1	Cell type-specific connectome predicts distributed
2	working memory activity in the mouse brain
3	Xingyu Ding ^{1,*} , Sean Froudist-Walsh ^{1,2,*} , Jorge Jaramillo ^{1,3,*} , Junjie Jiang ^{1,4} ,
4	and Xiao-Jing Wang ^{1,+}
5	¹ Center for Neural Science, New York University, New York, NY 10003, USA
6	$^2\mathrm{Bristol}$ Computational Neuroscience Unit, School of Engineering Mathematics
7	and Technology, University of Bristol, Bristol BS8 1UB, UK
8	$^{3}\mathrm{Campus}$ Institute for Dynamics of Biological Networks, Goettingen, Germany
9	$^4\mathrm{The}$ Key Laboratory of Biomedical Information Engineering of Ministry of
10	Education, Institute of Health and Rehabilitation Science, School of Life Science
11	and Technology, Research Center for Brain-inspired Intelligence, Xi'an Jiaotong
12	University, No.28, West Xianning Road, Xi'an, 710049, Shaanxi, P. R. China.
13	*co-first authors
14	⁺ lead contact: xjwang@nyu.edu
15	October 31, 2023

16 Abstract

Recent advances in connectome and neurophysiology make it possible to probe whole-brain 17 mechanisms of cognition and behavior. We developed a large-scale model of the mouse 18 multiregional brain for a cardinal cognitive function called working memory, the brain's 19 ability to internally hold and process information without sensory input. The model is built 20 on mesoscopic connectome data for inter-areal cortical connections and endowed with a 21 macroscopic gradient of measured parvalbumin-expressing interneuron density. We found 22 that working memory coding is distributed yet exhibits modularity; the spatial pattern of 23 mnemonic representation is determined by long-range cell type-specific targeting and density 24 of cell classes. Cell type-specific graph measures predict the activity patterns and a core 25 subnetwork for memory maintenance. The model shows numerous self-sustained internal 26 states (each engaging a distinct subset of areas). This work provides a framework to interpret 27 large-scale recordings of brain activity during cognition, while highlighting the need for cell 28 type-specific connectomics. 29

30 Introduction

In contrast to our substantial knowledge of local neural computation, such as orientation 31 selectivity in the primary visual cortex or the spatial map of grid cells in the medial entorhinal 32 cortex, much less is understood about distributed processes in multiple interacting brain 33 regions underlying cognition and behavior. This has recently begun to change, as advances in 34 new technologies enable neuroscientists to probe neural activity at single-cell resolution and 35 on a large-scale by electrical recording or calcium imaging of behaving animals (Jun et al. 36 2017; Steinmetz et al. 2019; Stringer et al. 2019; Musall et al. 2019; Steinmetz et al. 2021), 37 ushering in a new era of neuroscience investigating distributed neural dynamics and brain 38 functions (Wang 2022).

To be specific, consider a core cognitive function called working memory, the ability 40 to temporally maintain information in mind without external stimulation (Baddeley 2012). 41 Working memory has long been studied in neurophysiology using delay-dependent tasks, where 42 stimulus-specific information must be stored in working memory across a short time period 43 between a sensory input and a memory-guided behavioral response (Fuster and Alexander 44 1971; Funahashi et al. 1989; Goldman-Rakic 1995; Wang 2001). Delay-period mnemonic 45 persistent neural activity has been observed in multiple brain regions, suggesting distributed 46 working memory representation (Suzuki and Gottlieb 2013; Leavitt et al. 2017; Christophel 47 et al. 2017; Xu 2017; Dotson et al. 2018). Connectome-based computational models of the 48 macaque cortex found that working memory activity depends on interareal connectivity 49 (Murray et al. 2017; Jaramillo et al. 2019), macroscopic gradients of synaptic excitation 50 (Wang 2020; Mejias and Wang 2022) and dopamine modulation (Froudist-Walsh et al. 2021). 51 Mnemonic neural activity during a delay period is also distributed in the mouse brain 52 (Liu et al. 2014; Schmitt et al. 2017; Guo et al. 2017; Bolkan et al. 2017; Gilad et al. 2018). 53 The new recording and imaging techniques as well as optogenetic methods for causal analysis 54 (Yizhar et al. 2011), that are widely applicable to behaving mice, hold promise for elucidating 55 the circuit mechanism of distributed brain functions in rodents. Recurrent synaptic excitation 56 represents a neural basis for the maintenance of persistent neural firing (Goldman-Rakic 57 1995; D. J. Amit 1995; Wang 2021). In the monkey cortex, the number of spines (sites 58 of excitatory synapses) per pyramidal cell increases along the cortical hierarchy, consistent 59 with the idea that mnemonic persistent activity in association cortical areas including the 60 prefrontal cortex is sustained by recurrent excitation stronger than in early sensory areas. 61 Such a macroscopic gradient is lacking in the mouse cortex (Gilman et al. 2017; Ballesteros-62 Yáñez et al. 2010), raising the possibility that the brain mechanism for distributed working 63 memory representations may be fundamentally different between mice and monkeys. 64

In this paper we report a cortical mechanism of distributed working memory that does not depend on a gradient of synaptic excitation. We developed an anatomically-based model of the mouse brain for working memory, built on the recently available mesoscopic connectivity data of the mouse thalamocortical system (Oh et al. 2014; Gămănuţ et al. 2018; Harris

et al. 2019; Kim et al. 2017). Our model is validated by capturing large-scale neural activity 69 observed in recent mouse experiments (Guo et al. 2017; Gilad et al. 2018). Using this model, 70 we found that a decreasing gradient of synaptic inhibition mediated by parvalbumin (PV) 71 positive GABAergic cells (Kim et al. 2017; Fulcher et al. 2019; Wang 2020) and long-range 72 excitatory connections shape the distributed pattern of working memory representation. 73 Moreover, the engagement of inhibition through local and long range projections determines 74 the stability of the local circuits, further emphasizing the importance of inhibitory circuits. 75 A focus of this work is to examine whether anatomical connectivity can predict the 76 emergent large-scale neural activity pattern underlying working memory. Interestingly, traditional graph-theory measures of inter-areal connections, which ignore cell types of 78 projection targets, are uncorrelated with activity patterns. We propose new cell type-79 specific graph theory measures to overcome this problem, and differentiate contributions of 80 cortical areas in terms of their distinct role in loading, maintaining, and reading out the 81 content of working memory. Through computer-simulated perturbations akin to optogenetic 82 inactivations, a core subnetwork was uncovered for the generation of persistent activity. This 83 core subnetwork can be predicted based on the cell type-specific interareal connectivity, 84 highlighting the necessity of knowing the cell type targets of interareal connections in order 85 to relate anatomy with physiology and behavior. This work provides a computational and theoretical platform for cross-scale understanding of cognitive processes across the mouse 87 cortex. 88

⁸⁹ Results

⁹⁰ A decreasing gradient of PV interneuron density from sensory to ⁹¹ association cortex

Our large-scale circuit model of the mouse cortex uses inter-areal connectivity provided by 92 anatomical data within the 43-area parcellation in the common coordinate framework v3 atlas 93 (Oh et al. 2014) (Fig. 1A, Fig. 1 - supplement 1A). The model is endowed with area-to-area 94 variation of parvalbumin-expressing interneurons (PV) in the form of a gradient measured 95 from the qBrain mapping platform (Fig. 1 - supplement 1B) (Kim et al. 2017). The PV cell 96 density (the number of PV cells per unit volume) is divided by the total neuron density, to 97 give the PV cell fraction, which better reflects the expected amount of synaptic inhibition 98 mediated by PV neurons (Fig. 1B-C, neuron density is shown in Fig. 1 - supplement 1C). 99 Cortical areas display a hierarchy defined by mesoscopic connectome data acquired using 100 anterograde fluorescent tracers (Oh et al. 2014) (Fig. 1D-E). In Fig. 1F, the PV cell fraction 101 is plotted as a function of the cortical hierarchy, which shows a moderate negative correlation 102 between the two. Therefore, primary sensory areas have a higher density of PV interneurons 103 than association areas, although the gradient of PV interneurons does not align perfectly 104 with the cortical hierarchy. 105



Figure 1: Anatomical basis of the multi-regional mouse cortical model. (A). Flattened view of mouse cortical areas. Figure adapted from (Harris et al. 2019). (B). Normalized PV cell fraction for each brain area, visualized on a 3d surface of the mouse brain. Five areas are highlighted : VISp, Primary somatosensory area, barrel field (SSp-bfd), primary motor (MOp), MOs and PL. (C). The PV cell fraction for each cortical area, ordered. Each area belongs to one of five modules, shown in color. (Harris et al. 2019). (D). Hierarchical position for each area on a 3d brain surface. Five areas are highlighted as in (B), and color represents the hierarchy position. (E). Hierarchical positions for each cortical area. The hierarchical position is normalized and the hierarchical position of VISp is set to be 0. As in C), the colors represent the module that an area belongs to. (F). Correlation between PV cell fraction and hierarchy (Pearson correlation coefficient r = -0.35, p < 0.05).

¹⁰⁶ A whole-mouse cortex model with a gradient of interneurons

In our model, each cortical area is described by a local circuit (Fig. 2A), using a mean-field 107 reduction (Wong and Wang 2006) of a spiking neural network (Wang 2002). We use a 108 version of this model that has two excitatory neural pools selective for different stimuli and a 109 shared inhibitory neural pool to describe each cortical area. The model makes the following 110 assumptions. First, local inhibitory strength is proportional to PV interneuron density across 111 the cortex. Second, the inter-areal long-range connection matrix is given by the anterograde 112 tracing data (Oh et al. 2014; Knox et al. 2018; Wang et al. 2020). Third, targeting is 113 biased onto inhibitory cells for top-down compared with bottom-up projections. Therefore, 114 feedforward connections have a greater net excitatory effect than feedback connections, which 115 is referred to as counterstream inhibitory bias (CIB) (Mejias and Wang 2022; Javadzadeh 116 and Hofer 2022; Wang 2022). Briefly, we assume that long-range connections are scaled by 117 a coefficient that is based on the hierarchy of the source and target areas. According to 118 the CIB assumption, long-range connections to inhibitory neurons are stronger for feedback 119 connections and weaker for feedforward connections, while the opposite holds for long range 120 connections to excitatory neurons. 121

¹²² Distributed working memory activity depends on the gradient of ¹²³ inhibitory neurons and the cortical hierarchy

We simulated the large-scale network to perform a simple visual delayed response task that requires one of two stimuli to be held in working memory. We shall first consider the case in which the strength of local recurrent excitation is insufficient to generate persistent activity when parcellated areas are disconnected from each other. Consequently, the observed distributed mnemonic representation must depend on long-range interareal excitatory connection loops. Later in the paper we will discuss the network model behavior when some local areas are capable of sustained persistent firing in isolation.

The main question is: when distributed persistent activity emerges after a transient visual 131 input is presented to the primary visual cortex (VISp), what determines the spatial pattern 132 of working memory representation? After we remove the external stimulus, the firing rate in 133 area VISp decreases rapidly to baseline. Neural activity propagates throughout the cortex 134 after stimulus offset (Fig. 2B). Neural activities in the higher visual cortical areas (e.g. VISrl 135 and VISpl) show similar dynamics to VISp. In stark contrast, many frontal and lateral areas 136 (including prelimbic (PL), infralimbic (ILA), secondary motor (MOs) and ventral agranular 137 insula (AIv) areas) sustained a high firing rate during the delay period (Fig. 2B). Areas that 138 are higher in the cortical hierarchy show elevated activity during the delay period (Fig. 2C). 139 This persistent firing rate could last for more than 10 seconds and is a stable attractor state 140 of the network (Inagaki et al. 2019). 141

The cortical hierarchy and PV fraction predict the delay period firing rate of each cortical area (Fig. 2C-E). Thus the activity pattern of distributed working memory depends on both



Figure 2: Distributed working memory activity depends on the gradient of PV interneurons and the cortical hierarchy. (A). Model design of the large-scale model for distributed working memory. Top, connectivity map of the cortical network. Each node corresponds to a cortical area and an edge is a connection, where the thickness of the edge represents the strength of the connection. Only strong connections are shown (without directionality for the sake of clarity). Bottom, local and long-range circuit design. Each local circuit contains two excitatory populations (red), each selective to a particular stimulus and one inhibitory population (blue). Long-range connections are scaled by mesoscopic connectivity strength (Oh et al. 2014) and follows counterstream inhibitory bias (CIB) (Mejias and Wang 2022). (B). The activity of 6 selected areas during a working memory task is shown. A visual input of 500ms is applied to area VISp, which propagates to the rest of the large-scale network. (C). Delay period firing rate for each area on a 3d brain surface. Similar to Fig. 1B, the positions of 5 areas are labeled. (D). Delay-period firing rate is positively correlated with cortical hierarchy (r = 0.91, p < 0.05). (E). Delay-period firing rate is negatively correlated with PV cell fraction (r = -0.43, p < 0.05).

local and large-scale anatomy. The delay activity pattern has a stronger correlation with hierarchy (r = 0.91) than with the PV fraction (r = -0.43). The long-range connections thus play a predominant important role in defining the persistent activity pattern.

Activity in early sensory areas such as VISp displays a rigorous response to the transient input but returns to a low firing state after stimulus withdrawal. In contrast, many frontal areas show strong persistent activity. When the delay period firing rates are plotted versus hierarchy, we observe a gap in the distribution of persistent activity (Fig. 2D) that marks an abrupt transition in the cortical space. This leads to the emergence of a subnetwork of areas capable of working memory representations.

We also used our circuit model to simulate delayed response tasks with different sensory 153 modalities (Fig. 2 - supplement 1), by stimulating primary somatosensory area SSp-bfd and 154 primary auditory area AUDp. The pattern of delay period firing rates for these sensory 155 modalities is similar to the results obtained for the visual task: sensory areas show transient 156 activity, while frontal and lateral areas show persistent activity after stimulus withdrawal. 157 Moreover, the cortical hierarchy could predict the delay period firing rate of each cortical area 158 well (r = 0.89, p < 0.05), while the PV cell fraction could also predict the delay period firing 159 rate of each cortical area with a smaller correlation coefficient (r = -0.4, p < 0.05). Our model 160 thus predicts that working memory may share common activation patterns across sensory 161 modalities, which is partially supported by cortical recordings during a memory-guided 162 response task (Inagaki et al. 2018). 163

We explored the potential contributions of PV gradients and CIB in determining spatially-164 patterned activity across the cortex. To evaluate the importance of the PV gradient, we 165 replaced the PV gradient across areas with a constant value (Fig. 3A(ii)). As compared 166 to the model with a PV gradient (Fig. 3A(i)), we found that, during the delay period, the 167 number of cortical areas displaying persistent activity is diminished, but the abrupt transition 168 in delay period firing rates remains. This quantitative difference depends on the constant 169 value used to scale inhibition from PV cells across areas (Fig. 3 - supplement 1A, 1B). Next, 170 we performed the analogous manipulation on the CIB by scaling feedforward and feedback 171 projections with a constant value across areas, thus effectively removing the CIB. In this 172 case, the firing rate of both sensory and association areas exhibit high firing rates during the 173 delay period (Fig. 3A(iii)). Thus, CIB may be particularly important in determining which 174 areas exhibit persistent activity. 175

To further explore the model parameter space and better understand the interplay between 176 PV gradient and CIB, we systematically varied two critical model parameters: i) the base 177 local inhibitory weight $g_{EI,0}$ onto excitatory neurons, which sets the minimal inhibition for 178 each cortical area and ii) the scaling factor $g_{EI,scaling}$, which refers to how strongly the PV 179 gradient is reflected in the inhibitory weights. We created heatmaps that show the number of 180 areas with persistent activity during the delay period as a function of these parameters: in 181 Fig. 3C, we simulate the network with both CIB and PV gradient, while in Fig. 3D and 182 Fig. 3E we simulate networks when PV gradient or CIB is removed, respectively. In each of 183

these networks, we identify two regimes based on specific values for $g_{EI,0}$ and $g_{EI,scaling}$: a reference regime (used throughout the rest of the paper) and an alternative regime.

If we remove the PV gradient in the alternative parameter regime, persistent activity is 186 lost (Fig. 3B(ii)). In contrast, if we remove CIB the model still exhibits an abrupt transition 187 in firing rate activity (Fig. 3B(iii)). In this regime, a strong correlation and piece-wise linear 188 relationship between firing rate and PV cell fraction was uncovered that did not exist when 189 CIB was present. This observation led to a model prediction: if PV cell fraction is not 190 strongly correlated with delay firing rate across cortical areas (e.g., Fig. 3A(i) or Fig. 3B(i)), 191 this suggests the existence of a CIB mechanism at play. Importantly, the model without CIB 192 exhibits the abrupt transition in delay-period firing rates provided it is in a regime where 193 some areas exhibit 'independent' persistent activity: persistent activity that is generated 194 due to local recurrence and thus independent of long-range recurrent loops. The parameter 195 regime where some areas exhibit 'independent' persistent activity is quantified by varying the 196 base value of local inhibitory connections (Fig. 3 - supplement 1C). To conclude, the model 197 results suggest that CIB may be present in a large-scale brain network if the PV cell fraction 198 is not strongly correlated with the delay firing rate. Furthermore, CIB may be particularly 199 important in the regime where local connections are not sufficient to sustain independent 200 persistent activity. 201

Next, we evaluated the stability of the baseline state for the three conditions described 202 above: i) original with PV gradient and CIB, ii) after removal of CIB, and iii) after removal of 203 PV gradient. The heatmaps obtained after varying the base inhibitory strength and inhibitory 204 scaling factor were qualitatively the same across the three conditions, as shown by the blue 205 shaded squares in Fig. 3C-E. There are some regimes, such as the one depicted on the lower 206 left corner, where all areas exhibit persistent activity and there is no stable baseline: a regime 207 that is not biologically-realistic for a healthy brain. Thus, while PV and CIB shape the 208 distribution of delay firing rates across cortical areas, they don't qualitatively influence the 209 system's baseline stability. However, the inclusion of both PV gradient and CIB in the model 210 (Fig. 3C) results in a more robust system, i.e., a far wider set of parameters can produce 211 realistic persistent activity (Fig. 3D, 3E). 212

Local and long-range projections modulate the stability of the baseline state in the cortex

The stability of the baseline state for any given cortical area may have contributions from local inhibition or from long-range projections that target local inhibitory circuits. We found that individual local networks without long-range connections are stable without local inhibition (Fig. 4A, see methods for theoretical calculation of stability in a local circuit). However, in the full network with long-range connections, setting either the long-range connections to inhibitory neurons or local inhibition to zero made the network's baseline state unstable, and individual areas rose to a high firing rate (Fig. 4B). Thus, inhibition from local and ²²² long-range circuits contribute to the baseline stability of cortical areas.

Motivated by the results on stability, we investigated whether the large-scale network 223 model operates in the inhibitory stabilized network (ISN) regime (Tsodyks et al. 1997; Sanzeni 224 et al. 2020), whereby recurrent excitation is balanced by inhibition to maintain stability of 225 the baseline state. First, we examined whether individual brain areas (i.e., without long range 226 projections) may operate in this regime. We found a parameter set in which the baseline 227 firing rate is stable only when local inhibition is intact: when inhibition is removed, the 228 stable baseline state disappears, which suggests that the local circuits are ISNs (Fig. 4C and 229 see stability analysis in the Methods section). In the full neural network with long-range 230 connections, similar analysis as in Fig. 4B shows that the network becomes unstable if 231 long-range projections onto inhibitory interneurons are removed. (Fig. 4D). Thus we propose 232 that the network is also in a 'global' inhibitory stabilized network (ISN) regime, whereby 233 long-range connections to inhibitory neurons are necessary to maintain a stable baseline 234 state. Second, Second, we examined whether the ISN regime is consistent with distributed 235 working memory patterns in the cortex (Fig. 2). In the regime with increased local excitatory 236 connections but without long-range projections, some local circuits could reach a high stable 237 state when an external input is applied, demonstrating the bistability of those areas (Fig. 238 4E). When we considered the full network with long-range projections, the network exhibits a 239 graded firing rate pattern after transient stimulation of VISp, showing that the interconnected 240 ISN networks are compatible with bistability of a subset of cortical areas (Fig. 4F) 241

In summary, we have shown that distinct local and long-range inhibitory mechanisms shape the pattern of working memory activity and stability of the baseline state.

²⁴⁴ Thalamocortical interactions maintain distributed persistent activity

To investigate how thalamocortical interactions affect the large-scale network dynamics, 245 we designed a thalamocortical network similar to the cortical network (Fig. 5A). Several 246 studies have shown that thalamic areas are also involved in the maintenance of working 247 memory (Bolkan et al. 2017; Guo et al. 2017; Schmitt et al. 2017). However, the large-scale 248 thalamocortical mechanisms underlying memory maintenance are unknown. We set the 249 strength of connections between the thalamus and cortex using data from the Allen Institute 250 (Oh et al. 2014) (Fig. 5 - supplement 1). All thalamocortical connections in the model are 251 mediated by AMPA synapses. There are no recurrent connections in the thalamus within or 252 across thalamic nuclei (Jones 2007). The effect of thalamic reticular nucleus neurons was 253 included indirectly as a constant inhibitory current to all thalamic areas (Crabtree 2018; 254 Hádinger et al. 2023). Similarly to cortical areas, the thalamus is organized along a measured 255 hierarchy (Harris et al. 2019). For example, the dorsal part of the lateral geniculate nucleus 256 (LGd) is lower than the cortical area VISp in the hierarchy, consistent with the fact that 257 LGd sends feedforward inputs to VISp. Thalamocortical projections in the model are slightly 258 more biased toward excitatory neurons in the target area if they are feedforward projections 259



Figure 3: The role of PV inhibitory gradient and hierarchy-based counter inhibitory bias (CIB) in determining persistent activity patterns in the cortical network. (A(i)). Delay firing rate as a function of PV cell fraction with both CIB and PV gradient present (r = -0.42, p < 0.05). This figure panel is the same as Fig. 2E. (A(i)). Delay firing rate as a function of hierarchy after removal of PV gradient (r = 0.85, p < 0.05). (A(iii)). Delay firing rate as a function of PV cell fraction after removal of CIB (r = -0.74, p < 0.05). (B(i)). Delay firing rate as a function of PV cell fraction with both CIB and PV gradient present, in the alternative regime (r = -0.7, p < 0.05). (B(ii)). Delay firing rate as a function of hierarchy after removal of PV gradient, in the alternative regime (r =0.95, p < 0.05). (B(iii)). Delay firing rate as a function of PV cell fraction after removal of CIB, in the alternative regime (r = -0.84, p < 0.05). (C)-(E). Number of areas showing persistent activity (color coded) as a function of the local inhibitory gradient ($g_{EI,scaling}$, X axis) and the base value of the local inhibitory gradient $(g_{EI,0}, Y \text{ axis})$ for the following scenarios: (C) CIB and PV gradient, (D) with PV gradient replaced by a constant value, and (E) with CIB replaced by a constant value. The reference regime is located at the top left corner of the heatmap (green dot) and corresponds to A(i)-A(iii), while the alternative regime is located at the lower right corner (purple dot) and corresponds to B(i)-B(iii). The yellow dashed lines separate parameters sets for which none of the areas show 'independent' persistent activity (above the line) from parameter sets for which some areas are capable of maintaining persistent activity without input from other areas (below the line). Blue shaded squares in the heatmap mark the absence of a stable baseline.



Figure 4: Local and long-range projections modulate the baseline stability of individual cortical areas. Steady state firing rates are shown as a function of hierarchy for different scenarios: (A) without long-range connections in the reference regime ($g_{E,self} = 0.4nA$, $g_{EI,0} = 0.192nA$), (B) with long-range connections in the reference regime ($\mu_{EE} = 0.1nA$), (C) without long-range connections and increased local excitatory connections ($g_{E,self} = 0.6nA$, $g_{EI,0} = 0.5nA$), and (D) with long-range connections (increased long-range connections to excitatory neurons, $\mu_{EE} = 0.19nA$) and increased local excitatory connections to excitatory neurons, $\mu_{EE} = 0.19nA$) and increased local excitatory connections. (E) Firing rate as a function of hierarchy when external input given to each area, showing bistability for a subset of areas (parameters as in (C) and 'with local inh'). (F) Firing rate as a function of hierarchy when external input is applied to area VISp (parameters as in (D) and 'with local inh + with long-range E-I').

²⁶⁰ and towards inhibitory neurons if they are feedback.

Here, we weakened the strength of cortical interareal connections as compared to the cortex model of Fig. 2. Now, persistent activity can still be generated (Fig. 5B, blue) but is maintained with the help of the thalamocortical loop, as observed experimentally (Guo et al. 2017). Indeed, in simulations where the thalamus was inactivated, the cortical network no longer showed sustained activity (Fig. 5B, red).

In the thalamocortical model, the delay activity pattern of the cortical areas correlates 266 with the hierarchy, again with a gap in the firing rate separating the areas engaged in 267 persistent activity from those that do not (Fig. 5B, Fig. 5C). Sensory areas show a low 268 delay firing rate, and frontal areas show strong persistent firing. Unlike the cortex, the firing 269 rate of thalamic areas continuously increases along the hierarchy (Fig. 5E). On the other 270 hand, cortical dynamics in the thalamocortical and cortical models show many similarities. 271 Early sensory areas do not show persistent activity in either model. Many frontal and lateral 272 areas show persistent activity and there is an abrupt transition in cortical space in the 273 thalamocortical model, like in the cortex only model. Quantitatively, the delay firing pattern 274 of the cortical areas is correlated with the hierarchy and the PV fraction (Fig. 5C, Fig. 275 5D). Furthermore, the delay period firing rate of cortical areas in the thalamocortical model 276 correlates well with the firing rate of the same areas in the cortical model (Fig. 5F). This 277 comparison suggests that the cortical model captures most of the dynamical properties in 278 the thalamocortical model; therefore in the following analyses, we will mainly focus on the 279 cortex-only model for simplicity. 280

²⁸¹ Cell type-specific connectivity measures predict distributed persistent ²⁸² firing patterns

Structural connectivity constrains large-scale dynamics (Mejias and Wang 2022; Froudist-283 Walsh et al. 2021; Cabral et al. 2011). However, we found that standard graph theory 284 measures could not predict the pattern of delay period firing across areas. There is no 285 significant correlation between input strength and delay period firing rate (r = 0.25, p = 0.25, 286 Fig. 6A(i), A(ii)) and input strength cannot predict which areas show persistent activity 287 (prediction accuracy = 0.51, Fig. 6A(iii)). We hypothesized that this is because currently 288 available connectomic data used in this model do not specify the type of neurons targeted by 289 the long-range connections. For instance, when two areas are strongly connected with each 290 other, such a loop would contribute to the maintenance of persistent activity if projections are 291 mutually excitatory, but not if one of the two projections predominantly targets inhibitory 292 PV cells. Therefore, cell type-specificity of interareal connections must be taken into account 293 in order to relate the connectome with the whole-brain dynamics and function. To examine 294 this possibility, we introduced a cell type projection coefficient (see Calculation of network 295 structure measures in the Methods), which is smaller with a higher PV cell fraction in the 296 target area (Fig. 6 - supplement 1). The cell type projection coefficient also takes cell 297



Figure 5: Thalamocortical interactions help maintain distributed persistent activity. (A). Model schematic of the thalamocortical network. The structure of the cortical component is the same as our default model in Fig. 2A, but with modified parameters. Each thalamic area includes two excitatory populations (red square) selective to different stimuli. Long range projections between thalamus and cortex also follow the counterstream inhibitory bias rule as in the cortex. Feedforward projections target excitatory neurons with stronger connections and inhibitory neurons with weaker connections; the opposite holds for feedback projections. (B). The activity of 6 sample cortical areas in a working memory task is shown during control (blue) and when thalamic areas are inhibited in the delay period (red). Black dashes represent the external stimulus applied to VISp. Red dashes represent external inhibitory input given to all thalamic areas. (C). Delay period firing rate of cortical areas in the thalamocortical network. The activity pattern has a positive correlation with cortical hierarchy (r = 0.78, p < 0.05). (D). Same as (C) but plotted against PV cell fraction. The activity pattern has a negative correlation with PV cell fraction, but it is not significant (r = -0.26, p = 0.09). (E). Delay firing rate of thalamic areas in thalamocortical network. The firing rate has a positive correlation with thalamic hierarchy (r = 0.94, p < 0.05). (F). Delay period firing rate of cortical areas in thalamocortical network has a positive correlation with delay firing rate of the same areas in a cortex-only model (r = 0.77, p < 0.05). Note that only the areas showing persistent activity in both models are considered for correlation analyses.

type targets of long range connections into account, which, in our model, is quantified by counterstream inhibitory bias (CIB). As a result, the modified cell type-specific connectivity measures increase if the target area has a low density of PV interneurons and/or if long-range connections predominantly target excitatory neurons in the target area.

We found that cell type-specific graph measures accurately predict delay-period firing rates. The cell type-specific input strength of the early sensory areas is weaker than the raw input strength (Fig. 6B(i)). The firing rate across areas is positively correlated with cell type-specific input strength (Fig. 6B(ii)). Cell type-specific input strength also accurately predicts which areas show persistent activity (Fig. 6B(iii)). Similarly, we found that the cell type-specific eigenvector centrality, but not standard eigenvector centrality (Newman 2018), was a good predictor of delay period firing rates (Fig. 6 - supplement 2).

³⁰⁹ A core subnetwork for persistent activity across the cortex

Many areas show persistent activity in our model. However, are all active areas equally important in maintaining persistent activity? When interpreting large-scale brain activity, we must distinguish different types of contribution to working memory. For instance, inactivation of an area like VISp impairs performance of a delay-dependent task because it is essential for a (visual) "input" to access working memory; on the other hand a "readout" area may display persistent activity only as a result of sustained inputs from other areas that form a "core", which are causally important for maintaining a memory representation.

We propose four types of areas related to distributed working memory: input, core, 317 readout, and nonessential (Fig. 7A). External stimuli first reach input areas, which then 318 propagate activity to the core and non-essential areas. Core areas form recurrent loops and 319 support distributed persistent activity across the network. By definition, disrupting any of 320 the core areas would affect persistent activity globally. The readout areas also show persistent 321 activity. Yet, inhibiting readout areas has little effect on persistent activity elsewhere in 322 the network. We can assign the areas to the four classes based on three properties: a) the 323 effect of inhibiting the area during stimulus presentation on delay activity in the rest of the 324 network; b) the effect of inhibiting the area during the delay period on delay activity in the 325 rest of the network; c) the delay activity of the area itself on trials without inhibition. 326

In search of a core working memory subnetwork in the mouse cortex, in model simulations 327 we inactivated each area either during stimulus presentation or during the delay period, akin 328 to optogenetic inactivation in mice experiments. The effect of inactivation was quantified 329 by calculating the decrement in the firing rate compared to control trials for the areas that 330 were not inhibited (Fig. 7B). The VISp showed a strong inhibition effect during the stimulus 331 period, as expected for an Input area. We identified seven areas with a substantial inhibition 332 effect during the delay period (Fig. 7C), which we identify as a core for working memory. 333 Core areas are distributed across the cortex. They include frontal areas PL, ILA, medial part 334 of the orbital area (ORBm), which are known to contribute to working memory (Liu et al. 335



Figure 6: Cell type-specific connectivity measures are better at predicting firing rate pattern than nonspecific ones. (A(i)). Delay period firing rate (orange) and input strength for each cortical area. Input strength of each area is the sum of connectivity weights of incoming projections. Areas are plotted as a function of their hierarchical positions. Delay period firing rate and input strength are normalized for better comparison. (A(ii)). Input strength does not show significant correlation with delay period firing rate for areas showing persistent activity in the model (r = 0.25, p = 0.25). (A(iii)). Input strength cannot be used to predict whether an area shows persistent activity or not (prediction accuracy = 0.51). (B(i)). Delay period firing rate (orange) and cell type-specific input strength for each cortical area. Cell type-specific input strength considers how the long-rang connections target different cell types and is the sum of modulated connectivity weights of incoming projections. Same as (A(i)), areas are sorted according to their hierarchy and delay period firing rate and input strength are normalized for better comparison. (B(ii)). Cell type-specific input strength has a strong correlation with delay period firing rate of cortical areas showing persistent activity (r = 0.89, p < 0.05). Inset: Comparison of the correlation coefficient for raw input strength and cell type-specific input strength. (B(iii)). Cell type-specific input strength predicts whether an area shows persistent activity or not (prediction accuracy = 0.95). Inset: comparison of the prediction accuracy for raw input strength and cell type-specific input strength.



Figure 7: A core subnetwork generates persistent activity across the cortex. (A). We propose four different types of areas. Input areas (red) are responsible for coding and propagating external signals, which are then propagated through synaptic connections. Core areas (blue) form strong recurrent loops and generate persistent activity. Readout areas (green) inherit persistent activity from core areas. Nonessential areas (purple) may receive inputs and send outputs but they do not affect the generation of persistent activity. (B). Delay period firing rate for cortical areas engaged in working memory (Y axis) after inhibiting different cortical areas during the delay period (X axis). Areas in the X axis and Y axis are both sorted according to hierarchy. Firing rates of areas with small firing rate (<1Hz) are partially shown (only RSPv and RSPd are shown because their hierarchical positions are close to areas showing persistent activity). (C). The average firing rate for areas engaged in persistent activity under each inhibition simulation. The X axis shows which area is inhibited during the delay period, and the Y axis shows the average delay period activity for all areas showing persistent activity. Note that when calculating the average firing rate, the inactivated area was excluded in order to focus on the inhibition effect of one area on other areas. Average firing rates on the Y axis are normalized using the average firing in control (no inhibition) simulation. (D). Classification of 4 types of areas based on their delay period activity after stimulus- and delay-period inhibition (color denotes the type for area, as in A). The inhibition effect, due to either stimulus or delay period inhibition, is the change of average firing rate normalized by the average firing rate in the control condition. Areas with strong inhibition effect during stimulus period are classified as Input areas; areas with strong inhibition effect during delay period and strong delay period firing rate are classified as Core areas; areas with weak inhibition effect during delay period but strong delay period firing rate during control are classified as Readout areas; areas with weak inhibition effect during delay period and weak delay period firing rate during control are classified as Nonessential areas.

³³⁶ 2014; Bolkan et al. 2017). Other associative and sensory areas (AId, VISpm, ectorhinal area ³³⁷ (ECT), gustatory area (GU)) are also in the core. Similarly, we used the above criteria to ³³⁸ classify areas as Readout or Non-essential (Fig. 7D).

We have defined a core area for working memory maintenance as a cortical area that, first, exhibits persistent activity, and second, removal of this area (e.g., experimentally via a lesion or opto-inhibition) significantly affects persistent activity in other areas. It is possible, however, that effects on persistent activity at the network level only arise after lesioning two or more areas. Thus, we proceeded with inhibiting two, three, and four readout areas concurrently (Fig. 7- supplement 1A), as by definition, inhibiting any single readout area will not exhibit a strong inhibition effect.

We first inhibited pairs of readout areas and evaluated the effect of this manipulation at a 346 network level. Specifically, for any given readout area A, we plotted the average firing rate of 347 the network when A was inhibited as part of an inhibited pair (see description in Methods). 348 After inhibiting a pair of readout areas, there was a decrement in the average firing rate of the 349 network (Fig. 7 - supplement 1A). The decrement became more pronounced as more readout 350 areas were inhibited, e.g., triplets and quadruplets, and when a combination of readout and 351 core areas were inhibited pairwise (Fig. 7 - supplement 1B). This analysis demonstrates 352 that readout areas also play a role in maintaining distributed persistent activity: we may 353 define 'second-order core areas' as those readout areas that have a strong inhibition effect 354 only when inhibited concurrently with another area, while third-order and fourth-order core 355 areas are analogously defined via triplet and quadruplet inhibition, respectively. We note 356 that the effects of silencing pairs, triplets and quadruplets of readout areas remain smaller 357 than those seen after silencing single core areas listed above. We also tested the effect of 358 inhibiting all core areas during the delay period (Fig. 7 - supplement 1C). After inhibiting all 359 core areas, some readout areas lost persistent firing. Moreover, there was a 48% decrement 360 in the average firing rate compared with a 15% decrement for a single core area and a 3%361 decrement for a single readout area. Thus, the pattern of persistent activity is more sensitive 362 to perturbations of core areas, which underscores the classification of some cortical areas into 363 core vs readout.

The core subnetwork can be identified by the presence of strong excitatory loops

Inhibition protocols across many areas are computationally costly. We sought a structural indicator that is easy to compute and is predictive of whether an area is engaged in working memory function. Such an indicator could also guide the interpretation of large-scale neural recordings in experimental studies. In the dynamical regime where individual cortical areas do not show persistent activity independently, distributed working memory patterns must be a result of long-range recurrent loops across areas. We thus introduced a quantitative measurement of the degree to which each area is involved in long-range recurrent loops (Fig. 374 8A).

The core subnetwork can be identified by the presence of strong loops between excitatory 375 cells. Here we focus on length-2 loops (Fig. 8A); the strength of a loop is the product of 376 two connection weights for a reciprocally connected pair of areas; and the loop strength 377 measure of an area is the sum of the loop strengths of all length-2 loops that the area is part 378 of. Results were similar for longer loops (Fig. 8B, also see Fig. 8 - supplement 1 for results 379 of longer loops). The raw loop strength had no positive statistical relationship to the core 380 working memory subnetwork (Fig. 8C(i), Fig. 8C(ii)). We then defined cell type-specific loop 381 strength (see Methods). The cell type-specific loop strength is the loop strength calculated 382 using connectivity multiplied by the cell type projection coefficient. The cell type-specific loop 383 strength, but not the raw loop strength, predicts which area is a core area with high accuracy 384 (Fig. 8D(i), Fig. 8D(ii), prediction accuracy = 0.93). This demonstrates that traditional 385 connectivity measures are informative but not sufficient to explain dynamics during cognition 386 in the mouse brain. Cell type-specific connectivity, and new metrics that account for such 387 connectivity, are necessary to infer the role of brain areas in supporting large-scale brain 388 dynamics during cognition. 389

To better demonstrate our cell type-specific connectivity measures, we have implemented 390 two other measures for comparison: a) a loop-strength measure that adds a 'sign' without 391 further modification, and b) a loop strength measure that takes hierarchical information - and 392 not PV information- into account. These two graph-theoretic measures can be used to predict 393 delay firing rate during a sensory working memory task, thus highlighting the importance of 394 hierarchical information, which distinguishes excitatory from inhibitory feedback (Fig. 6 -395 supplement 3). On the other hand, the prediction of the core areas greatly depends on cell-type 396 specificity: the sign-only and 'no-PV' mechanisms do not reliably predict whether an area is 397 a core area or not, especially in the case of calculating with length 3 loops, demonstrating 398 the importance of cell-type specific connectivity measures. (Fig. 8 - supplement 2). 399

⁴⁰⁰ Multiple attractor states emerge from the mouse mesoscopic connec-⁴⁰¹ tome and local recurrent interactions

Different tasks lead to dissociable patterns of internally sustained activity across the brain, 402 described dynamically as distinct attractor states. Generally, attractor states may enable 403 computations such as decision making and working memory (Wang 1999; Wang 2002; Mejias et 404 al. 2016). Specifically, a given task may be characterized by a specific attractor landscape and 405 thereby define different core areas for working memory, as introduced above. We developed a 406 protocol to identify multiple attractor states, then analyzed the relationship between network 407 properties and the attractor states (Fig. 9A-C). For different parameters, the number of 408 attractors and the attractor patterns change. Two parameters are especially relevant here. 409 These are the long-range connection strength (μ_{EE}) and local excitatory connection strength 410 $(g_{E,self})$. These parameters affect the number of attractors in a model of the macaque cortex 411



Figure 8: The core subnetwork can be identified structurally by the presence of strong excitatory loops. (A). Distribution of length-2 loops. X axis is the single loop strength of each loop (product of connectivity strengths within loop) and Y axis is their relative frequency. (B). Loop strengths of each area calculated using different length of loops (e.g., length 3 vs length 2) are highly correlated (r = 0.96, p < 0.05). (C(i)). Loop strength (blue) is plotted alongside Core Areas (orange), a binary variable that takes the value 1 if the area is a Core Area, 0 otherwise. Areas are sorted according to their hierarchy. The loop strength is normalized to a range of (0, 1) for better comparison. (C(ii)). A high loop strength value does not imply that an area is a core area. Blue curve shows the logistic regression curve fits to differentiate the core areas versus non core areas. (D(i)). Same as (C), but for cell type-specific loop strength. (D(ii)). A high cell type-specific loop strength predicts that an area is a core area (prediction accuracy = 0.93). Same as (C), but for cell type-specific loop strength.

(Mejias and Wang 2022). Increasing the long range connection strength decreases the number 412 of attractors (Fig. 9D). Stronger long-range connections implies that the coupling between 413 areas is stronger. If areas are coupled with each other, the activity state of an area will be 414 highly correlated to that of its neighbors. This leads to less variability and fewer attractors. 415 To quantify how the patterns of attractors change for different parameters, two quantities 416 are introduced. The attractor fraction is the fraction of all detected attractor states to which 417 an area belongs. An area "belongs" to an attractor state if it is in a high activity state in that 418 attractor. The *attractor size* is defined by the number of areas belonging to that attractor. 419 As we increased the long-range connection strength, the attractor size distribution became bimodal. The first mode corresponded to large attractors, with many areas. The second 421 mode corresponded to small attractors, with few areas (Fig. 9D). 422

When the local excitatory strength is increased, the number of attractors increased as 423 well (Fig. 9E). In this regime some areas are endowed with sufficient local reverberation to 424 sustain persistent activity even when decoupled from the rest of the system, therefore the 425 importance of long-range coupling is diminished and a greater variety of attractor states is 426 enabled. This can be understood by a simple example of two areas 1 and 2, each capable of 427 two stimulus-selective persistent activity states; even without coupling there are $2 \times 2 = 4$ 428 attractor states with elevated firing. Thus, local and long-range connection strength have 429 opposite effects on the number of attractors. 430



Figure 9: Multiple attractors coexist in the mouse working memory network. (A-C) Example attractor patterns with a fixed parameter set. Each attractor pattern can be reached via different external input patterns applied to the brain network. Delay activity is shown on a 3D brain surface. Color represents the firing rate of each area. (D-E) The distribution of attractor fractions (left) and number of attractors as a function of size (right) for different parameter combinations are shown. Attractor fraction of an area is the ratio between the number of attractors that include the area and the total number of identified attractors. In (D), local excitatory strengths are fixed $(g_{E,self} = 0.44 \text{ nA})$ while long-range connection strengths vary in the range $\mu_{EE} = 0.01-0.05$ nA. Left and right panels of (D) show one specific parameter $\mu_{EE} = 0.03$ nA. Inset panel of (D) shows the number of attractors under different long-range connection strengths while $g_{E,self}$ is fixed at 0.44 nA. In (E), long range connection strengths are fixed ($\mu_{EE} = 0.02 \text{ nA}$) while local excitatory strengths varies in the range $g_{E,self} = 0.4-0.44$ nA. Left and right panels of (E) show one specific parameter $g_{E,self} = 0.43$ nA. Inset panel of (E) shows the number of attractors under different local excitatory strengths, while μ_{EE} is fixed at 0.02 nA. (F). Prediction of the delay period firing rate using input strength and cell type-specific input strength for each attractor state identified under $\mu_{EE} = 0.04$ nA and $g_{E,self}$ = 0.44 nA. 143 distinct attractors were identified and the average correlation coefficient using cell type-specific input strength is better than that using input strength. (G). A example attractor state identified under the parameter regime $\mu_{EE} = 0.03$ nA and $g_{E,self} = 0.44$ nA. The 5 areas with persistent activity are shown in red. (H). Effect of single area inhibition analysis for the attractor state in (G). For a regime where 5 areas exhibit persistent activity during the delay period, inactivation of the premotor area MOs yields a strong inhibition effect (<0.95 orange dashed line) and is therefore a Core area for the attractor state in (G). (I). Cell type-specific loop strength (blue) is plotted alongside Core areas (orange) for the attractor state in (G). Only 5 areas with persistent activity are used to calculate the loop strength. Loop strength is normalized to be within the range of 0 and 1. High cell type-specific loop measures predict that an area is a Core area (prediction accuracy is 100% correct). The number of areas is limited, so prediction accuracy is very high.

431 The cell type-specific input strength predicted firing rates across many attractors. In an

example parameter regime ($\mu_{EE} = 0.04$ nA and $g_{E,self} = 0.44$ nA), we identified 143 attractors. We correlated the input strength and cell type-specific input strength with the many attractor firing rates (Fig. 9F). The raw input strength is weakly correlated with activity patterns. The cell type-specific input strength is strongly correlated with activity across attractors. This

⁴³⁷ in many scenarios. These results further prove the importance of having cell type-specific

shows that the cell type-specific connectivity measures are better at predicting the firing rates

438 connectivity for modeling brain dynamics.

436

Different attractor states rely on distinct subsets of core areas. In one example attractor, we 439 found 5 areas that show persistent activity: VISa, VISam, FRP, MOs and ACAd (Fig. 9G) 440 (parameter regime, $\mu_{EE} = 0.03$ nA and $g_{E,self} = 0.44$ nA). We repeated the previous inhibition 441 analysis to identify core areas for this attractor state. Inhibiting one area, MOs, during the 442 delay had the strongest effect on delay activity in the other parts of the attractor (Fig. 9H). 443 MOs also showed strong persistent activity during delay period. This is consistent with its 444 role in short-term memory and planning (Li et al. 2015; Inagaki et al. 2019). According to 445 our definition, MOs is a core area for this attractor. To calculate a loop strength that was 446 specific to this attractor, we only examined connections between these five areas. The cell 447 type-specific loop strength was strongest in area MOs (Fig. 9I). Thus, we can identify likely 448 core areas for individual attractor states from cell type-specific structural measures. This 449 also demonstrates that different attractor states can be supported by distinct core areas. 450

451 Discussion

We developed a connectome-based dynamical model of the mouse brain. The model was capable of internally maintaining sensory information across many brain areas in distributed activity in the absence of any input. To our knowledge this is the first biologically-based model of the entire mouse cortex and the thalamocortical system that supports a cognitive function, in this case working memory. Together with our recent work (Mejias and Wang 2022; Froudist-Walsh et al. 2021; Froudist-Walsh et al. 2023), it provides an important reference point to study the differences between mice and monkeys.

Our main findings are threefold. First, the mnemonic activity pattern is shaped by the 459 differing densities of PV interneurons across cortical areas. Areas with a high PV cell fraction 460 encoded information only transiently. Those with low PV cell fraction sustained activity 461 for longer periods. Thus, the gradient of PV cells (Kim et al. 2017) has a definitive role in 462 separating rapid information processing in sensory areas from sustained mnemonic information 463 representation in associative areas of the mouse cortex. This is consistent with the view that 464 each local area operates in the "inhibition-stabilizing regime" where recurrent excitation alone 465 would lead to instability but the local network is stabilized by feedback inhibition, which may 466 arise from long-range excitatory inputs to inhibitory neurons. This consistent with the regime 467 of the primary visual cortex (R. J. Douglas et al. 1995; Murphy and Miller 2009). Second, we 468 deliberately considered two different dynamical regimes: when local recurrent excitation is not 469 sufficient to sustain persistent activity and when it is. In the former case, distributed working 470 memory must emerge from long-range interactions between parcellated areas. Thereby the 471 concept of synaptic reverberation (Lorente de Nó 1933; P. S. Goldman-Rakic 1995; Wang 472 2001; Wang 2021) is extended to the large-scale global brain. Note that currently it is unclear 473 whether persistent neural firing observed in a delay dependent task is generated locally or 474 depends on long-distance reverberation among multiple brain regions. Our work made the 475 distinction explicit and offers specific predictions to be tested experimentally. Third, presently 476 available connectomic data are not sufficient to account for neural dynamics and distributed 477 cognition, and we propose cell type-specific connectomic measures that are shown to predict 478 the observed distributed working memory representations. Our model underscores that, 479 although connectome databases are an invaluable resource for basic neuroscience, they should 480 be supplemented with cell-type-specific information. 481

We found that recurrent loops within the cortex and the thalamocortical network aided in 482 sustaining activity throughout the delay period (Guo et al. 2017; Schmitt et al. 2017). The 483 presence of thalamocortical connections had a similar effect on the model as cortico-cortical 484 projections, with the distinct contributions of the thalamus to large-scale dynamics still to be 485 uncovered (Shine et al. 2018; Jaramillo et al. 2019). The specific pattern of cortico-cortical 486 connections was also critical to working memory. However, standard graph theory measures 487 based on the connectome were unable to predict the pattern of working memory activity. 488 By focusing on cell type-specific interactions between areas, we were able to reveal a core of 489

cortical areas. The core is connected by excitatory loops, and is responsible for generating
a widely distributed pattern of sustained activity. This clarifies the synergistic roles of the
connectome and gradients of local circuit properties in producing a distributed cognitive
function.

Previous large-scale models of the human and macaque cortex have replicated functional 494 connectivity (Deco et al. 2014; Demirtaş et al. 2019; Honey et al. 2007; Schmidt et al. 2018; 495 Shine et al. 2018; Cabral et al. 2011; Wang et al. 2019) and propagation of information 496 along the cortical hierarchy (Chaudhuri et al. 2015; Joglekar et al. 2018; Diesmann et al. 497 1999). More recently, large-scale neural circuit models have been developed specifically to reproduce neural activity during cognitive tasks (Mejias and Wang 2022; Froudist-Walsh 499 et al. 2021; Klatzmann et al. 2022). These models consider the fact that in the macaque 500 cortex, the density of spines on pyramidal cells increases along the cortical hierarchy (Elston 501 and Rosa 1998; Elston 2007; Chaudhuri et al. 2015). In a large-scale model of the macaque 502 cortex (Chaudhuri et al. 2015), it was shown that this 'excitatory gradient' (Wang 2020) 503 is correlated with the distribution of intrinsic timescales in the cortex (Murray et al. 2014) 504 and is consistent with spatially distributed working memory patterns (Mejias and Wang 505 2022; Froudist-Walsh et al. 2021). Such excitatory gradients based on spine count are less 506 pronounced, and may be entirely absent in the rodent cortex (Ballesteros-Yáñez et al. 2010; 507 Gilman et al. 2017). However, there are gradients of synaptic inhibition in the mouse cortex 508 (Kim et al. 2017; Wang 2021). Kim et al., showed that the ratio of SST+ neurons to PV+ 509 neurons is low for early sensory areas and motor areas, while it is high in association areas 510 such as the frontal cortex. We have used this gradient of inhibition in our model to show 511 that spatially distributed persistent-activity patterns in the mouse cortex do not require 512 gradients of recurrent excitation. In our model, the PV gradient and CIB may be particularly 513 important to maintain the stability of an otherwise highly excitable cortical area. Along these 514 lines, we predict that local recurrency in the mouse early sensory areas is higher than in the 515 primate. Consistent with this claim, both the spine density and the number of excitatory and 516 inhibitory synapses in layer 2/3 pyramidal neurons in area V1 are higher in mouse compared 517 to macaque (Fig. 5A in Gilman et al. 2017, Fig. 1A in Wildenberg et al. 2021). 518

Other anatomical properties at the area and single cell level may be informative of the 519 differences in computational and/or cognitive abilities between rodents and macaques. In the 520 language of network theory, the macaque cortex is a densely connected graph at an inter-area 521 level, with the connectivity spanning five orders of magnitude (Markov et al. 2014a), which is 522 more than what is expected for small-world networks (Bassett and Bullmore 2017). Critically, 523 the mouse 'connectome' (e.g., (Oh et al. 2014; Harris et al. 2019; Knox et al. 2018) has 524 even denser area-to-area connections. In the visual cortex, individual neurons target more 525 cortical areas in the mouse (Siu et al. 2021) and they have more inhibitory and excitatory 526 synapses (Wildenberg et al. 2021). Thus, connectivity in the mouse is denser at both the 527 area and single-cell levels (at least for primary visual cortex). We propose that there is a 528 greater functional specialization in the primate cortex which is afforded by the sparser and 529

more targeted patterns of connectivity at the single-cell and area levels. Other differences to explore in future computational models include the ratio of NMDA to AMPA-mediated synaptic currents, which is approximately constant in the mouse cortex (Myme et al. 2003) but varies along the cortical hierarchy in primates (Yang et al. 2018; Klatzmann et al. 2022), as well as hierarchy, which is defined based on feedforward and feedback projections in the mouse (Harris et al. 2019) and primate (Markov et al. 2014a).

We found that traditional graph theory metrics of connectivity were unable to predict the 536 working memory activity in the mouse brain. This may be due to the almost fully connected 537 pattern of interareal connectivity in the mouse cortex (Gămănuţ et al. 2018). This implies 538 that, qualitatively, all areas have a similar set of cortical connections. In our model, we 539 allowed the cell type target of interareal connections to change according to the relative 540 position of the areas along the cortical hierarchy. Specifically, feedforward connections had 541 a greater net excitatory effect than feedback connections, a hypothesis which we refer to 542 as CIB. This preferential targeting of feedback projections serves to stabilize the otherwise 543 excitable activity of sensory areas (Mejias and Wang 2022), and is consistent with recent 544 experiments that report long-range recruitment of GABAergic neurons in early sensory areas 545 (Campagnola et al. 2022; Shen et al. 2022; Naskar et al. 2021). Our model predicts that 546 if there is a weak correlation between PV cell density and delay firing rate across cortical 547 areas, then the CIB mechanism is at play. Moreover, the model results suggest that CIB is 548 particularly important in the regime where local connections are not sufficient to sustain 549 spatially-patterned persistent activity. We also showed that there are parameter regimes 550 where CIB becomes less important, provided there is a gradient of synaptic inhibition as in the 551 mouse cortex ((Kim et al. 2017), but see (Nigro et al. 2022)). Notably, the model's resilience 552 to parameter variations in inhibitory connection strengths is significantly enhanced when 553 both the PV gradient and CIB are present. Given that working memory is a fundamental 554 cognitive function observed across many individual brains with anatomical differences, the 555 inclusion of multiple inhibitory mechanisms that allow for connectivity variations might confer 556 evolutionary advantages. Although there is some evidence for similar inhibitory gradients 557 in humans (Burt et al. 2018) and macaque (Torres-Gomez et al. 2020), the computational 558 consequences of differences across species remain to be established. 559

To conclude, the manner in which long-range recurrent interactions affect neural dynamics 560 depends not only on the existence of excitatory projections per se, but also on the target 561 neurons' cell type. Thus, for some cortical areas afferent long-range excitatory connections 562 promote working memory-related activity while for some others, e.g., early sensory areas, it 563 does not. Moreover, the existence of long-range interactions is consistent with potentially 564 distinct dynamical regimes. For example, in one regime some areas exhibit independent 565 persistent activity, i.e., local recurrent interactions are sufficient to sustain a memory state 566 for these areas, while others do not. In this regime CIB is not required for the existence of 567 distributed persistent activity patterns. In another regime, none of the areas can sustain a 568 memory state without receiving long-range input. These two regimes are functionally distinct 569

⁵⁷⁰ in terms of their robustness to perturbation as well as in the number of attractors that they ⁵⁷¹ can sustain. These regimes may be identified via perturbation analysis in future experimental ⁵⁷² and theoretical work.

By introducing cell type-specific graph theory metrics, we were able to predict the pattern 573 and strength of delay period activity with high accuracy. Moreover, we demonstrated how 574 cell type-specific graph-theory measures can accurately identify the core subnetwork, which 575 can also be identified independently using a simulated large-scale optogenetic experiment. We 576 found a core subnetwork of areas that, when inhibited, caused a substantial drop in activity 577 in the remaining cortical areas. This core working memory subnetwork included frontal 578 cortical areas with well documented patterns of sustained activity during working memory 579 tasks, such as prelimbic (PL), infralimbic (ILA) and medial orbitofrontal cortex (ORBm) 580 (Schmitt et al. 2017; Liu et al. 2014; Wu et al. 2020). However, the core subnetwork for the 581 visual working memory task we assessed was distributed across the cortex. It also included 582 temporal and higher visual areas, suggesting that long-range recurrent connections between 583 the frontal cortex and temporal and visual areas are responsible for generating persistent 584 activity and maintaining visual information in working memory in the mouse. 585

Some of the areas that were identified as core areas in our model have been widely studied in other tasks. For example, the gustatory area exhibits delay-period preparatory activity in a taste-guided decision-making task and inhibition of this area during the delay period impairs behavior (Vincis et al. 2020).

The core visual working memory subnetwork generates activity that is then inherited by 590 many readout areas, which also exhibit persistent activity. However, inhibiting readout areas 591 only mildly affects the activity of other areas (Fig. 7 and Fig. 7 - supplement 1). The 592 readout areas in our model were a mixture of higher visual areas, associative areas and 593 premotor areas of cortex. We also concluded that MOs is a readout area and not a core area. 594 This finding may be surprising considering previous studies that have shown this area to 595 be crucial for short-term memory maintenance, planning, and movement execution during 596 a memory-guided response task (Guo et al. 2017; Guo et al. 2014; Inagaki et al. 2019; Li 597 et al. 2015; Wu et al. 2020; Voitov and Mrsic-Flogel 2022). This task has shown to engage, 598 not only ALM, but a distributed subcortical-cortical network that includes the thalamus, 599 basal ganglia and cerebellum (Svoboda and Li 2018). We note that in the version of the 600 memory-guided response task studied by Svoboda and others, short-term memory is conflated 601 with movement preparation. In our task, we proposed to study the maintenance of sensory 602 information independent of any movement preparation as in delayed match-to-sample tasks 603 and variations thereof. It is for this behavioral context that we found that MOs is not a core 604 visual working memory area. We emphasize that readout areas are not less important than 605 core areas as readout areas can use the stored information for further computations and thus 606 some readout areas are expected to be strongly coupled to behavior. Indeed, there is evidence 607 for a differential engagement of cortical networks depending on the task design (Jonikaitis 608 et al. 2023) and on effectors (Kubanek and Snyder 2015). If ALM is indeed a readout area 609

for sensory working memory tasks, (e.g., (Schmitt et al. 2017)), then the following prediction 610 arises. Inhibiting ALM should have a relatively small effect on sustained activity in core 611 areas (such as PL) during the delay period. In contrast, inhibiting PL and other core areas 612 may disrupt sustained activity in ALM. Even if ALM is not part of the core for sensory 613 working memory, it could form part of the core for motor preparation tasks (Fig. 9G). We 614 found a high cell-type-specific loop strength for area ALM, like that in core areas, which 615 supports this possibility (Fig. 9I). Furthermore, we found some attractor states for which the 616 MOs was classified as a core area, that do not contain area PL. This result is supported by a 617 recent study that found no behavioral effect after PL inhibition in a motor planning task 618 (Wang et al. 2021). Therefore, the core subnetwork required for generating persistent activity 619 is likely task-dependent. Future modeling work may help elucidate the biological mechanisms 620 responsible for switching between attractor landscapes for different tasks. 621

Neuroscientists are now observing task-related neural activity at single-cell resolution across 622 much of the brain (Stringer et al. 2019; Steinmetz et al. 2019). This makes it important to 623 identify ways to distinguish the core areas for a function from those that display activity that 624 serves other purposes. We show that a large-scale inhibition protocol can identify the core 625 subnetwork for a particular task. We further show how this core can be predicted based on 626 the interareal loops that target excitatory neurons. Were such a cell type-specific interareal 627 connectivity dataset available, it may help interpretation of large-scale recording experiments. 628 This could also focus circuit manipulation on regions most likely to cause an effect on the 629 larger network activity and behavior. Our approach identifies the brain areas that work 630 together to support working memory. It also identifies those that benefit from such activity 631 to serve other purposes. Our simulation and theoretical approach is therefore ideally suited 632 to understand the large-scale anatomy, recording and manipulation experiments which are at 633 the forefront of modern systems neuroscience. 634

Neuroscience has rapidly moved into a new era of investigating large-scale brain circuits. 635 Technological advances have enabled the measurement of connections, cell types and neural 636 activity across the mouse brain. We developed a model of the mouse brain and theory of 637 working memory that is suitable for the large-scale era. Previous reports have emphasized the 638 importance of gradients of dendritic spine expression and interareal connections in sculpting 639 task activity in the primate brain (Mejias and Wang 2022; Froudist-Walsh et al. 2021). 640 Although these anatomical properties from the primate cortex are missing in the mouse brain 641 (Gămănuț et al. 2018; Gilman et al. 2017), other properties such as interneuron density (Kim 642 et al. 2017) may contribute to areal specialization. Indeed, our model clarifies how gradients 643 of interneurons and cell type-specific interactions define large-scale activity patterns in the 644 mouse brain during working memory, which enables sensory and associative areas to have 645 complementary contributions. Future versions of the large-scale model may consider different 646 interneuron types to understand their contributions to activity patterns in the cortex (Kim 647 et al. 2017; Meng et al. 2023; Froudist-Walsh et al. 2021; Wang et al. 2004; Tremblay et al. 648 2016; Nigro et al. 2022), the role of interhemispheric projections in providing robustness for 649

short-term memory encoding (Li et al. 2016a), and the inclusions of populations with tuning to various stimulus features and/or task parameters that would allow for switching across tasks (Yang et al. 2019). Importantly, these large-scale models may be used to study other important cognitive computations beyond working memory, including learning and decision making (Abbott et al. 2017; Abbott et al. 2020).

Acknowledgements

We thank Daniel P. Bliss and Ulises Pereira for support with analysis tools at the beginning of the project, and members of the Wang Lab at New York University for discussions related to the project.

Declaration of Interests

⁶⁶⁰ No competing interests declared.

661 Methods

⁶⁶² Anterograde tracing, connectivity data

We used the mouse connectivity map from Allen institute (Oh et al. 2014) to constrain our 663 large-scale circuit model of the mouse brain. The Allen Institute measured the connectivity 664 among cortical and subcortical areas using an anterograde tracing method. In short, they 665 injected virus and expressed fluorescent protein in source areas and performed fluorescent 666 imaging in target areas to measure the strength of projections from source areas. Unlike 667 retrograde tracing methods used in other studies (Markov et al. 2014b), the connectivity strength measured using this method does not need to be normalized by the total input or 669 output strength. This means that connectivity strength between any two areas is comparable. 670 The entries of the connectivity matrix from the Allen Institute can be interpreted as propor-671 tional to the total number of axonal fibers projecting from unit volume in one area to unit 672 volume in another area. Before incorporating the connectivity into our model, we normalized 673 the data as follows. In each area, we model the dynamics of an "average" neuron, assuming 674 that the neuron receive inputs from all connected areas. Thus, we multiplied the connectivity 675 matrix by the volume Vol_j of source area j and divided by the average neuron density d_i in 676 target area i: 677

$$W_{norm,ij} = W_{raw,ij} \frac{Vol_j}{d_i} \tag{1}$$

where $W_{raw,ij}$ is the raw, i.e., original, connection strength from unit volume in source area jto unit volume in target area i, Vol_j is the volume of source area j (Wang et al. 2020), and d_i is the neuron density in source area i (Erö et al. 2018). $W_{norm,ij}$ is the matrix that we use to set the long rang connectivity in our circuit model. We can define the cortico-thalamic connectivity $W_{ct,norm,ij}$ and thalamo-cortical connectivity $W_{tc,norm,ij}$ in a similar manner, except that we didn't apply the normalization to thalamic connectivity due to not having enough neuron density data.

⁶⁸⁵ Interneuron density along the cortex

Kim and colleagues measured the density of typical interneuron types in the brain (Kim et al. 2017). They expressed fluorescent proteins in genetically labeled interneurons and counted the number of interneurons using fluorescent imaging. We took advantage of these interneuron density data and specifically used the PV cell fraction to set local and long-range inhibitory weights.

The PV cell density of all layers is first divided by the total neuron density d_i in the area i, to give the PV cell fraction $PV_{raw,i}$, which better reflects the expected amount of synaptic inhibition mediated by PV neurons. The PV cell fraction is then normalized across the whole cortex.

$$PV_i = \frac{PV_{raw,i} - min(PV_{raw,i})}{max(PV_{raw,i}) - min(PV_{raw,i})}$$
(2)

⁶⁹⁵ $PV_{raw,i}$ is the PV cell fraction in area *i*, and PV_i is the normalized value of PV_{raw} , which ⁶⁹⁶ will be used in subsequent modeling.

⁶⁹⁷ Hierarchy in the cortex

The concept of hierarchy is important for understanding the cortex. Hierarchy can be 698 defined based on mapping corticocortical long range connections onto feedforward or feedback 699 connections (Felleman and Essen 1991; Markov et al. 2014a; Harris et al. 2019). Harris and 700 colleagues measured the corticocortical projections and target areas in a series of systematic 701 experiments in mice (Harris et al. 2019). Projection patterns were clustered into multiple 702 groups and the label "feedforward" or "feedback" was assigned to each group. Feedforward 703 and feedback projections were then used to determine relative hierarchy between areas. For 704 example, if the projections from area A to area B are mostly feedforward, then area B has 705 a higher hierarchy than area A. This optimization process leads to a quantification of the 706 relative hierarchy of cortical areas $h_{raw,i}$. We defined the normalized hierarchy value h_i as 707

$$h_i = \frac{h_{raw,i} - \min(h_{raw,i})}{\max(h_{raw,i}) - \min(h_{raw,i})}$$
(3)

where $h_{raw,i}$ is the raw, i.e., original hierarchical ordering from (Harris et al. 2019). Due to data acquisition issues, 6 areas did not have a hierarchy value assigned to them (SSp-un, AUDv, GU, VISC, ECT, PERI) (Harris et al. 2019). We estimated hierarchy through a weighted sum of the hierarchy value of 37 known areas, while the weight is determined through the connectivity strength. The parameters α_h and β_h are selected so that $h_{i,estimate}$ are close to h_i for areas with known hierarchy.

$$h_{i,estimate} = \alpha_h \frac{\sum_{j=1}^{37} W_{raw,ij} h_j}{\sum_{j=1}^{37} W_{raw,ij}} + \beta_h$$
(4)

For the thalamocortical model, we also used the hierarchy value for thalamic areas (Harris et al. 2019). The hierarchy of thalamic areas are comparable to cortical areas, so in order to use it in the model, we also normalized them.

$$h_{th,i} = \frac{h_{th,raw,i} - \min(h_{raw,i})}{\max(h_{raw,i}) - \min(h_{raw,i})}$$
(5)

⁷¹⁷ To estimate the hierarchy value of thalamic areas with missing values, we used the known⁷¹⁸ hierarchy value of the thalamic area next to the missing one as a replacement.

719 Description of the local circuit

Our large-scale circuit model includes 43 cortical areas. Each area includes two excitatory populations, labeled A and B, and one inhibitory population, C. The two excitatory populations are selective to different stimuli. The synaptic dynamics between populations are based on previous firing rate models of working memory (Wang 1999; Wong and Wang 2006). The equations that define the dynamics of the synaptic variables are

$$\frac{dS_A}{dt} = -\frac{S_A}{\tau_N} + \gamma (1 - S_A) r_A \tag{6}$$

$$\frac{dS_B}{dt} = -\frac{S_B}{\tau_N} + \gamma (1 - S_B) r_B \tag{7}$$

$$\frac{dS_C}{dt} = -\frac{S_C}{\tau_G} + \gamma_I r_C \tag{8}$$

where S_A and S_B are the NMDA synaptic variables of excitatory populations A and B, while S_C is the GABA synaptic variable of the inhibitory population C. r_A , r_B and r_C are the firing rates of populations A, B and C, respectively. τ_N and τ_G are the time constants of NMDA and GABA synaptic conductances. γ and γ_I are the parameters used to scale the contribution of presynaptic firing rates. The total currents received I_i (i = A, B, C) are given by

$$I_A = g_{E,self} S_A + g_{E,cross} S_B - g_{EI} S_C + I_{0A} + I_{LR,A} + x_A(t)$$
(9)

$$I_B = g_{E,self} S_B + g_{E,cross} S_A - g_{EI} S_C + I_{0B} + I_{LR,B} + x_B(t)$$
(10)

$$I_C = g_{IE}S_A + g_{IE}S_B - g_{II}S_C + I_{0C} + I_{LR,C} + x_C(t)$$
(11)

In these equations, $g_{E,self}$, $g_{E,cross}$ denote the connection strength between excitatory neurons with same or different selectivity, respectively. These connection strengths are the same for different areas, since there is no significant gradient for excitatory strength in mice. g_{IE} are the connection strengths from excitatory to inhibitory neurons, while g_{EI} , and g_{II} are connection strengths from inhibitory to excitatory neurons and from inhibitory to inhibitory neurons, respectively. These connections will be scaled by PV cell fraction PV_i in the corresponding area. We will discuss the details in the next section. I_{0i} (i = A, B, C) are constant background ⁷³³ currents to each population. $I_{LR,i}$ (i = A, B, C) are the long range (LR) currents received ⁷³⁴ by each population. The term $x_i(t)$ where i = A, B, C represents noisy contributions from ⁷³⁵ neurons external to the network. It is modeled as an Ornstein-Uhlenbeck process:

$$\tau_{noise} \frac{dx_i}{dt} = -x_i + \sqrt{\tau_{noise}} \sigma_i \zeta_i(t) \tag{12}$$

where $\zeta_i(t)$ is Gaussian white noise, τ_{noise} describes the time constant of external AMPA synapses and σ_i sets the strength of the noise for each population. $\sigma_A = \sigma_B = 5pA$ while $\sigma_C = 0pA$.

The steady state firing rate of each population is calculated based on a transfer function $\phi_i(I)$ of input current received by each population I_i (i = A, B, C) given by

$$\phi_{A,B}(I_{A,B}) = \frac{aI_{A,B} - b}{1 - exp[-d(aI_{A,B} - b)]}$$
(13)

$$\phi_C(I_C) = \left[\frac{1}{g_I}(c_1 I - c_0) + r_0\right]^+ \tag{14}$$

Note that the transfer functions $\phi_i(t)$ are the same for two excitatory populations. x^+ denotes the positive part of the function x. The firing rate of each population follows equations:

$$\tau_r \frac{dr_{A,B}}{dt} = -r_{A,B} + \phi_{A,B}(I_{A,B})$$
(15)

$$\tau_r \frac{dr_C}{dt} = -r_C + \phi_C(I_C) \tag{16}$$

⁷³⁹ Interneuron gradient and local connections

We scaled local interneuron connectivity with the interneuron density that was obtained using fluorescent labeling (Kim et al. 2017). Specifically, local I-I connections and local I-E connections are scaled by the interneuron density by setting the connection strength $g_{k,i}(k = EI, II)$ as a linear function of PV cell fraction PV_i in area *i*.

$$g_{EI,i} = g_{EI,0}(1 + g_{EI,scaling}PV_i) \tag{17}$$

$$g_{II,i} = g_{II,0}(1 + g_{II,scaling}PV_i) \tag{18}$$

where $g_{k,0}$ (k = EI, II) is the base value of I to E connections and $g_{k,scaling}$ (k = EI, II) is the scaling factor of PV value. $g_{k,0}$ also accounts for the inhibition of other cell types not explicitly considered in this study.

⁷⁴³ Hierarchy and long range connections

Long range (LR) connections between areas are scaled by connectivity data from the Allen Institute (Oh et al. 2014). We consider long-range connections that arise from excitatory neurons because most long-range connections in the cortex correspond to excitatory connections (Petreanu et al. 2009). Long-range connections will target excitatory populations in other brain areas with the same selectivity (Zandvakili and Kohn 2015) and will also target inhibitory neurons. These long-range connections are given by the following equations:

$$I_{A,B,LR,i} = \Sigma_j \mu_{EE} W_{E,ij} S_{A,B,j} \tag{19}$$

$$I_{C,LR,i} = \Sigma_j \mu_{IE} W_{I,ij} (S_{A,j} + S_{B,j})$$

$$\tag{20}$$

where W_E is the normalized long-range connectivity to excitatory neurons, and W_I is the normalized long-range connectivity to inhibitory neurons. μ_{EE} and μ_{IE} are coefficients scaling the long-range E to E and E to I connection strengths, respectively.

Here, we assume that the long-range connections will be scaled by a coefficient that is based 747 on the hierarchy of source and target area. To quantify the difference between long-range 748 feedforward and feedback projections, we introduce m_{ij} to measure the "feedforwardness" of projections between two areas. According to our assumption of counterstream inhibitory bias 750 (CIB), long-range connections to inhibitory neurons are stronger for feedback connections 751 and weaker for feedforward connections, while the opposite holds for long range connections 752 to excitatory neurons. Following this hypothesis, we define m_{ij} as a sigmoid function of 753 the difference between the hierarchy value of source and target areas. For feedforward 754 projections, $m_{ij} > 0.5$; for feedback projections, $m_{ij} < 0.5$. Excitatory and inhibitory 755 long-range connection strengths are implemented by multiplying the long-range connectivity 756 strength W_{ij} by m_{ij} and $(1 - m_{ij})$, respectively: 757

$$m_{ij} = \frac{1}{1 + e^{-\beta(h_i - h_j)}} \tag{21}$$

$$W_{E,ij} = m_{ij}W_{ij} \tag{22}$$

$$W_{I,ij} = (1 - m_{ij})W_{ij}$$
(23)

with

$$W_{scale,ij} = (W_{norm,ij})^{k_{scale}}$$
(24)

$$W_{ij} = \frac{W_{scale,ij}}{max(W_{scale,ij})} \tag{25}$$

The connectivity $W_{norm,ij}$ is then rescaled to translate the broad range of connectivity values (over five orders of magnitude) to a range more suitable for our firing rate models. k_{scale} is the coefficient used for this scaling. $k_{scale} < 1$ effectively makes the range much smaller than the original normalized connectivity $W_{norm,ij}$. After that, the scaled connectivity $W_{scale,ij}$ is then normalized so that the maximum value is fixed at 1.

⁷⁶³ Simulations of replacing the PV gradient and CIB

In order to demonstrate the importance of PV gradient and CIB, we replace the PV gradient value/CIB with the average value accordingly in the simulation. Specifically, we replace PV

gradient with the average PV cell fraction.

$$PV_{mean} = \frac{\sum_{i} PV_i}{n_{areas}} \tag{26}$$

$$g_{EI,i} = g_{EI,0}(1 + g_{EI,scaling}PV_{mean}) \tag{27}$$

$$g_{II,i} = g_{II,0}(1 + g_{II,scaling}PV_{mean})$$

$$\tag{28}$$

We also replace CIB with its average value 0.5, which means there is no bias to inhibitory cells for all long range connections.

$$m_{ij} = 0.5$$
 (29)

$$W_{E,ij} = 0.5 W_{ij}$$
 (30)

$$W_{I,ij} = 0.5W_{ij} \tag{31}$$

For the simulations of varying the local inhibitory connection strengths, we specifically change the value of $g_{EI,0}$ and $g_{EI,scaling}$ for $g_{EI,i}$. For each combination of parameters of $g_{EI,0}$ and $g_{EI,scaling}$, simulations are performed for default parameters (no changes to PV gradient or CIB), PV average (PV gradient is replaced by average value) and CIB average (CIB is replaced by average value). The average firing rate of all areas and number of areas showing persistent activity are quantified for each parameter combination.

In other simulations, we varied the parameters $g_{EI,0}$ and $g_{EI,scaling}$ with long range connections μ_{EE} and μ_{IE} set to be 0. This enabled us to discover the range of parameter values for which individual areas were capable of maintaining persistent activity without input from other areas. In practice, the only key parameter that determines this behaviour is the smallest inhibitory connection strength of any area, $g_{EI,i} = g_{EI,0}$.

⁷⁷⁵ Simulations and theoretical calculation of the baseline stability of ⁷⁷⁶ the network

In the simulation focusing on the stability of the baseline state of the network, there was no external input provided to any of the areas apart from noise (Eq. 12). The steady firing rate of each area after 10 s is recorded as a measure of the baseline stability.

We tested the baseline stability on five different scenarios (Fig. 4A-B) : In (1) and (2) we set the long-range connections μ_{EE} and μ_{IE} to zero since we focus on the local network. In (2), we also set the local inhibitory connections $g_{EI,0}$ to zero. In (3) - (5) the long-range connections are intact. In (4), we set the long-range connection to inhibitory neurons μ_{IE} to zero. In (5), we set the local inhibitory connections $g_{EI,0}$ to zero.

We analytically calculated the stability of baseline state for a local circuit when the long range connections μ_{EE} and μ_{IE} are set to zero, which means $I_{LR,A}$, $I_{LR,B}$, $I_{LR,C}$ are zero in Eqs 9-11. In Eqs 15 and 16, we assume that r_A , r_B and r_C reach their steady states instantaneously, since its time constant τ_r is much smaller than the time constant of NMDA synaptic variable τ_N in Eqs 6 and 7. Thus, we can express the firing rate r_A , r_B and r_C as

⁷⁹⁰ functions of synaptic variables S_A , S_B and s_C (Eqs. 13-16):

$$r_A = \phi_A(g_{E,self}S_A + g_{E,cross}S_B - g_{EI}S_C + I_{0A})$$
(32)

$$r_B = \phi_A (g_{E,self} S_B + g_{E,cross} S_A - g_{EI} S_C + I_{0A}) \tag{33}$$

$$r_C = \phi_C (g_{IE} S_A + g_{IE} S_B - g_{II} S_C + I_{0C}) \tag{34}$$

where ϕ_A and ϕ_C have the same form as Eqs 13 and 14.

Then we can insert Eqs 32-34 into Eqs 6-8 to obtain a differential equation for S_A , S_B and S_C .

$$\frac{dS_A}{dt} = -\frac{S_A}{\tau_N} + \gamma (1 - S_A) \phi_A (g_{E,self} S_A + g_{E,cross} S_B - g_{EI} S_C + I_{0A})$$
(35)

$$\frac{dS_B}{dt} = -\frac{S_B}{\tau_N} + \gamma (1 - S_B) \phi_A (g_{E,self} S_B + g_{E,cross} S_A - g_{EI} S_C + I_{0A})$$
(36)

$$\frac{dS_C}{dt} = -\frac{S_C}{\tau_G} + \gamma_I \phi_C (g_{IE}S_A + g_{IE}S_B - g_{II}S_C + I_{0C})$$
(37)

The steady state of S_A , S_B and S_C can be solved numerically by setting the left side of the above equations to be zero. We denote the right side of the equations as FA, FB and FC. Then we can calculate the Jacobian matrix and its eigenvalues.

$$J_{S_A,S_B,S_C} = \begin{bmatrix} \frac{dFA}{dS_A} & \frac{dFA}{dS_B} & \frac{dFA}{dS_C} \\ \frac{dFB}{dS_A} & \frac{dFB}{dS_B} & \frac{dFB}{dS_C} \\ \frac{dFC}{dS_A} & \frac{dFC}{dS_B} & \frac{dFC}{dS_C} \end{bmatrix}$$
(38)

⁷⁹⁷ If the real part of all the eigenvalues are negative, then that means the baseline state is stable. ⁷⁹⁸ The eigenvalues of the scenario (1) are -10.4, -12.5 and -229.8, while those of scenario (2), ⁷⁹⁹ where local inhibitory connections $g_{EI,0}$ are zero, are -7.4, -7.9, -232.3. These results coincide ⁸⁰⁰ with the simulation results of Fig. 4A.

We also considered an alternative parameter regime, where the local excitatory connections 801 $g_{E,self}$ is set to a higher level $g_{E,self} = 0.6nA$. The local inhibitory connections strength $g_{EI,0}$ 802 is also set to a higher level $g_{EI,0} = 0.5nA$ to balance the increased excitatory connections. 803 Under such alternative parameter regime, we performed similar analysis as the five different 804 scenarios in Fig. 4A-B. The results are shown in Fig. 4C-D. In simulations of a network 805 with intact long-range connections and increased local excitatory connections (in Fig. 4D 806 and also in Fig. 4F), we changed the long-range connections strength $\mu_{EE} = 0.19nA$. In Fig. 807 4C, when we gradually decrease the inhibitory connection strength $g_{EI,0}$ from 0.5nA to 0 808 (from blue dots to orange dots), analytical calculations demonstrate that the stable low firing 809 rate state disappears via a saddle node bifurcation at $g_{EI,0} = 0.175nA$ (for area AIp). This 810 demonstrates that, upon removal of inhibition, the high firing rate in Fig. 4C corresponds to 811 a distinct state and not simply a shift of the baseline state. 812

In the increased local excitatory connection regime, we further introduced temporary external input to each local brain areas and record its stable firing rate shown in Fig. 4E. In the simulation of Fig. 4F, we used the classic simulation protocol: an temporary external input is given to primary visual cortex and the delay period firing rate of each areas are recorded and shown.

⁸¹⁸ Thalamocortical network model

⁸¹⁹ Corticothalamic connectivity. We introduced thalamic areas in the network to examine their ⁸²⁰ effect on cortical dynamics. Each thalamic area includes 2 excitatory populations, A and B, ⁸²¹ with no inhibitory population. These two populations share the same selectivity with the ⁸²² corresponding cortical areas. Unlike cortical areas, there are no recurrent connections between ⁸²³ thalamic neurons (Sherman 2007). Thalamic currents have the following contributions (tc ⁸²⁴ stands for thalamocortical connections and ct for corticothalamic connections):

$$I_{th,A,B} = I_{ct,A,B} + I_{th,0,A,B} + I_{th,noise,A,B}$$

$$\tag{39}$$

where $I_{th,i}$ (i = A, B) is the total current received by each thalamic population, $I_{ct,i}$ (i = A, B)is the long range current from cortical areas to target thalamic area, $I_{th,0,i}$ (i = A, B) is the background current for each population, and $I_{th,noise,i}$ (i = A, B) is the noise input to thalamic population A and B, which we set to 0 in our simulations. $I_{ct,i}$ (i = A, B) has the following form:

$$I_{ct,A,B,i} = g_{ct} W_{ct,E,ij} S_{k,j} \tag{40}$$

where $W_{ct,E,ij}$ is the LR connectivity to thalamic neurons, and $S_{k,j}$ is the synaptic variable of population k (k = A, B) in cortical area j. Since all thalamic neurons are excitatory, we model corticothalamic projections as in the previous section:

$$m_{ct,ij} = \frac{1}{1 + e^{-\beta(h_{th,i} - h_j)}} \tag{41}$$

$$W_{ct,E,ij} = m_{ct,ij} W_{ct,ij} \tag{42}$$

(43)

where

$$W_{ct,scale,ij} = (W_{ct,norm,ij})^{k_{scale}}$$

$$\tag{44}$$

$$W_{ct,ij} = \frac{W_{ct,scale,ij}}{max(W_{ct,scale,ij})}$$
(45)

 $W_{ct,norm,ij}$ is the normalized connection strength from cortical area j to thalamic area i. $m_{ct,ij}$ is the coefficient quantifying how the long range connections target excitatory neurons based on cortical hierarchy h_j and thalamic hierarchy $h_{th,i}$.

⁸³⁶ The thalamic firing rates are described by:

$$\tau_r \frac{dr_{th,A,B}}{dt} = -r_{th,A,B} + \phi_{th,A,B}(I_{th,A,B}) \tag{46}$$

⁸³⁷ with the activation function for thalamic neurons given by:

$$\phi_{th,A,B}(I_{th,A,B}) = \frac{aI_{th,A,B} - b}{1 - exp[-d(aI_{th,A,B} - b)]}$$
(47)

⁸³⁸ Thalamic neurons are described by AMPA synaptic variables (Jaramillo et al. 2019):

$$\frac{dS_{th,A,B}}{dt} = -\frac{S_{th,A,B}}{\tau_A} + \gamma_A r_{th,A,B} \tag{48}$$

Thalamocortical connectivity. The connections from thalamic neurons to cortical neurons follow these equations

$$I_{tc,A,B,i} = g_{E,tc} W_{E,tc,ij} S_{th,A,B,j}$$

$$\tag{49}$$

$$I_{tc,C,i} = g_{I,tc} W_{I,tc,ij} (S_{th,A,j} + S_{th,B,j})$$
(50)

and connectivity

$$m_{tc,ij} = \frac{1}{1 + e^{-\beta(h_i - h_{th,j})}}$$
(51)

$$W_{E,tc,ij} = m_{tc,ij} W_{tc,ij} \tag{52}$$

$$W_{I,tc,ij} = (1 - m_{tc,ij})W_{tc,ij}$$
 (53)

and connectivity matrix

$$W_{tc,scale,ij} = (W_{tc,norm,ij})^{k_{scale}}$$
(54)

$$W_{tc,ij} = \frac{W_{tc,scale,ij}}{max(W_{tc,scale,ij})}$$
(55)

⁸³⁹ The thalamocortical input is added to the total input current of each cortical population.

$$I_A = g_{E,self} S_A + g_{E,cross} S_B + g_{EI} S_C + I_{0A} + I_{LR,A} + I_{tc,A} + x_A(t)$$
(56)

$$I_B = g_{E,self} S_B + g_{E,cross} S_A + g_{EI} S_C + I_{0B} + I_{LR,B} + I_{tc,B} + x_B(t)$$
(57)

$$I_C = g_{IE}S_A + g_{IE}S_B + g_{II}S_C + I_{0C} + I_{LR,C} + I_{tc,C} + x_C(t)$$
(58)

⁸⁴⁰ Calculation of network structural measures

We considered three types of structural measures. The first one is input strength. Input strength of area i is the summation of the connection strengths onto node i. It quantifies the total external input onto area i.

$$W_{input,i} = \sum_{j=1}^{n} W_{ij} \tag{59}$$

The second one is eigenvector centrality (Newman 2018). Eigenvector centrality of area i is the ith element of the leading eigenvector of the connectivity matrix. It quantifies how many areas are connected with the target area i and how important these neighbors are. W is a matrix where each element is W_{ij} .

$$W = Q\Lambda Q^{-1} \tag{60}$$

$$C_{eig,i} = q_{i1} \tag{61}$$

The third structural measure is loop strength, which quantifies how each area is involved in strong recurrent loops. We first define the strength of a single loop k

$$L_k = \prod_{A_i, A_j \in loop_k} W_{ij} \tag{62}$$

and then the loop strength S_{A_i} of a single area A_i

$$S_{A_i} = \sum_{A_i \in loop_k} L_k \tag{63}$$

We now focus on cell type-specific structural measures. Cell type specificity is introduced via a coefficient k_{cell} that scales all long range connection strengths (cell type projection coefficient):

$$k_{cell} = m_{ij} - PV_i(1 - m_{ij})$$
 (64)

Thus, we can define cell type-specific connectivity as:

$$W_{cell,ij} = (m_{ij} - PV_i(1 - m_{ij}))W_{ij}$$
(65)

⁸⁵⁰ The cell type-specific connectivity is further normalized so that the maximum value is 1.

$$\tilde{W}_{ij} = \frac{W_{cell,ij}}{max(W_{cell,ij})} \tag{66}$$

and cell type-specific input strength could be defined as:

$$W_{input,i,cellspec} = \sum_{j=1}^{n} \tilde{W}_{ij} \tag{67}$$

Similarly, cell type-specific eigenvector centrality is defined as

$$\tilde{W} = \tilde{Q}\tilde{\Lambda}\tilde{Q}^{-1} \tag{68}$$

$$C_{eig,i,cellspec} = \tilde{q}_{i1} \tag{69}$$

where \tilde{W} is a matrix where each element is \tilde{W}_{ij} and the cell type-specific loop strength is defined as:

$$L_{k,cellspec} = \prod_{A_i, A_j \in loop_k} \tilde{W}_{ij} \tag{70}$$

$$S_{A_i,cellspec} = \sum_{A_i \in loop_k} L_{k,cellspec}$$
(71)

As a comparison, we also calculated the sign-only loop strength and no PV loop strength. We can define sign-only connectivity as:

$$W_{signonly,ij} = sgn(m_{ij} - (1 - m_{ij}))W_{ij}$$

$$\tag{72}$$

where sgn(x) is the sign function, which returns positive or negative values based on the sign

of x. The major difference between sign only connectivity and cell type specific connectivity

is that the strength of long range projection bias are not considered except the sign of it.

We can also define no-PV connectivity as:

$$W_{noPV,ij} = (m_{ij} - (1 - m_{ij}))W_{ij}$$
(73)

The difference between no-PV connectivity and cell-type specific connectivity is that the different strengths of local connections for each area are not considered in the no-PV connectivity.

We also used the sign-only and no-PV variants of connectivity measures to predict the delay period firing rate and classify core areas. This enabled us to compare these simplified measures to the cell-type-specific connectivity measures.

⁸⁶² Stimulation protocol and inhibition analysis

The model is simulated using an stochastic differential equation solver: Euler-Maruyama method. We write customized program using Python to implement this numerical method. The time step is set to be dt, and all the firing rates, synaptic variables and currents are initialized to be zero.

We simulate a working memory task by applying an external current I_{stim} to one of the 867 excitatory populations, which represents a sensory (e.g., visual) stimulus that is to be 868 remembered across a delay period. The external current is a pulsed input with start time T_{on} 869 and offset time T_{off} . Without losing generality, we assume that the external input is provided 870 to population A. In most of the simulations in this study, we simulate a visual working 871 memory task, with the external input applied to VISp. The simulation duration is T_{trial} and 872 we used a time step of dt. The delay period is defined as the duration between the offset time 873 T_{off} and trial end T_{trial} . In order to obtain a stable firing rate, the delay period firing rate is 874 calculated by averaging the firing rate from 2 seconds until the end of the delay period to 875 0.5 seconds until end. Firing rate, PV cell fraction, and hierarchy are plotted on a 3d brain 876 surface using the website scalable brain atlas (https://scalablebrainatlas.incf.org/index.php). 877 We apply inhibition analysis to understand the robustness of attractors and, more importantly, 878 to investigate which areas play an important role in maintaining the attractor state. Excitatory 879 input was applied to the inhibitory population I to simulate opto-genetic inhibition. The 880 external input I_{inh} is strong as compared to I_{stim} and results in an elevated firing rate of the 881 inhibitory population, which in turn decreases the firing rate of the excitatory populations. 882 Usually the inhibition is applied to a single area. When inhibition is applied during the 883 stimulus period, its start and end times are equal to T_{on} and T_{off} , respectively. When 884 inhibition is applied during delay period, its start time is later than T_{off} to allow the system 885 settle to a stable state. Thus, the onset of inhibition starts 2 seconds after T_{off} and lasts 886 until the end of trial. In the case of thalamocortical network simulations, we inhibit thalamic 887

areas by introducing a hyperpolarizing current to both excitatory populations, since we do not have inhibitory populations in thalamic areas in the model.

To quantify the effect of single area or multiple areas inhibition, we calculate the average firing 890 rate of areas that satisfy two conditions: i) the area shows persistent activity before inhibition 891 and ii) the area does not receive inhibitory input. The ratio between such average firing 892 rate after inhibition and before inhibition is used to quantify the overall effect of inhibition. 893 If the ratio is lower than 100%, this suggests that inhibiting certain area(s) disrupts the 894 maintenance of the attractor state. Note that the inhibition effect is typically not very strong, 895 and only in rare cases, inhibition of a single area leads to loss of activity of other areas (Fig. 7B, Fig. 7C). To quantify such differences, we use a threshold of 10% to differentiate them. 897 We will use (relatively) "weak inhibition effect" and "strong inhibition effect" to refer to 898 them afterwards. 899

We used the three measures to classify areas into 4 types (Fig. 7D): i) inhibition effect during delay period, ii) inhibition effect during stimulus period, and iii) delay period firing rate. Areas with strong inhibition effect during stimulus period are classified as input areas; areas with strong inhibition effect during delay period and strong delay period firing rate are classified as core areas; areas with weak inhibition effect during delay period but strong firing rate are classified as readout areas; areas with weak inhibition effect during delay period and weak firing rate during delay period are classified as nonessential areas.

As an extension of the single area inhibition study, we focus on the role of readout areas. A 907 pair of readout areas is randomly chosen and inhibited during the delay period under a similar 908 protocol as the single area inhibition study. The inhibition effect, i.e., the decrement of the 909 delay period firing rates of other non inhibited areas, is first quantified for each inhibition 910 pair (A_i, A_i) . Next, the inhibition effect is averaged one more time for each area A_i across 911 all inhibition pairs that includes the area $((A_i, A_j), \text{ where } j \neq i)$. An anologous procedure is 912 performed for triplets and quadruplets of readout areas. Additionally, we also calculate the 913 mean inhibition effect between pair of areas, which are both selected from core areas, both 914 selected from readout areas, or we chose one area from core areas, one area from readout 915 areas. 916

⁹¹⁷ Simulation of multiple attractors

Multiple attractors coexist in the network and its properties and number depends on the 918 connectivity and dynamics of each node. In this study we did not try to capture all the 919 possible attractors in the network, but rather compare the number of attractors for different 920 networks. Here we briefly describe the protocol used to identify multiple attractors in the 921 network. We first choose k areas and then generate a subset of areas as the stimulation areas. 922 We cover all possible subsets, which means we run 2^k simulations in total. The external 923 stimulus is given to all areas in the subset simultaneously with same strength and duration. 924 The delay period activity is then quantified using a similar protocol as the standard simulation 925

protocol. The selection of k areas corresponds to a qualitative criterion. First we choose 926 the areas with small PV fraction or high hierarchy, since these areas are more likely to show 927 persistent activity. Second, the number of possible combination grows exponentially as we 928 increase k, and if we use k = 43, the number of combinations is around 8.8e+12, which is 920 beyond our simulation power. As a trade-off between the simulation power and coverage of 930 areas, we choose k = 18, which correspond to 2.6e+5 different combinations of stimulation. 931 For each parameter setting, we run 2.6e+5 simulations to capture possible attractor patterns. 932 For each attractor pattern, a binary vector is generated by thresholding delay firing rate using 933 a firing rate threshold of 5Hz. An attractor pattern is considered distinct if and only if the binary vector is different from all identified attractors. In these way we can identify different 935 attractors in the simulation. We also apply same simulation pipeline to identify attractors 936 for different parameters. Specifically we change the long range connectivity strength μ_{EE} 937 and local excitatory connections $g_{E,self}$. 938

⁹³⁹ Data availability

The manuscript constitutes a computational study, so no experimental data has been generated.

⁹⁴¹ The simulation and analysis code will be available in GitHub upon publication.

942 References

- Abbott, L. F. et al. (2020). "The mind of a mouse". Cell 182, pp. 1372–1376.
- Lorente de Nó, R. (1933). "Vestibulo-ocular reflex arc". Arch. Neurol. Psych. 30, pp. 245–291.
- Abbott, Larry F. et al. (2017). "An International Laboratory for Systems and Computational

946 Neuroscience". *Neuron* 96.6, pp. 1213–1218.

- Baddeley, Alan (2012). "Working Memory: Theories, Models, and Controversies." Annual *review of psychology* 63.1, pp. 1–29. pmid: 21961947.
- Ballesteros-Yáñez, Inmaculada, Ruth Benavides-Piccione, Jean-Pierre Bourgeois, Jean-Pierre
 Changeux, and Javier DeFelipe (2010). "Alterations of Cortical Pyramidal Neurons in
 Mice Lacking High-Affinity Nicotinic Receptors". Proceedings of the National Academy of
 Sciences 107.25, pp. 11567–11572.
- Bassett, Danielle S. and Edward T. Bullmore (2017). "Small-World Brain Networks Revisited".
 The Neuroscientist 23.5, pp. 499–516.
- ⁹⁵⁵ Bolkan, Scott S., Joseph M. Stujenske, Sebastien Parnaudeau, Timothy J. Spellman, Caroline
- Rauffenbart, Atheir I. Abbas, Alexander Z. Harris, Joshua A. Gordon, and Christoph
- ⁹⁵⁷ Kellendonk (2017). "Thalamic Projections Sustain Prefrontal Activity during Working
- Memory Maintenance". Nature Neuroscience 20.7 (7), pp. 987–996.
- Burt, Joshua B., Murat Demirtaş, William J. Eckner, Natasha M. Navejar, Jie Lisa Ji,
- William J. Martin, Alberto Bernacchia, Alan Anticevic, and John D. Murray (2018).

- "Hierarchy of Transcriptomic Specialization across Human Cortex Captured by Structural
 Neuroimaging Topography". Nature Neuroscience 21.9, pp. 1251–1259.
- ⁹⁶³ Cabral, Joana, Etienne Hugues, Olaf Sporns, and Gustavo Deco (2011). "Role of Local Network
- Oscillations in Resting-State Functional Connectivity". NeuroImage 57.1, pp. 130–139.
- ⁹⁶⁵ Campagnola, Luke et al. (2022). "Local Connectivity and Synaptic Dynamics in Mouse and
- Human Neocortex". Science 375.6585, eabj5861.
- ⁹⁶⁷ Chaudhuri, Rishidev, Kenneth Knoblauch, Marie-Alice Gariel, Henry Kennedy, and Xiao-Jing
 ⁹⁶⁸ Wang (2015). "A Large-Scale Circuit Mechanism for Hierarchical Dynamical Processing

⁹⁶⁹ in the Primate Cortex". Neuron 88.2, pp. 419–431.

- 970 Christophel, Thomas B., P. Christiaan Klink, Bernhard Spitzer, Pieter R. Roelfsema, and
- John-Dylan Haynes (2017). "The Distributed Nature of Working Memory". Trends in
 Cognitive Sciences 21.2, pp. 111–124.
- 973 Crabtree, John W. (2018). "Functional Diversity of Thalamic Reticular Subnetworks". Frontiers
 974 in Systems Neuroscience 12, p. 41.
- D. J. Amit (1995). "The Hebbian paradigm reintegrated: local reverberations as internal
 representations". *Behav. Brain Sci.* 18, pp. 617–626.
- Deco, Gustavo, Adrián Ponce-Alvarez, Patric Hagmann, Gian Luca Romani, Dante Mantini,
 and Maurizio Corbetta (2014). "How Local Excitation–Inhibition Ratio Impacts the Whole
 Brain Dynamics". Journal of Neuroscience 34.23, pp. 7886–7898. pmid: 24899711.
- Demirtaş, Murat, Joshua B. Burt, Markus Helmer, Jie Lisa Ji, Brendan D. Adkinson, Matthew
 F. Glasser, David C. Van Essen, Stamatios N. Sotiropoulos, Alan Anticevic, and John D.
- Murray (2019). "Hierarchical Heterogeneity across Human Cortex Shapes Large-Scale
 Neural Dynamics". Neuron 101.6, 1181–1194.e13.
- ⁹⁸⁴ Diesmann, Markus, Marc-Oliver Gewaltig, and Ad Aertsen (1999). "Stable Propagation of

⁹⁸⁵ Synchronous Spiking in Cortical Neural Networks". *Nature* 402.6761 (6761), pp. 529–533.

⁹⁸⁶ Dotson, Nicholas M., Steven J. Hoffman, Baldwin Goodell, and Charles M. Gray (2018).

- "Feature-Based Visual Short-Term Memory Is Widely Distributed and Hierarchically
 Organized". Neuron 99.1, 215–226.e4.
- Egorov, Alexei V., Bassam N. Hamam, Erik Fransén, Michael E. Hasselmo, and Angel
 A. Alonso (2002). "Graded Persistent Activity in Entorhinal Cortex Neurons". Nature
 420.6912, pp. 173–178.
- ⁹⁹² Elston, Guy N. (2007). "Specialization of the Neocortical Pyramidal Cell during Primate
 ⁹⁹³ Evolution". In: *Evolution of Nervous Systems*. Elsevier, pp. 191–242.
- Elston, Guy N. and Marcello G.P. Rosa (1998). "Complex Dendritic Fields of Pyramidal Cells
 in the Frontal Eye Field of the Macaque Monkey: Comparison with Parietal Areas 7a and
- ⁹⁹⁶ LIP". *NeuroReport* 9.1, pp. 127–131.
- ⁹⁹⁷ Erlich, Jeffrey C., Max Bialek, and Carlos D. Brody (2011). "A Cortical Substrate for
 ⁹⁹⁸ Memory-Guided Orienting in the Rat". Neuron 72.2, pp. 330–343.
- ⁹⁹⁹ Erö, Csaba, Marc-Oliver Gewaltig, Daniel Keller, and Henry Markram (2018). "A Cell Atlas
 ¹⁰⁰⁰ for the Mouse Brain". Frontiers in Neuroinformatics 12, p. 84.

- Felleman, D J and D C Van Essen (1991). "Distributed Hierarchical Processing in the Primate
 Cerebral Cortex." *Cerebral Cortex* 1.1, pp. 1–47. pmid: 1822724.
- 1003 Froudist-Walsh, Sean, Daniel P. Bliss, Xingyu Ding, Lucija Rapan, Meiqi Niu, Kenneth
- Knoblauch, Karl Zilles, Henry Kennedy, Nicola Palomero-Gallagher, and Xiao-Jing Wang
 (2021). "A Dopamine Gradient Controls Access to Distributed Working Memory in the
- Large-Scale Monkey Cortex". *Neuron* 109.21, 3500–3520.e13.
- Froudist-Walsh, Sean, Ting Xu, Meiqi Niu, Lucija Rapan, Ling Zhao, Daniel S. Margulies,
 Karl Zilles, Xiao-Jing Wang, and Nicola Palomero-Gallagher (2023). "Gradients of Neurotransmitter Receptor Expression in the Macaque Cortex". Nature Neuroscience 26.7,
 pp. 1281–1294.
- Fulcher, Ben D., John D. Murray, Valerio Zerbi, and Xiao-Jing Wang (2019). "Multimodal
 Gradients across Mouse Cortex". *Proceedings of the National Academy of Sciences* 116.10,
 pp. 4689–4695. pmid: 30782826.
- Funahashi, Shintaro, Charles J. Bruce, and Patricia S. Goldman-Rakic (1989). "Mnemonic
 Coding of Visual Space in the Monkey's Dorsolateral Prefrontal Cortex". Journal of
 Neurophysiology 61.2, pp. 331–349.
- Fuster, Joaquin M. and Garrett E. Alexander (1971). "Neuron Activity Related to Short-Term
 Memory". Science 173.3997, pp. 652–654. pmid: 4998337.
- Gămănuţ, Răzvan, Henry Kennedy, Zoltán Toroczkai, Mária Ercsey-Ravasz, David C. Van
 Essen, Kenneth Knoblauch, and Andreas Burkhalter (2018). "The Mouse Cortical Connectivity Profiles". Neuron 97.3, 698–715.e10.
- Gao, Zhenyu, Courtney Davis, Alyse M Thomas, Michael N Economo, Amada M Abrego,
 Karel Svoboda, Chris I De Zeeuw, and Nuo Li (2018). "A Cortico-Cerebellar Loop for
 Motor Planning." *Nature* 39, p. 1062. pmid: 30333626.
- Gilad, Ariel, Yasir Gallero-Salas, Dominik Groos, and Fritjof Helmchen (2018). "Behavioral
 Strategy Determines Frontal or Posterior Location of Short-Term Memory in Neocortex".
 Neuron 99.4, 814–828.e7.
- ¹⁰²⁹ Gilman, Joshua P., Maria Medalla, and Jennifer I. Luebke (2017). "Area-Specific Features
- of Pyramidal Neurons—a Comparative Study in Mouse and Rhesus Monkey". Cerebral
 Cortex 27.3, pp. 2078–2094.
- 1032 Goldman-Rakic, P. S (1995). "Cellular Basis of Working Memory". Neuron 14.3, pp. 477–485.
- 1033 Guo, Zengcai V., Hidehiko K. Inagaki, Kayvon Daie, Shaul Druckmann, Charles R. Gerfen, and
- Karel Svoboda (2017). "Maintenance of Persistent Activity in a Frontal Thalamocortical
 Loop". Nature 545.7653 (7653), pp. 181–186.
- Guo, Zengcai V., Nuo Li, Daniel Huber, Eran Ophir, Diego Gutnisky, Jonathan T. Ting,
 Guoping Feng, and Karel Svoboda (2014). "Flow of Cortical Activity Underlying a Tactile
 Decision in Mice". Neuron 81.1, pp. 179–194.

- Hádinger, Nóra, Emília Bősz, Boglárka Tóth, Gil Vantomme, Anita Lüthi, and László Acsády 1039
- (2023). "Region-Selective Control of the Thalamic Reticular Nucleus via Cortical Layer 5 1040 Pyramidal Cells". Nature Neuroscience 26.1, pp. 116–130. 1041
- Harris, Julie A. et al. (2019). "Hierarchical Organization of Cortical and Thalamic Connectiv-1042 ity". Nature 575.7781 (7781), pp. 195–202. 1043
- Harvey, Christopher D., Philip Coen, and David W. Tank (2012). "Choice-Specific Sequences 1044 in Parietal Cortex during a Virtual-Navigation Decision Task". Nature 484,7392, pp. 62–68. 1045 pmid: 22419153. 1046
- Honey, Christopher J., Rolf Kötter, Michael Breakspear, and Olaf Sporns (2007). "Network 1047 Structure of Cerebral Cortex Shapes Functional Connectivity on Multiple Time Scales". 1048 Proceedings of the National Academy of Sciences 104.24, pp. 10240–10245. pmid: 17548818. 1049
- Inagaki, Hidehiko K., Lorenzo Fontolan, Sandro Romani, and Karel Svoboda (2019). "Discrete 1050 Attractor Dynamics Underlies Persistent Activity in the Frontal Cortex". Nature 566.7743 1051 (7743), pp. 212–217. 1052
- Inagaki, Hidehiko K., Miho Inagaki, Sandro Romani, and Karel Svoboda (2018). "Low-1053 Dimensional and Monotonic Preparatory Activity in Mouse Anterior Lateral Motor 1054 Cortex". The Journal of Neuroscience 38.17, pp. 4163–4185. 1055
- Jaramillo, Jorge, Jorge F. Mejias, and Xiao-Jing Wang (2019). "Engagement of Pulvino-1056 cortical Feedforward and Feedback Pathways in Cognitive Computations". Neuron 101.2, 1057 321-336.e9. 1058
- Javadzadeh, Mitra and Sonja B. Hofer (2022). "Dynamic Causal Communication Channels 1059 between Neocortical Areas". Neuron 0.0. 1060
- Joglekar, Madhura R., Jorge F. Mejias, Guangyu Robert Yang, and Xiao-Jing Wang (2018). 1061 "Inter-Areal Balanced Amplification Enhances Signal Propagation in a Large-Scale Circuit 1062 Model of the Primate Cortex". Neuron 98.1, 222–234.e8. pmid: 29576389.
- Jones, Edward G. (2007). "Neuroanatomy: Cajal and after Cajal". Brain Research Reviews 1064 55.2, pp. 248–255. 1065
- Jonikaitis, Donatas, Behrad Noudoost, and Tirin Moore (2023). Dissociating the Contributions 1066 of Frontal Eye Field Activity to Spatial Working Memory and Motor Preparation. preprint. 1067 Neuroscience. 1068
- Jun, James J. et al. (2017). "Fully Integrated Silicon Probes for High-Density Recording of 1069 Neural Activity". Nature 551.7679 (7679), pp. 232–236. 1070
- Kim, Yongsoo, Guangyu Robert Yang, Kith Pradhan, Kannan Umadevi Venkataraju, Mi-1071 hail Bota, Luis Carlos García del Molino, Greg Fitzgerald, Keerthi Ram, Miao He, 1072 1073 Jesse Maurica Levine, Partha Mitra, Z. Josh Huang, Xiao-Jing Wang, and Pavel Osten (2017). "Brain-Wide Maps Reveal Stereotyped Cell-Type-Based Cortical Architecture and 1074
- Subcortical Sexual Dimorphism". Cell 171.2, 456–469.e22. 1075

1063

Klatzmann, Ulysse, Sean Froudist-Walsh, Daniel P. Bliss, Panagiota Theodoni, Jorge F. 1076 Mejias, Meiqi Niu, Lucija Rapan, Nicola Palomero-Gallagher, Claire Sergent, Stanislas 1077

- Dehaene, and Xiao-Jing Wang (2022). "A Connectome-Based Model of Conscious Access 1078 in Monkey Cortex". bioRxiv, p. 2022.02.20.481230. 1079
- Knox, Joseph E., Kameron Decker Harris, Nile Graddis, Jennifer D. Whitesell, Hongkui 1080 Zeng, Julie A. Harris, Eric Shea-Brown, and Stefan Mihalas (2018). "High-Resolution 1081 Data-Driven Model of the Mouse Connectome". Network Neuroscience 3.1, pp. 217–236. 1082 pmid: 30793081.
- Kopec, Charles D., Jeffrey C. Erlich, Bingni W. Brunton, Karl Deisseroth, and Carlos D. Brody 1084 (2015). "Cortical and Subcortical Contributions to Short-Term Memory for Orienting 1085 Movements". Neuron 88.2, pp. 367–377. pmid: 26439529. 1086

1083

- Kubanek, Jan and Lawrence H. Snyder (2015). "Reward-Based Decision Signals in Parietal 1087 Cortex Are Partially Embodied". The Journal of Neuroscience 35.12, pp. 4869–4881. 1088
- Leavitt, Matthew L., Diego Mendoza-Halliday, and Julio C. Martinez-Trujillo (2017). "Sus-1089 tained Activity Encoding Working Memories: Not Fully Distributed". Trends in Neuro-1090 sciences 40.6, pp. 328–346. 1091
- Li, Nuo, Tsai-Wen Chen, Zengcai V. Guo, Charles R. Gerfen, and Karel Svoboda (2015). 1092
- "A Motor Cortex Circuit for Motor Planning and Movement". Nature 519.7541 (7541), 1093 pp. 51–56. 1094
- Li, Nuo, Kayvon Daie, Karel Svoboda, and Shaul Druckmann (2016a). "Robust Neuronal 1095 Dynamics in Premotor Cortex during Motor Planning". Nature 532.7600, pp. 459–464. 1096
- Li, Nuo, Kayvon Daie, Karel Svoboda, and Shaul Druckmann (2016b). "Robust Neuronal 1097 Dynamics in Premotor Cortex during Motor Planning". Nature 532.7600, pp. 459–464. 1098
- Liu, Ding, Xiaowei Gu, Jia Zhu, Xiaoxing Zhang, Zhe Han, Wenjun Yan, Qi Cheng, Jiang 1099 Hao, Hongmei Fan, Ruiqing Hou, Zhaoqin Chen, Yulei Chen, and Chengyu T. Li (2014). 1100 "Medial Prefrontal Activity during Delay Period Contributes to Learning of a Working 1101 Memory Task". Science 346.6208, pp. 458–463.
- 1102 Markov, Nikola T, Julien Vezoli, Pascal Chameau, Arnaud Falchier, René Quilodran, Cyril 1103
- Huissoud, Camille Lamy, Pierre Misery, Pascale Giroud, Shimon Ullman, Pascal Barone, 1104 Colette Dehay, Kenneth Knoblauch, and Henry Kennedy (2014a). "Anatomy of Hierarchy: 1105 Feedforward and Feedback Pathways in Macaque Visual Cortex." Journal of Comparative 1106 Neurology 522.1, pp. 225–259. pmid: 23983048. 1107
- Markov, Nikola T. et al. (2014b). "A Weighted and Directed Interareal Connectivity Matrix 1108 for Macaque Cerebral Cortex". Cerebral Cortex 24.1, pp. 17–36. 1109
- Mejias, Jorge F., John D. Murray, Henry Kennedy, and Xiao-Jing Wang (2016). "Feedforward 1110
- and Feedback Frequency-Dependent Interactions in a Large-Scale Laminar Network of 1111 the Primate Cortex". Science Advances 2.11, e1601335. 1112
- Mejias, Jorge F. and Xiao-Jing Wang (2022). "Mechanisms of Distributed Working Memory 1113 in a Large-Scale Network of Macaque Neocortex". eLife 11, e72136. 1114
- Meng, John Hongyu, Benjamin Schuman, Bernardo Rudy, and Xiao-Jing Wang (2023). 1115
- "Mechanisms of Dominant Electrophysiological Features of Four Subtypes of Layer 1 1116 Interneurons". The Journal of Neuroscience 43.18, pp. 3202–3218. 1117

- Murphy, B. K. and K. D. Miller (2009). "Balanced amplification: a new mechanism of selective
 amplification of neural activity patterns". *Neuron* 61, pp. 635–648.
- ¹¹²⁰ Murray, John D, Alberto Bernacchia, David J Freedman, Ranulfo Romo, Jonathan D Wallis,
- 1121 Xinying Cai, Camillo Padoa-Schioppa, Tatiana Pasternak, Hyojung Seo, Daeyeol Lee,
- and Xiao-Jing Wang (2014). "A Hierarchy of Intrinsic Timescales across Primate Cortex".

¹¹²³ Nature Neuroscience 17.12, pp. 1661–1663.

- Murray, John D., Jorge Jaramillo, and Xiao-Jing Wang (2017). "Working Memory and
 Decision-Making in a Frontoparietal Circuit Model". Journal of Neuroscience 37.50,
 pp. 12167–12186. pmid: 29114071.
- Musall, S., M. T. Kaufman, A. L. Juavinett, S. Gluf, and A. K. Churchland (2019). "Singletrial neural dynamics are dominated by richly varied movements". *Nat. Neurosci.* 22, pp. 1677–1686.
- ¹¹³⁰ Myme, Chaelon I. O., Ken Sugino, Gina G. Turrigiano, and Sacha B. Nelson (2003). "The

1131 NMDA-to-AMPA Ratio at Synapses Onto Layer 2/3 Pyramidal Neurons Is Conserved

- Across Prefrontal and Visual Cortices". *Journal of Neurophysiology* 90.2, pp. 771–779.
- ¹¹³³ Naskar, Shovan, Jia Qi, Francisco Pereira, Charles R. Gerfen, and Soohyun Lee (2021). "Cell-
- Type-Specific Recruitment of GABAergic Interneurons in the Primary Somatosensory
 Cortex by Long-Range Inputs". *Cell Reports* 34.8, p. 108774.
- ¹¹³⁶ Newman, M. E. J. (2018). *Networks*. Second edition. Oxford: Oxford University Press.
- ¹¹³⁷ Nigro, Maximiliano José, Kasper Kjelsberg, Laura Convertino, Rajeevkumar Raveendran
- Nair, and Menno P. Witter (2022). Enrichment of Specific GABAergic Neuronal Types in
 the Mouse Perirhinal Cortex. preprint. Neuroscience.
- Oh, Seung Wook et al. (2014). "A Mesoscale Connectome of the Mouse Brain". Nature 508.7495 (7495), pp. 207–214.
- ¹¹⁴² P. S. Goldman-Rakic (1995). "Cellular basis of working memory". Neuron 14, pp. 477–485.
- Petreanu, Leopoldo, Tianyi Mao, Scott M Sternson, and Karel Svoboda (2009). "The Subcellular Organization of Neocortical Excitatory Connections." *Nature* 457.7233, pp. 1142–1145.
- 1145 pmid: 19151697.
- Pinto, Lucas, Kanaka Rajan, Brian DePasquale, Stephan Y. Thiberge, David W. Tank, and
 Carlos D. Brody (2019). "Task-Dependent Changes in the Large-Scale Dynamics and
 Necessity of Cortical Regions". Neuron 104.4, 810–824.e9.
- R. J. Douglas, C. Koch, M. Mahowald, K. M. Martin, and H. H. Suarez (1995). "Recurrent
 excitation in neocortical circuits". *Science* 269, pp. 981–985.
- ¹¹⁵¹ Sanzeni, Alessandro, Bradley Akitake, Hannah C Goldbach, Caitlin E Leedy, Nicolas Brunel,
 ¹¹⁵² and Mark H Histed (2020). "Inhibition Stabilization Is a Widespread Property of Cortical
- ¹¹⁵³ Networks". *eLife* 9, e54875.
- Schmidt, Maximilian, Rembrandt Bakker, Kelly Shen, Gleb Bezgin, Markus Diesmann, and
 Sacha Jennifer van Albada (2018). "A Multi-Scale Layer-Resolved Spiking Network Model
- of Resting-State Dynamics in Macaque Visual Cortical Areas". *PLOS Computational*
- ¹¹⁵⁷ *Biology* 14.10, e1006359.

- Schmitt, L. Ian, Ralf D. Wimmer, Miho Nakajima, Michael Happ, Sima Mofakham, and
 Michael M. Halassa (2017). "Thalamic Amplification of Cortical Connectivity Sustains
 Attentional Control". *Nature* 545.7653 (7653), pp. 219–223.
- Shen, Shan, Xiaolong Jiang, Federico Scala, Jiakun Fu, Paul Fahey, Dmitry Kobak, Zhenghuan
 Tan, Na Zhou, Jacob Reimer, Fabian Sinz, and Andreas S. Tolias (2022). "Distinct
 Organization of Two Cortico-Cortical Feedback Pathways". Nature Communications 13.1,
- 1164 p. 6389.
- Sherman, S Murray (2007). "The Thalamus Is More than Just a Relay". Current opinion in
 neurobiology 17.4, pp. 417–422.
- Shine, James M, Matthew J Aburn, Michael Breakspear, and Russell A Poldrack (2018).
 "The Modulation of Neural Gain Facilitates a Transition between Functional Segregation
 and Integration in the Brain". *eLife* 7. Ed. by Gustavo Deco, e31130.
- ¹¹⁷⁰ Siu, Caitlin, Justin Balsor, Sam Merlin, Frederick Federer, and Alessandra Angelucci (2021).
- "A Direct Interareal Feedback-to-Feedforward Circuit in Primate Visual Cortex". Nature
 Communications 12.1, p. 4911.
- 1173 Steinmetz, Nicholas A., Peter Zatka-Haas, Matteo Carandini, and Kenneth D. Harris (2019).
- "Distributed Coding of Choice, Action and Engagement across the Mouse Brain". Nature 576.7786 (7786), pp. 266–273.
- Steinmetz, Nicholas A. et al. (2021). "Neuropixels 2.0: A Miniaturized High-Density Probe
 for Stable, Long-Term Brain Recordings". *Science* 372.6539, eabf4588.
- Stringer, Carsen, Marius Pachitariu, Nicholas Steinmetz, Charu Bai Reddy, Matteo Carandini,
 and Kenneth D. Harris (2019). "Spontaneous Behaviors Drive Multidimensional, Brainwide
 Activity". Science 364.6437, eaav7893.
- ¹¹⁸¹ Suzuki, Mototaka and Jacqueline Gottlieb (2013). "Distinct Neural Mechanisms of Distractor
- ¹¹⁸² Suppression in the Frontal and Parietal Lobe". *Nature Neuroscience* 16.1 (1), pp. 98–104.
- 1183 Svoboda, Karel and Nuo Li (2018). "Neural Mechanisms of Movement Planning: Motor Cortex
- and Beyond". Current Opinion in Neurobiology 49, pp. 33–41.
- Torres-Gomez, Santiago, Jackson D Blonde, Diego Mendoza-Halliday, Eric Kuebler, Michelle
 Everest, Xiao Jing Wang, Wataru Inoue, Michael O Poulter, and Julio Martinez-Trujillo
 (2020). "Changes in the Proportion of Inhibitory Interneuron Types from Sensory to
 Executive Areas of the Primate Neocortex: Implications for the Origins of Working
 Memory Representations". Cerebral Cortex 30.8, pp. 4544–4562.
- Tremblay, Robin, Soohyun Lee, and Bernardo Rudy (2016). "GABAergic Interneurons in the
 Neocortex: From Cellular Properties to Circuits". Neuron 91.2, pp. 260–292.
- Tsodyks, Misha V., William E. Skaggs, Terrence J. Sejnowski, and Bruce L. McNaughton
 (1997). "Paradoxical Effects of External Modulation of Inhibitory Interneurons". The
- Journal of Neuroscience 17.11, pp. 4382–4388.
- ¹¹⁹⁵ Vincis, Roberto, Ke Chen, Lindsey Czarnecki, John Chen, and Alfredo Fontanini (2020).
- "Dynamic Representation of Taste-Related Decisions in the Gustatory Insular Cortex of
- ¹¹⁹⁷ Mice". Current Biology 30.10, 1834–1844.e5.

- Voitov, Ivan and Thomas D. Mrsic-Flogel (2022). "Cortical Feedback Loops Bind Distributed 1198 Representations of Working Memory". Nature, pp. 1–9. 1199
- Wang, Peng, Ru Kong, Xiaolu Kong, Raphaël Liégeois, Csaba Orban, Gustavo Deco, Martijn 1200
- P. van den Heuvel, and B. T. Thomas Yeo (2019). "Inversion of a Large-Scale Circuit 1201 Model Reveals a Cortical Hierarchy in the Dynamic Resting Human Brain". Science 1202
- Advances 5.1, eaat 7854. 1203
- Wang, Quanxin et al. (2020). "The Allen Mouse Brain Common Coordinate Framework: A 1204 3D Reference Atlas". Cell 181.4, 936–953.e20. pmid: 32386544. 1205
- Wang, X.-J. (2002). "Probabilistic decision making by slow reverberation in cortical circuits". 1206 Neuron 36, pp. 955–968. 1207
- Wang, X.-J., J. Tegnér, C. Constantinidis, and P. S. Goldman-Rakic (2004). "Division of Labor 1208 among Distinct Subtypes of Inhibitory Neurons in a Cortical Microcircuit of Working 1209 Memory". Proceedings of the National Academy of Sciences 101.5, pp. 1368–1373. 1210

Wang, Xiao-Jing (1999). "Synaptic Basis of Cortical Persistent Activity: The Importance of 1211 NMDA Receptors to Working Memory". Journal of Neuroscience 19.21, pp. 9587–9603. 1212

- Wang, Xiao-Jing (2001). "Synaptic Reverberation Underlying Mnemonic Persistent Activity". 1213
- Trends in Neurosciences 24.8, pp. 455–463. 1214
- Wang, Xiao-Jing (2020). "Macroscopic Gradients of Synaptic Excitation and Inhibition in 1215 the Neocortex". Nature Reviews Neuroscience 21.3, pp. 169–178. 1216
- Wang, Xiao-Jing (2021). "50 years of mnemonic persistent activity: Quo vadis?" Trends in 1217 *Neurosci.* 44, pp. 888–902. 1218
- Wang, Xiao-Jing (2022). "Theory of the Multiregional Neocortex: Large-Scale Neural Dynamics 1219 and Distributed Cognition". Annual Review of Neuroscience 45.1, pp. 533–560. 1220
- Wang, Yu, Xinxin Yin, Zhouzhou Zhang, Jiejue Li, Wenyu Zhao, and Zengcai V. Guo (2021). 1221
- "A Cortico-Basal Ganglia-Thalamo-Cortical Channel Underlying Short-Term Memory". 1222 Neuron 109.21, 3486-3499.e7. 1223
- Wildenberg, Gregg A., Matt R. Rosen, Jack Lundell, Dawn Paukner, David J. Freedman, 1224 and Narayanan Kasthuri (2021). "Primate Neuronal Connections Are Sparse in Cortex as 1225 Compared to Mouse". Cell Reports 36.11, p. 109709. 1226
- Wong, Kong-Fatt and Xiao-Jing Wang (2006). "A Recurrent Network Mechanism of Time 1227 Integration in Perceptual Decisions". Journal of Neuroscience 26.4, pp. 1314–1328. pmid: 1228 16436619.
- Wu, Zheng, Ashok Litwin-Kumar, Philip Shamash, Alexei Taylor, Richard Axel, and Michael 1230 N. Shadlen (2020). "Context-Dependent Decision Making in a Premotor Circuit". Neuron 1231 106.2, 316-328.e6. 1232
- Xu, Yaoda (2017). "Reevaluating the Sensory Account of Visual Working Memory Storage". 1233
- Trends in Cognitive Sciences 21.10, pp. 794–815. 1234

1229

- Yang, Guangyu Robert, Madhura R. Joglekar, H. Francis Song, William T. Newsome, and 1235
- Xiao-Jing Wang (2019). "Task Representations in Neural Networks Trained to Perform 1236
- Many Cognitive Tasks". Nature Neuroscience 22.2, pp. 297–306. 1237

Yang, Sheng-Tao, Min Wang, Constantinos D Paspalas, Johanna L Crimins, Marcus T
Altman, James A Mazer, and Amy F T Arnsten (2018). "Core Differences in Synaptic
Signaling Between Primary Visual and Dorsolateral Prefrontal Cortex". *Cerebral Cortex*28.4, pp. 1458–1471.

Yizhar, Ofer, Lief E. Fenno, Thomas J. Davidson, Murtaza Mogri, and Karl Deisseroth (2011).
"Optogenetics in Neural Systems". Neuron 71.1, pp. 9–34.

- Zandvakili, Amin and Adam Kohn (2015). "Coordinated Neuronal Activity Enhances Corticocortical Communication." Neuron 87.4, pp. 827–839. pmid: 26291164.
- ¹²⁴⁶ Zhang, Xiaoxing, Wenjun Yan, Wenliang Wang, Hongmei Fan, Ruiqing Hou, Yulei Chen,
 ¹²⁴⁷ Zhaoqin Chen, Chaofan Ge, Shumin Duan, Albert Compte, and Chengyu T Li (2019).
 ¹²⁴⁸ "Active Information Maintenance in Working Memory by a Sensory Cortex". *eLife* 8,
- 1249 e43191.
- 1250 Zhu, Jia, Qi Cheng, Yulei Chen, Hongmei Fan, Zhe Han, Ruiqing Hou, Zhaoqin Chen, and

1251 Chengyu T. Li (2020). "Transient Delay-Period Activity of Agranular Insular Cortex

- ¹²⁵² Controls Working Memory Maintenance in Learning Novel Tasks". Neuron 105.5, 934–
- 1253 946.e5.



Figure 1 - Supplement 1. Anatomical details of the mouse cortex. (A). Connectivity matrix depicting cortico-cortical connections between 43 cortical areas. Areas are sorted according to their hierarchy. (B). The raw PV cell density for each cortical area (Y axis), with areas sorted (X axis). Each area belongs to one of five modules, shown in color (see also Fig. 1). (Harris et al. 2019). (C). Neuron density for each cortical area with same sorted order as (B). The data is from Erö et al. 2018.



Figure 2 - Supplement 1. Example simulation for different sensory modalities. The simulation protocol is the same as the default one in Fig. 2, except that the external input is applied to primary sensory areas related to two other sensory modalities: somatosensory and auditory. (A). The activity of 6 selected areas during the working memory task is shown. A somatosensory input of 500ms is applied to primary somatosensory area SSp-bfd, which propagates to the rest of the large-scale network. (B). Similar to the simulation where a primary visual area is stimulated (Fig. 2D), delay period firing for somatosensory stimulation is positively correlated with cortical hierarchy (r = 0.89, p < 0.05). (C). Delay period firing rate is moderately correlated with PV cell fraction (r = -0.4, p < 0.05). (D), (E) and (F) are similar to (A), (B) and (C) except that the input is given to primary auditory area AUDp. (E). Delay period firing rate is moderately correlated with cortical hierarchy (r = 0.89, p < 0.05). (F). Delay period firing rate is moderately correlated with cortical hierarchy (r = 0.89, p < 0.05).



Figure 3 - Supplement 1. Dependence of persistent activity on inhibitory model parameters (A). The maximum firing rate of all areas depends on the constant PV cell fraction in models without a gradient of PV. Average PV cell fraction from the anatomical data is shown as an orange dot. (B). Same as (A), except for the number of areas showing persistent activity. (C). Firing rate during the delay period for local circuits without long-range projections as a function of base inhibitory strength. If the base inhibitory strength is larger than a threshold (0.155, marked by the dashed line), none of the areas show independent persistent activity.



Figure 5 - Supplement 1. Anatomical data of thalamus and cortical connectivity. (A). Connectivity matrix of corticothalamic connections: 43 cortical areas to 40 thalamic areas. (B). Connectivity matrix of thalamocortical connections: 40 thalamic areas to 43 cortical areas.



Figure 6 - Supplement 1. Details of cell type-specific connectivity measures. (A). The matrix of cell type projection coefficients between cortical areas. The cell type projection coefficient is given by the formula $k_{cell} = m_{ij} - PV_i(1 - m_{ij})$. (B). The matrix of connectivity strengths, modified by cell type projection coefficient between cortical areas. The modified connectivity strength is given by $\tilde{W}_{ij} = (m_{ij} - PV_i(1 - m_{ij}))W_{ij}$.



Figure 6 - Supplement 2. Cell type-specific eigenvector centrality measures are better at predicting firing rate patterns than raw eigenvector centrality measures. The analysis is the same as in Fig. 6, where we compared cell type-specific input strength and raw input strength. Eigenvector centrality (EC, eigencentrality) of area i is the ith element of the leading eigenvector of the connectivity matrix. It quantifies how many areas are connected with the target area i and how important are these neighbors. Details are in the Methods section. (A(i)). Delay period firing rate (orange) and eigenvector centrality for each cortical area (blue). (A(ii)). Eigenvector centrality does not show a significant correlation with delay period firing rate for areas showing persistent activity in the model (r = 0.24, p = 0.29). (A(iii)). Eigenvector centrality cannot be used to predict whether an area shows persistent activity or not (prediction accuracy = 0.46). (B(i)). Delay period firing rate (orange) and cell type-specific eigenvector centrality for each cortical area (blue). (B(ii)). Cell type-specific eigenvector centrality has a strong correlation with the firing rate of cortical areas showing persistent activity (r = 0.94, p < 0.05). (B(iii)). Cell type-specific eigenvector centrality predicts whether an area shows persistent activity or not (prediction accuracy = 0.79). (C). Comparison of the correlation coefficient r for raw eigenvector centrality and cell type-specific eigenvector centrality in predicting delay firing rate. Raw input strength and cell type-specific input strength are also included for comparison. (D). Comparison of the prediction accuracy for raw eigenvector centrality and cell type-specific eigenvector centrality. Raw input strength and cell type-specific input strength are also included for comparison.



Figure 6 - Supplement 3. Sign only input strength measure and noPV input strength measure predict firing rate well. (A1). Delay period firing rate (orange) and sign only input strength for each cortical areas. (A2). Sign only input strength has a strong correlation with delay period firing rate of cortical areas showing persistent activity. (r = 0.90, p < 0.05) (B1). Delay period firing rate (orange) and noPV input strength for each cortical areas. (B2). noPV input strength has a strong correlation with delay period firing rate of cortical areas showing persistent activity (r = 0.90, p < 0.05).



Figure 7 - Supplement 1. Multiple-area inhibition experiments demonstrate the relative importance for core and readout areas in maintaining network-level persistent activity. (A). The x-axis shows readout areas that are inhibited as part of a pair (blue), triplet (orange), or quadruplet (green). For any given readout area A, the y-axis shows the average firing rate of all cortical areas that exhibit persistent activity when A was inhibited as part of the inhibited pair (triplet, quadruplet). The decrement in delay period activity is stronger as more areas are inhibited. (B). Bar plots showing the average firing rate of the network after inhibition of pair-wise combinations of core and readout areas. For example, the bar plot for 'readout-readout' is the average firing rate for all readout-readout areas pairs and is corresponding to the blue curve in (A). Dashed line in (A) and (B) denotes a threshold below which we consider an 'inhibition effect' to be significant. (C). Delay period firing rates as a function of hierarchy after inhibition of all core areas during the delay period. Although some readout areas show persistent firing, there is 48% decrement in average firing rate.



Figure 8 - Supplement 1. Cell type-specific loop strengths (Length 3 loops) are also better at predicting firing rate patterns than raw loop measures. Loop strengths (length 3 loops or L3) is calculated using similar method as loop strengths (length 2 loops). The only difference is we considered loops with length 3 (eg. A1->A2->A3->A1). The analysis is the same as in Fig. 7, where we compared cell type-specific loop strengths (length 2 loops) and raw loop strengths. (A(i)). Loop strength (blue) is plotted alongside Core Areas (orange), a binary variable that takes the value 1 if the area is indeed a Core Area, 0 otherwise. Loop strength is normalized to a range of (0, 1) for better comparison. (A(ii)). A high loop strength value does not imply that an area is a Core Area. (B(i)). Same as (A), but for cell type-specific loop strength. (B(ii)). High cell type-specific loop measures predicts that an area is a Core Area (prediction accuracy = 0.95). Same as (A), but for cell type-specific loop strength.



Figure 8 - Supplement 2. (A1). Relationship between core areas (orange) and length 2 (sign only) loop strength. Areas are sorted according to their hierarchy. Whether an area is a core area is represented in either 0 or 1. (A2). High loop strength is a good predictor of whether an area is a core area. Blue curve shows the logistic regression analysis used to differentiate the core areas versus non core areas (prediction accuracy = 0.83). (B1) and (B2). Same as (A1) and (A2), but with length 3 sign only loop strength. Length 3 sign only loop strength does not show a positive relationship with core areas (prediction accuracy = 0.83) (C1). (C2). Same as (A1) and (A2), except for comparing whether an area is a core area (orange) and length 2 noPV loop strength. Length 2 noPV loop strength predicts the core areas. prediction accuracy = 0.90 (D1). (D2). Same to A1 and A2, except for comparing whether an area is a core area (orange) and length 3 noPV loop strength. Length 3 noPV loop strength does not show a positive relationship with core areas area is a core area is a core area (orange) and length 3 noPV loop strength. Length 3

Supporting literature
(Kopec et al. 2015; Guo et al. 2014; Li et al. 2016b),
(Inagaki et al. 2019; Erlich et al. 2011; Guo et al. 2017),
(Gilad et al. 2018; Gao et al. 2018),
(Wu et al. 2020; Voitov and Mrsic-Flogel 2022)
(Liu et al. 2014; Schmitt et al. 2017),
(Bolkan et al. 2017)
(Wu et al. 2020)
(Harvey et al. 2012; Voitov and Mrsic-Flogel 2022)
(Zhu et al. 2020)
(Gilad et al. 2018)
(Pinto et al. 2019)
(Egorov et al. 2002)
(Zhang et al. 2019; Wu et al. 2020)
(Guo et al. 2017)
(Schmitt et al. 2017; Bolkan et al. 2017)
(Kopec et al. 2015)
(Gao et al. 2018)

Table 1: Supplementary experimental evidence. The listed literature include experiments that provide supporting evidence for working memory activity in cortical and subcortical brain areas in the mouse or rat. These studies show either that a given area is involved in working memory tasks and/or exhibit delay period activity. Area name corresponds to what has been reported in the literature. Some areas do not correspond exactly to the names from the Allen common coordinate framework.

MAIN TEXT: DISTRIBUTED WORKING MEMORY IN THE MOUSE BRAIN

Parameter	Description	$\mathbf{Task}/\mathbf{Figure}$	Value
	Cortical circuit parameters		
$ au_{NMDA}$	NMDA synapse time constant	All figures	60 ms
τ_{GABA}	GABA synapse time constant	All figures	5 ms
$ au_{AMPA}$	AMPA synapse time constant	All figures	2 ms
τ_{rates}	neuron time constant	All figures	2 ms
τ_{noise}	noise time constant	All figures	2 ms
a, b, d	parameters in excitatory F-I curve.	All figures	140 Hz/nA, 54 Hz, 308 ms
g_I, c_1, c_0, r_{0I}	parameters in inhibitory F-I curve.	All figures	4, 615 Hz/nA, 177 Hz, 5.5 Hz
γ	parameters in NMDA excitatory synaptic equations.	All figures	1.282
γ_I	parameters in GABA synaptic equations.	All figures	2
γ_A	parameters in AMPA excitatory synaptic equations.	All figures	2
$q_{E,self}$	local self excitatory connections	Figures 1-8	0.4 nA
$q_{E,cross}$	local cross population excitatory connections	All figures	10.7 pA
q_{IE}	local E to I connections	All figures	0.2656 nA
QEI 0, QEI scalina	local I to E connection strength	All figures	0.192 nA, 0.83
q _{II} ₀ , q _{II} scaling	local I to I connection strength	All figures	0.105 nA, 0.714
In A. InB	background current for excitatory neurons	All figures	0.305 nA
-0A, -0B	background current for inhibitory neurons	All figures	0.26 nA
σ_{A}, σ_{B}	standard deviation of excitatory noise current	All figures	5 pA
σ_A, σ_B	standard deviation of inhibitory noise current	All figures	0 pA
r_{0E}	background current for excitatory neurons	All figures	5 Hz
rot	background current for excitatory neurons	All figures	5 5 Hz
/ 01	long range E to E connection strength	Figures $1 2 3 4 6 7 8$	0.1 nA
PEE UID	long range E to L connection strength	Figures 1, 2, 3, 4, 6, 7, 8	0.167 nA
B	parameters in $m_{\rm eff}$	All figures	2 42
k.	parameters for scaling the connectivity matrix	All figures	0.3
α_{scale}	parameters for estimation of hierarchy	All figures	133 - 0.22
a_h, p_h	ovtornal stimulus strongth	All figures	0.5 n
I stim I	external input to inhibitory neurons	All figures	5 n A
T_{inh}	external input to inhibitory neurons	All figures	
T_{on}	stimulus start time	All former	2 S
T_{off}	simulation time for each trial	All figures	2.5 S
1 trial	simulation time for each trial	All ngures	10 s
at	simulation time step	All ngures	0.5 ms
	Thalamocortical network	Dimuna E	0.01 ~ A
μ_{EE}	long range E to E connection strength	Figure 5	0.01 nA
μ_{IE}	long range E to I connection strength	Figure 5	0.0167 nA
g_{ct}	cortico thalamic connections strength	Figure 5	0.32 nA
$g_{E,tc}$	thaiamo-cortical connections to excitatory neurons	Figure 5	0.6 nA
$g_{I,tc}$	thalamo-cortical connections to inhibitory neurons	Figure 5	1.38 nA
	Simulation of multiple attractors		
μ_{EE}	long range E to E connection strength	Figure 9	0.01, 0.02, 0.03, 0.04, 0.05 nA
μ_{IE}	long range E to I connection strength	Figure 9	0.0167, 0.033, 0.05, 0.066, 0.083 n.
$g_{E,self}$	local self excitatory connections	Figure 9	0.4, 0.41, 0.2, 0.43, 0.44 nA

 $Table \ 2: \ Parameters \ for \ numerical \ simulations$

1254