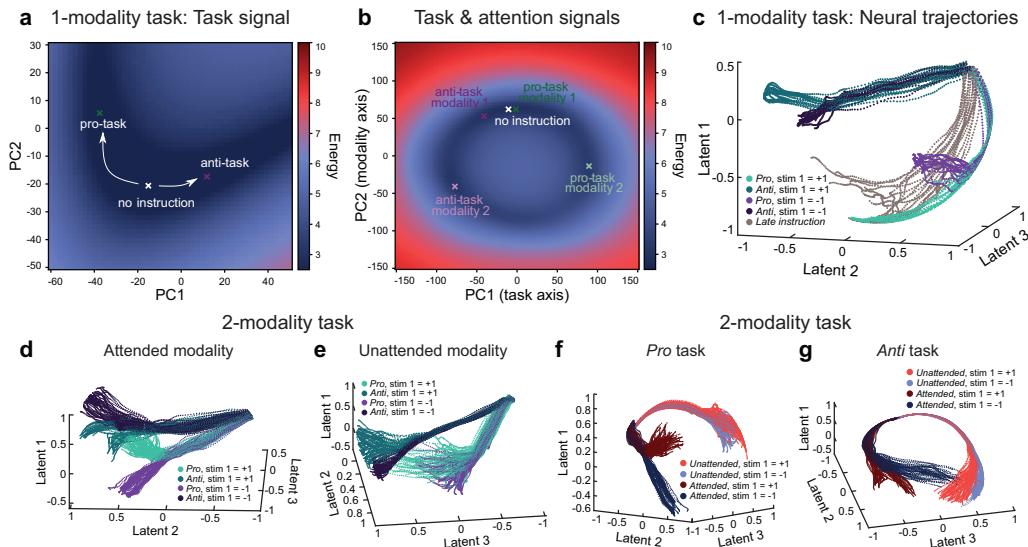


Fig. 3 | Task and attention representations in RNN state space. **a)** Heat map of kinetic energy in the one-modality RNN during the instruction period, when no task information is available. The middle cross indicates the local energy minimum for the no-task condition, while the left and right crosses indicate the new local energy minima for the pro-DMS and anti-DMS conditions, respectively. **b)** Same as **a**, for the two-modality RNN. Clockwise from the left-most point, minima indicate the following conditions: (anti-DMS and first modality attended), (anti-DMS and second modality attended), (pro-DMS and first modality attended), and (pro-DMS and second modality attended). **c)** Latent space one-modality RNN trajectories during instruction and first stimulus periods for all combinations of task instructions and stimulus identities. **d)** Same as **c**, for two-modality networks, highlighting the attended stimulus modality. **e)** Same as **c**, for two-modality networks, highlighting the unattended stimulus modality. **f)** Latent space two-modality RNN trajectories for pro-DMS trials, splitting trajectories by the identity of attended versus unattended cues. **g)** Same as **f**, for anti-DMS trials.

163 We performed minimization both in principal component space and in native space (with the
 164 *sklearn* module *optimize.minimize*), finding qualitatively similar results. For this analysis, we

165 defined principal components using the average activity of neurons in the early instruction period as
166 features. We found that the one-modality models defined an inverted U-shape energy landscape, in
167 which each branch contains one local minimum, for either DMS or anti-DMS inputs. Additionally, a
168 third local minimum exists in the middle, for when the task instruction has not yet been presented
169 to the network (fig. 3a and Supplementary Figure 2). Notably, each task’s local minimum was
170 roughly equidistant from the central local minimum. Interestingly, the example 2-modality network
171 defined a ring-like energy landscape containing 4 distinct local minima for each combination of task
172 instruction and attention cue (fig. 3b and Supplementary Figure 2).

173 Therefore, the top-down task input signals reshape the energy landscape of the network and
174 nudge neural trajectories to a preferred path prior to stimulus presentation. Subsequently, the
175 same stimulus inputs can lead to separate trajectories due to their distinct initial conditions. We
176 confirmed this by selecting an example network and visualizing its latent network trajectories
177 during early instruction and first stimulus periods, as obtained through the unsupervised manifold
178 discovery method CEBRA [16]. The presence of task instructions distinguished trajectories, which
179 further bifurcated upon receiving bottom-up sensory information about stimulus identity (fig. 3c).

180 We further used the CEBRA manifold representation to summarize the distinct effects of
181 sensory-based (attention cue) and sensory-independent top-down information (task cue) over neu-
182 ral trajectories. In the two-modality task, the network trajectories bifurcated based on the task
183 instruction (pro-DMS or anti-DMS) regardless of which modality was cued (fig. 3d,e). More specif-
184 ically, the first stimulus identity (-1 or +1) for the cued modality were coded in distinct trajectories
185 following the stimulus presentation (fig. 3d). On the other hand, the trajectories collapsed across
186 the two stimulus identities for the unattended modality, even though the task coding was preserved
187 (fig. 3e). We further confirmed this result by investigating whether attended versus unattended
188 modalities were differentially represented in the latent space for pro-DMS trials versus anti-DMS
189 trials (for both early and late instruction conditions; fig. 3f,g and Supplementary Figure 3, re-
190 spectively). For both pro-DMS and anti-DMS trials, the latent trajectories split depending on
191 the attention signal such that the attended trajectories further bifurcated based on the attended
192 stimulus identities. However, the trajectories for the unattended modality did not separate based
193 on the unattended stimulus identities. Taken together, these results highlight the differential uti-
194 lization of top-down signals and provide compelling evidence that our model can asymmetrically
195 employ distinct top-down signaling to optimize efficient coding of bottom-up sensory information.

196 **Inhibitory-to-inhibitory connections carry top-down information**

197 To better understand the circuit mechanisms required for flexible integration of the sensory signals
198 and task-specific top-down information, we first identified neurons selective to the task cue signal
199 (i.e., pro-DMS vs. anti-DMS) in the one-modality RNNs. Using the method described in the sup-
200 plementary material, we discovered four distinct subgroups of neurons within each one-modality
201 RNN: excitatory and inhibitory neurons tuned to the pro-DMS task cue signal, as well as excitatory
202 and inhibitory neurons tuned to the anti-DMS task context. Interestingly, we observed that reduc-
203 ing the synaptic inputs (by reducing incoming synaptic weights by 50%) to the inhibitory neurons
204 specific to the task context effectively eliminated the task cue signaling during the instruction
205 window prior to the presentation of the first stimulus. For instance, decreasing the synaptic inputs
206 to the inhibitory neurons biased toward the anti-DMS task in a one-modality RNN disrupted the
207 anti-DMS task signaling. Consequently, the network exhibited a tendency to perform the pro-DMS
208 task irrespective of the task context, as illustrated in (fig. 4a). By visualizing the neural dynamics
209 of this “lesioned” network using CEBRA, we confirmed that the neural trajectories were no longer
210 separated by the task cue signal during the instruction period (fig. 4b).

211 Motivated by the above observation, we systematically reduced synaptic inputs (by reducing
212 incoming synaptic weights by 50%) to the anti-DMS preferring neurons in all the RNNs trained
213 for the one-modality network model ($n = 20$ RNNs). This manipulation revealed that the largest
214 decrease in model performance for the anti-DMS task occurred when inhibitory synaptic inputs to
215 the anti-DMS preferring inhibitory neurons were decreased (indicated as I→I in fig. 4c). Addition-
216 ally, reducing the inhibitory inputs to the excitatory neurons biased toward the anti-DMS task also
217 resulted in a significant decrease in model performance (I→E in fig. 4c). Manipulating excitatory
218 synaptic inputs did not lead to significant changes (E→I and E→E in fig. 4c). Similar results
219 were observed when we lesioned neurons preferring the pro-DMS task context (see Supplementary
220 Figure 4).

221 Performing the above analysis in the two-modality RNNs showed that disrupting the inhibitory
222 input signaling to the inhibitory neurons specific to the anti-DMS task did not lead to a comparable
223 drop in performance (fig. 4d). Furthermore, we found that the networks remained robust to
224 manipulations of the synaptic inputs targeting the neurons selective for the attention signal. For
225 instance, when we reduced the synaptic inputs to the neurons tuned to the first modality signal,
226 we did not observe a complete loss of the attention signaling (fig. 4e).

227 The robustness of the two-modality networks to the above manipulations could be attributed

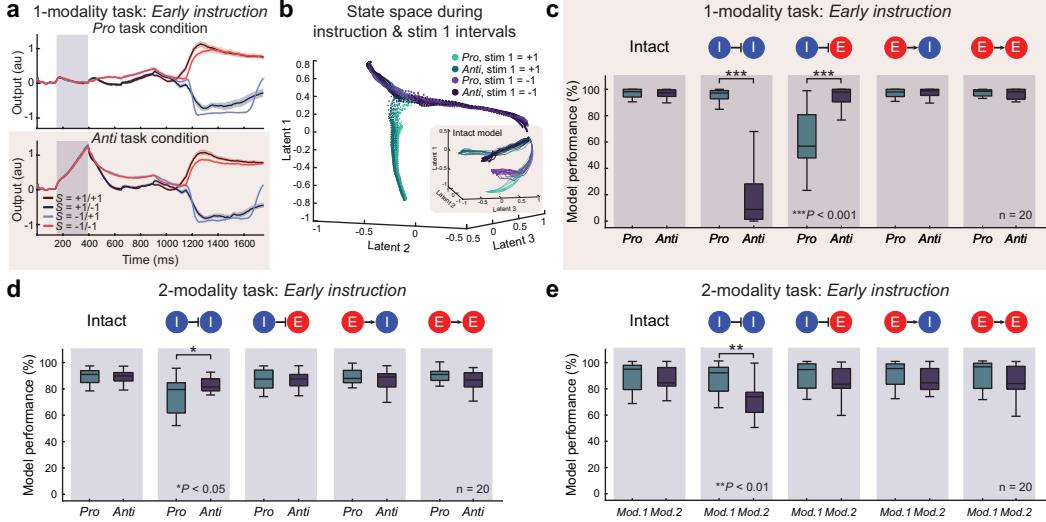


Fig. 4 | I→I connections encode top-down information. **a)** Reducing incoming inhibitory synaptic weights to inhibitory neurons selective to the anti-DMS cue by 50% in an example one-modality RNN led to the network performing the pro-DMS even when the anti-DMS cue was given (bottom). **b)** CEBRA neural trajectories of the lesioned model shown in **a** during the first instruction and stimulus periods. Inset, trajectories from the intact network (identical to fig. 3c). **c)** Performance of all 20 RNNs trained to perform the one-modality task when incoming synaptic weights (separated by excitatory and inhibitory weights) to anti-DMS preferring neurons (separated by excitatory and inhibitory) were reduced by 50%. **d)** Similar to **c**, for the two-modality RNNs. **e)** Similar to **d**, but when modality 1 attention cue selective neurons were lesioned. Statistical significance computed using the Wilcoxon signed rank test. Boxplot: central lines, median; bottom and top edges, lower and upper quartiles; whiskers, $1.5 \times$ interquartile range; outliers are not plotted.

228 to the emergence of neurons that multiplex or encode multiple top-down factors (i.e., task and
 229 attention cues) simultaneously (Supplementary Figure 5). These neurons exhibit the capacity to
 230 integrate and represent multiple sources of top-down information, contributing to the network's
 231 ability to maintain stable performance despite perturbations in a subgroup of neurons.

232 Hierarchical organization of inhibitory-to-inhibitory feedback connections

233 In this section, we extended our RNN model to introduce an additional constraint to simulate a
 234 cortical hierarchy. In contrast to the previous configuration where all neurons in a network received
 235 input signals, we imposed a constraint restricting the input signals to a specific subpopulation
 236 within the network, effectively creating a distinct group of “sensory” neurons (fig. 5a). The network
 237 size (N) remained unchanged at 1000 neurons, but only 200 neurons received the input signals.
 238 Training this model on the one-modality task revealed that both “sensory” neurons and “non-

239 “sensory” units encoded both bottom-up and top-down information (latent neural trajectories of
 240 the “sensory” units from an example model shown in fig. 5b).

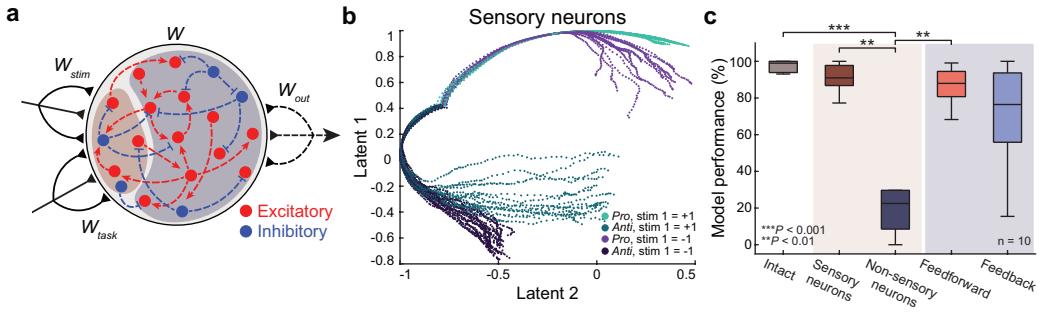


Fig. 5 | RNN model with input signaling constraints. **a)** Model schematic. The stimulus and task-related signals are only given to the first 200 units. **b)** CEBRA neural trajectories of the “sensory” neurons in an example RNN trained to perform the one-modality task. **c)** Anti-DMS task performance (across 10 RNNs) when I→I connections selective for the anti-DMS task signal were reduced by 50%. Statistical significance computed using Kruskal-Wallis test with Dunn’s post hoc test. Boxplot: central lines, median; bottom and top edges, lower and upper quartiles; whiskers, $1.5 \times$ interquartile range; outliers are not plotted.

241 We performed lesion studies on these networks, similar to the ones presented in the previous
 242 section. We found that selectively disrupting the inhibitory-to-inhibitory connections specific to the
 243 task cue within the non-sensory area had a pronounced effect on the model’s performance (fig. 5c).
 244 Lesioning I→I signaling from the non-sensory area to the sensory area also resulted in a significant
 245 decrease in the task performance. Notably, disrupting the inhibitory-to-inhibitory connections
 246 within the “sensory” area did not lead to any impairment in task performance, suggesting that
 247 these synapses do not participate in stimulus encoding. Given our results from the previous section
 248 highlighting the importance of I→I synapses in encoding of the task context, the current findings
 249 in fig. 5c can be attributed to the higher abundance of I→I connections within the non-sensory area
 250 (which was modeled after experimental studies showing the increase in diverse inhibitory neuronal
 251 populations along the cortical hierarchy [17]).

252 Discussion

253 In this work, we proposed a framework that employs a biologically plausible RNN model and
 254 variants of a delayed match-to-sample (DMS) task to study how bottom-up (i.e., encoding sensory
 255 stimuli) and top-down (i.e., task instructions, attention) are integrated to facilitate flexible decision
 256 making. While previous works modeled flexible dynamic learning in biological [1, 18] and artificial
 257 [15, 19] agents, our work revealed the mechanisms through which an artificial neural network meets

258 such demands dynamically under relevant biological constraints. Specifically, we constrained the
259 connectivity structure of our RNN model with Dale’s principle and enforced heterogeneous synaptic
260 timescales along with regionally specialized sensorimotor areas to model the neocortical hierarchy.

261 Using this model with important biological constraints, we demonstrated how our model dy-
262 namically coded signals containing sensory, instructional, and attentional information, and how
263 this coding could be selectively disrupted by lesioning specific subsets of neurons, which revealed
264 a privileged role for inhibitory projections in enabling instruction-dependent task performance.
265 These results provide a substrate for future work in developing agents which can dynamically gen-
266 eralize their performance across different sensory environments, accommodating for fluctuations
267 in top-down contextual signals. Additionally, our work provides a framework to generate testable
268 hypotheses to further investigate the neural mechanisms of flexible and dynamic decision making
269 in biological agents.

270 This work includes a number of limitations. While we enforced several elements of biological
271 fidelity, future work can still be done to further approximate neural network structure to cortical
272 hierarchical specialization, an approach which has been successful in the sensory domain with
273 convolutional neural networks [20, 21]. Additionally, this mechanistic approach must be extended
274 beyond the DMS paradigm to more complex, naturalistic tasks, involving high-dimensional state
275 and decision spaces. Another limitation of our study is the biological plausibility of gradient-
276 descent optimization. One way to address this is to compare our findings with RNNs trained
277 using alternative methods, including non-gradient descent approaches. One notable alternative
278 is the First-Order Reduced and Controlled Error (FORCE) learning algorithm, which has shown
279 success in training rate-based and spiking RNNs [22, 23]. However, it’s important to note that
280 FORCE learning currently does not support training the synaptic decay time constant, making
281 direct comparisons with our models challenging.

282 **References**

- 283 1. Mante, V., Sussillo, D., Shenoy, K. V. & Newsome, W. T. Context-dependent computation by recurrent
284 dynamics in prefrontal cortex. *Nature* **503**, 78–84 (2013).
- 285 2. Song, H. F., Yang, G. R. & Wang, X.-J. Training Excitatory-Inhibitory Recurrent Neural Networks
286 for Cognitive Tasks: A Simple and Flexible Framework. *PLOS Computational Biology* **12**, e1004792
287 (2016).
- 288 3. Miconi, T. Biologically plausible learning in recurrent neural networks reproduces neural dynamics
289 observed during cognitive tasks. *Elife* **6**, e20899 (2017).
- 290 4. Yang, G. R., Joglekar, M. R., Song, H. F., Newsome, W. T. & Wang, X.-J. Task representations in
291 neural networks trained to perform many cognitive tasks. *Nature Neuroscience* **22**, 297–306 (2019).
- 292 5. Werbos, P. J. Backpropagation through time: what it does and how to do it. *Proceedings of the IEEE*
293 **78**, 1550–1560 (1990).
- 294 6. Kim, R., Li, Y. & Sejnowski, T. J. Simple framework for constructing functional spiking recurrent
295 neural networks. *Proceedings of the national academy of sciences* **116**, 22811–22820 (2019).
- 296 7. Rungratsameetaweemana, N., Kim, R. & Sejnowski, T. J. Internal noise promotes heterogeneous synaptic
297 dynamics important for working memory in recurrent neural networks. *bioRxiv*. eprint: <https://www.biorxiv.org/content/early/2022/10/18/2022.10.14.512301.full.pdf>. <https://www.biorxiv.org/content/early/2022/10/18/2022.10.14.512301> (2022).
- 300 8. Hendry, S. H., Schwark, H., Jones, E. & Yan, J. Numbers and proportions of GABA-immunoreactive
301 neurons in different areas of monkey cerebral cortex. *Journal of Neuroscience* **7**, 1503–1519 (1987).
- 302 9. Sherwood, C. C., Raghanti, M. A., Stimpson, C. D., Bonar, C. J., de Sousa, A. A., Preuss, T. M. &
303 Hof, P. R. Scaling of inhibitory interneurons in areas V1 and V2 of anthropoid primates as revealed by
304 calcium-binding protein immunohistochemistry. *Brain, Behavior and Evolution* **69**, 176–195 (2007).
- 305 10. Alreja, A., Nemenman, I. & Rozell, C. J. Constrained brain volume in an efficient coding model explains
306 the fraction of excitatory and inhibitory neurons in sensory cortices. *PLOS Computational Biology* **18**,
307 e1009642 (2022).
- 308 11. Duarte, R., Seeholzer, A., Zilles, K. & Morrison, A. Synaptic patterning and the timescales of cortical
309 dynamics. *Current Opinion in Neurobiology* **43**, 156–165 (2017).
- 310 12. Zador, A. M. & Dobrunz, L. E. Dynamic synapses in the cortex. *Neuron* **19**, 1–4 (1997).
- 311 13. Bernardi, S., Benna, M. K., Rigotti, M., Munuera, J., Fusi, S. & Salzman, C. D. The geometry of
312 abstraction in the hippocampus and prefrontal cortex. *Cell* **183**, 954–967 (2020).
- 313 14. Sussillo, D. & Barak, O. Opening the black box: low-dimensional dynamics in high-dimensional recur-
314 rent neural networks. *Neural computation* **25**, 626–649 (2013).

- 315 15. Driscoll, L., Shenoy, K. & Sussillo, D. Flexible multitask computation in recurrent networks utilizes
316 shared dynamical motifs. *bioRxiv*, 2022–08 (2022).
- 317 16. Schneider, S., Lee, J. H. & Mathis, M. W. Learnable latent embeddings for joint behavioural and
318 neural analysis. *Nature*, 1–9 (2023).
- 319 17. Kim, Y., Yang, G. R., Pradhan, K., Venkataraju, K. U., Bota, M., García del Molino, L. C., Fitzgerald,
320 G., Ram, K., He, M., Levine, J. M., Mitra, P., Huang, Z. J., Wang, X.-J. & Osten, P. Brain-wide Maps
321 Reveal Stereotyped Cell-Type-Based Cortical Architecture and Subcortical Sexual Dimorphism. *Cell*
322 **171**, 456–469.e22. ISSN: 0092-8674 (2017).
- 323 18. Yang, G. R., Cole, M. W. & Rajan, K. How to study the neural mechanisms of multiple tasks. *Current*
324 *opinion in behavioral sciences* **29**, 134–143 (2019).
- 325 19. Yang, G. R., Joglekar, M. R., Song, H. F., Newsome, W. T. & Wang, X.-J. Task representations in
326 neural networks trained to perform many cognitive tasks. *Nature neuroscience* **22**, 297–306 (2019).
- 327 20. Yamins, D. L. & DiCarlo, J. J. Using goal-driven deep learning models to understand sensory cortex.
328 *Nature neuroscience* **19**, 356–365 (2016).
- 329 21. Nayebi, A., Bear, D., Kubilius, J., Kar, K., Ganguli, S., Sussillo, D., DiCarlo, J. J. & Yamins, D. L.
330 Task-driven convolutional recurrent models of the visual system. *Advances in neural information pro-*
331 *cessing systems* **31** (2018).
- 332 22. Sussillo, D. & Abbott, L. Generating Coherent Patterns of Activity from Chaotic Neural Networks.
333 *Neuron* **63**, 544 –557 (2009).
- 334 23. Nicola, W. & Clopath, C. Supervised learning in spiking neural networks with FORCE training. *Nature*
335 *Communications* **8**, 2208 (2017).

336 Acknowledgements

337 We gratefully acknowledge the support of the Swartz Foundation and Mission funding from the
338 cooperative agreement under the United States Army Research Laboratory. Any opinions, findings,
339 and conclusions or recommendations expressed in this material are those of the authors and do
340 not necessarily reflect the views, policies, or endorsements, either expressed or implied, of the
341 DEVCOM Army Research Laboratory or the U.S. Government.

342 Author contributions

343 All authors designed the study, performed the experiments, and wrote the paper.

344 Declaration of interests

345 The authors declare no competing interests.

Supplementary Figures

Supplementary Notes

Supplementary Table

	LIF	QIF
Membrane time constant (τ_m)	10 ms	10 ms
Absolute refractory period	2 ms	2 ms
Synaptic rise time (τ_r)	2 ms	2 ms
Constant bias current	-40 pA	0 pA
Spike threshold	-40 mV	30 mV
Spike reset voltage	-65 mV	-65 mV

Supplementary Table 1 | Parameter values used to construct LIF and QIF networks.