

Data-based large-scale models provide a window into the organization of cortical computations

Guozhang Chen*, Franz Scherr*, Wolfgang Maass¹

Institute of Theoretical Computer Science, Graz University of Technology,
Inffeldgasse 16b, Graz, Austria

* Contributed equally.

¹ To whom correspondence should be addressed; E-mail: maass@igi.tugraz.at.

Abstract

Understanding how the brain solves demanding computational tasks is one of the most exciting scientific challenges of our times. So far, recurrently connected artificial neural network models (RANNs) were primarily used to reverse-engineer brain computations. We show that it is now also feasible to reverse-engineer computations of detailed data-based large-scale models of cortical microcircuits. Furthermore, results of these analyses produce hypotheses that can readily be tested in biological experiments since they clarify from which neurons one needs to record and what type of information can be expected at specific time points during a trial. We apply this approach to a demanding visual processing task that has often been used in mouse experiments. Both the cortical microcircuit model and RANNs can solve this task as well as the mouse. But the resulting network dynamics matches only for the cortical microcircuit model experimental data on the sparseness of network activity and the impact of individual neurons on the network decision. Reverse-engineering of the computation in the cortical microcircuit model suggests that a particular subset of neurons causes a bifurcation of the network dynamics that triggers the network decision. Altogether, our results introduce a new type of neural network model for brain computations.

Short title: Reverse engineering of cortical computations

Teaser: Large-scale modeling on supercomputers provides new tools for understanding how the brain computes

1 Introduction

We tackle a central problem in computational neuroscience: How do neural networks of the neocortex compute? A major insight into brain function was the discovery that the mammalian neocortex is in first approximation a continuous 2D sheet consisting of rather stereotypical local circuits that consist of many different types of neurons that are located in 6 parallel layers (laminae) with numerous vertical connections between these layers and primarily local horizontal connections (Mountcastle 1998; Douglas and Martin 2004; Harris and Shepherd 2015). This architecture offers hope that one can understand how the neocortex solves demanding computational tasks by understanding the organization of computations in a representative patch of the 2D sheet that makes up the neocortex. Such representative patches are

35 usually referred to as cortical microcircuits, in spite of the fact that models for them are typically quite
36 large in order to provide a representative picture of interconnection patterns for many different types
37 of neurons. Intense research in several labs on the anatomy of cortical microcircuits (Mountcastle 1998;
38 Thomson and Lamy 2007; Markram et al. 2015) and subsequent further work at the Allen Institute has
39 recently culminated in a detailed publicly available model for a cortical microcircuit consisting of 51,978
40 neurons (Billeh et al. 2020). We refer to this model in the following as the Billeh et al. model.

41 We address the fundamental question how this cortical microcircuit model computes. More precisely,
42 since it is a model for a patch of V1 and comes together with a model for the LGN that serves as
43 a gateway for visual input from the retina to V1, we examine how the model of Billeh et al. solves a
44 demanding computational task that has frequently been used in experimental studies for visual processing
45 in the mouse (Garrett et al. 2020; Joshua H. Siegle et al. 2021): the visual-change-detection task. In
46 this task the subject receives a long sequence of natural images, with intermediate phases where just a
47 gray screen is shown. The task is to report after each image presentation whether it differs from the
48 previously shown image. This is a really demanding computational task since complex natural images
49 are shown. Furthermore, the computational performance is tested for novel images that never occurred
50 during training. Since a neural network is not likely to be able to retain all pixel values of the preceding
51 image in its working memory, it has to adopt a more sophisticated strategy that amounts to extracting
52 and retaining features of natural images that are generally useful for telling images apart, even for novel
53 images.

54 We trained both the data-based model of Billeh et al. and a generic RANN of the same size with
55 the same training method, stochastic gradient descent, to solve this computational task. Successful
56 application of stochastic gradient descent is less standard for a network of spiking neuron models that
57 have been fitted to biological data. But it can be made to work with a suitable modification of BPTT
58 (backpropagation through time) with pseudo-derivatives for spiking neurons as in (Bellec et al. 2018),
59 and further modifications from (Chen et al. 2022) for the more complex generalized leaky integrate-and-
60 fire (GLIF₃) neuron models from Billeh et al. that had been fitted to recordings from diverse neurons
61 from the Allen Brain Atlas Allen Institute 2018. Still, the application of stochastic gradient training to
62 such data-based cortical microcircuit models is computationally substantially more demanding than for
63 RANNs. But it becomes feasible through the use of advanced software, TensorFlow (Martin Abadi et al.
64 2015), and computer hardware (GPUs) that have been developed to support fast training of artificial
65 neural networks. Note that the connectivity structure of the cortical microcircuit model was not changed
66 through this training process, only the values of synaptic weights within a biologically reasonable range.

67 We adopt a biologically realistic convention for extracting the network decision for each image, change or
68 no-change, from the model: Network decisions have to be reported by a rather small subset of pyramidal
69 cells on layer 5 of the model that represent projections of V1 to subcortical targets. This modeling
70 convention is supported by experimental data of (Houweling and Brecht 2008; Marshell et al. 2019) and
71 others which showed that stimulation of just 1 or 2 pyramidal cells on layer 5 suffices for triggering a
72 behavioral response. This readout convention has a substantial impact on the computational analysis of
73 a cortical microcircuit model, since the customarily used linear readout from all neurons in the network
74 tends to mask the computational contribution of the network itself, as we will show.

75 We find that the computation for the visual-change-detection task achieves in the model a similar perfor-
76 mance as in-vivo. Furthermore, each computation engages the network in a way that has been reported
77 in numerous experimental data on cortical computation but which provides a stark contrast to typical
78 computations in generic recurrent artificial neural network models (RANNs): Neural activity is very
79 sparse, and therefore implemented in a very energy-efficient manner, with most neurons firing mostly
80 during a rather short phase within a trial. Furthermore, the temporal order of their peak activity is
81 different for different task conditions. In addition, a surprisingly small subset of neurons extracts from
82 the currently presented image the information whether it agrees with the previous one, and controls the

83 bifurcation of the trajectory of network states, thereby triggering the network decision. The sensitivity
84 of the network decision to the activity of very small subsets of neurons is another characteristic feature
85 of computations of biological neural networks (Houweling and Brecht 2008; Doron et al. 2014; Marshel
86 et al. 2019; Dalglish et al. 2020; Doron et al. 2020) which is reproduced by the model of Billeh et al. but
87 not by RANN control models of the same size that are trained for the same task. This result provides
88 evidence that the model of Billeh et al. operates in a critical regime, in spite of its very sparse activity,
89 where it is highly sensitive to even a few spikes of individual neurons in the network. A closer analysis
90 of these pivot neurons reveals that most of them have slowly changing internal variables, which most
91 biological neurons have according to the Allen Brain Atlas (Allen Institute 2018), that provide implicit
92 information about the preceding image while the network processes the current one.

93 We expect that the analysis and supercomputing methods that we present pave the way for research on
94 a new generation of models in computational neuroscience that integrate a substantially larger body of
95 experimental data, and can reproduce features of cortical computations that are difficult to reproduce in
96 artificial neural network models. Also, their predictions can be tested more directly through biological
97 experiments because neurons in the model can be immediately related to the specific types and laminar
98 locations of neurons that are examined in wetlab neuroscience. Going forth and back between detailed
99 modeling and biological experiments is likely to be needed to elucidate the computational function of
100 the neocortex. Since the neocortex achieves its superior computational performance with an energy
101 consumption that is by several orders of magnitudes lower than that of current computer hardware,
102 an understanding how the cortex is able to combine energy-efficient very sparse activity with superior
103 computational performance is likely to have also important technological implications.

104 2 Results

105 2.1 Setting up a data-based V1 model for reverse engineering of cortical 106 computations

107 The V1 model of (Billeh et al. 2020) represents one of the most comprehensive efforts to integrate the
108 available experimental data on the anatomy and neurophysiology of area V1 in mouse that is currently
109 available (Fig. 1A-C). It distinguishes 17 different neuron types (listed in each row and column of Fig. 1B).
110 These neuron types are further split into 111 different variations based on response profiles of individual
111 neurons from the Allen Brain Atlas (Allen Institute 2018), to which generalized leaky integrate-and-fire
112 (GLIF₃) neuron models with 3 internal variables had been fitted. These neuron models have in addition
113 to the membrane potential two further internal variables that model after-spike currents (Fig. 1D).
114 The resulting model for a patch of V1 receives visual input from an LGN model that consists of 2,589
115 filters (Billeh et al. 2020) that had been fitted to experimental data. This LGN model produces input
116 currents to neurons of the V1 model in a retinotopic and lamina-specific manner (Billeh et al. 2020). We
117 will refer in the following to this model of V1 in conjunction with the LGN model of (Billeh et al. 2020)
118 as the V1 model of Billeh et al.

119 We employed a data-driven noise model based on experimental data from area V1 of the awake mouse
120 (Stringer et al. 2019). This noise model was not present in (Billeh et al. 2020) and had subsequently
121 been introduced in (Chen et al. 2022). It models both quickly changing forms of noise and slower forms
122 of noise that contribute to experimentally found trial-to-trial variability (Methods). The resulting V1
123 model was trained to solve the visual-change-detection task, which has frequently been used in mouse
124 experiments (Garrett et al. 2020; Joshua H. Siegle et al. 2021). A sequence of natural images was
125 presented to the model, interleaved by periods without visual input. The subject had to report whenever

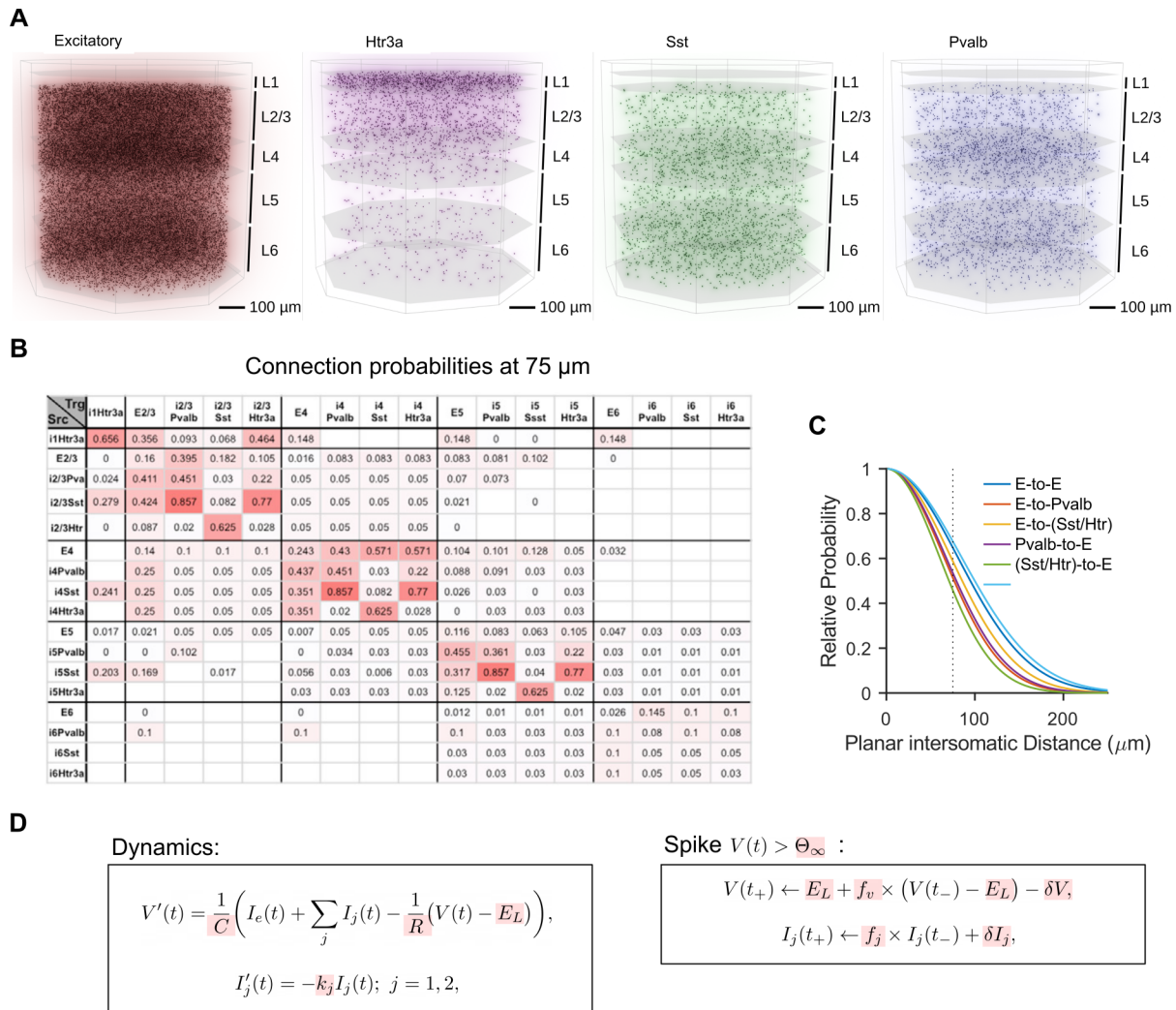
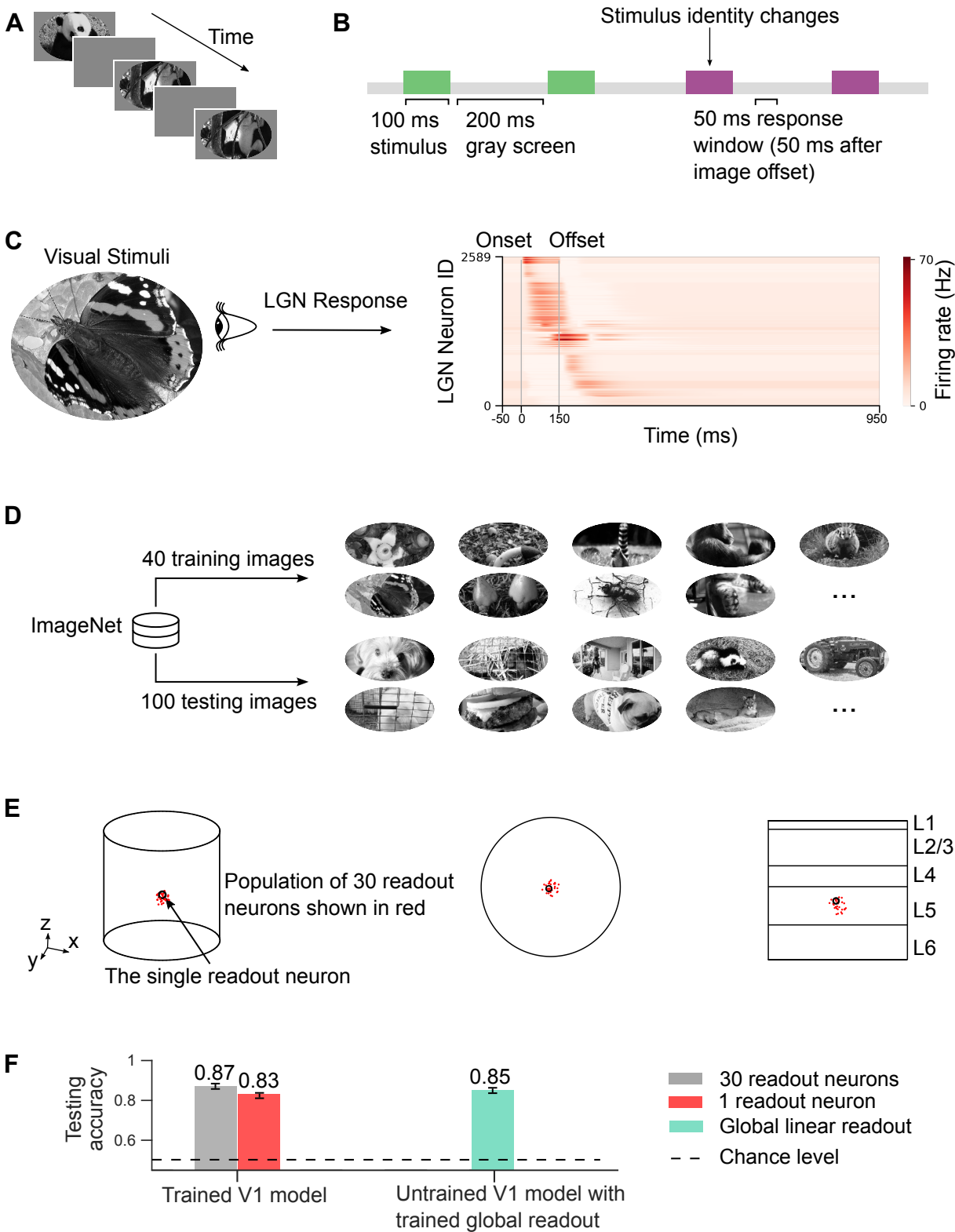


Figure 1: Structure and components of the data-based cortical microcircuit model of (Billeh et al. 2020). (A) Side view of the 3D architecture of the model, which consists of with 51,978 neurons from 1 excitatory and 3 inhibitory neuron classes (Htr3a, Sst, Pvalb) in a column with 800 μm diameter. These neurons are distributed over 5 laminar sheets (L1 to L6, where L2 and L3 are lumped together). (B) Base connection probabilities between these neuron classes on different laminae, which are valid if the horizontal distance between neurons is at most 75 μm. (C) Scaling function for connection probabilities in dependence of the class to which the pre- and postsynaptic neuron belongs. The probability of a synaptic connection between the two neurons is the product of the base connection probability from panel (B) and this scaling function. (D) Main equations defining the GLIF₃ neuron models of (Billeh et al. 2020) with 3 internal variables. Assignments to their parameters (highlighted in red) define the 111 neuron models of the networks, based on experimental data from the Allen Brain Atlas (Allen Institute 2018).



(caption next page)

Figure 2 (previous page): Visual-change-detection task. (A) A sequence of natural images is presented, interleaved with gray screen periods. (B) The visual-change-detection task requires to give during a 50 ms response window that begins 50 ms after image offset an output signal whenever the current image is different from the preceding one. (C) We used the LGN model from (Billeh et al. 2020) to transform visual stimuli into the outputs of 2,589 filters that model firing rates of LGN neurons and are connected to neurons in the V1 model according to data-based rules. These outputs provide input currents to V1 neurons using data-based rules. (D) We used 40 natural images for training and a separate set of 100 natural images for testing. Images were drawn from ImageNet (Russakovsky et al. 2015) and presented in grayscale because of the insensitivity of LGN model of Billeh et al. to color. To ensure the robustness of our findings, we trained 10 V1 models using 10 different sets of training data (each with its own 40-image set). We found that all results presented below were consistent across all models, demonstrating that our conclusions are not reliant on specific training datasets or coincidences. (E) The network was trained to provide the network decision through projection neurons within the network, either a single (black circle) or 30 (red dots) randomly selected pyramidal cells in layer 5 (within a sphere of 55 μm) depending on the experiment. Their task was to report an image change through an increased sum of firing rates during the response window. (F) Testing accuracy of the trained V1 model of Billeh et al. for reporting network decisions through increased firing activity of 1 or 30 projection neurons is shown on the left. One sees that the chosen number of projection neurons has little impact on the network performance. The green bar on the right shows a control result for the case when one uses instead of projection neurons from within the network a global linear readout from all neurons in the network. One sees that in case it is not even necessary to train the V1 model for the task: Training of the weights of the artificial global linear readout suffices, indicating that the computation within the V1 model can be effectively masked by such a global linear readout. The error bars, which are small (< 0.01), represent the standard error of the mean (SEM) across 10 models with different training datasets.

126 the most recently presented image differed from the one before (Fig. 2A-C). We randomly selected a pool
127 of 140 natural images from the Imagenet dataset (Deng et al. 2009) that we used as network inputs. We
128 used 40 of them for training, similar to the biological experiments of (Garrett et al. 2020), and 100 of
129 them for testing (Fig. 2D).

130 Projections neurons on layer 5 of cortical microcircuits extract computational results of the microcircuit
131 and transmit them to other brain areas, thereby triggering behavioral responses (Harris and Shepherd
132 2015; Marshel et al. 2019). Therefore we selected a set of pyramidal cells on layer 5 as readout neurons of
133 the V1 model (Fig. 2E). These readout neurons had the task to fire during a specific response window after
134 an image presentation if the preceding image differed from the one that had been presented before that.
135 The size and spatial distribution of this readout population had little impact on the results (Fig. 2F).
136 For simplicity, we used a single readout neuron unless stated otherwise.

137 A more common output convention in modeling brain computations is to use an external readout neuron
138 that receives synaptic input from all network neurons. We found that this convention is not suitable
139 for probing the computational capability of a network model. First, such global readout neurons that
140 receive synaptic inputs from all neurons in a large patch of the neocortex have not been found in the
141 brain. Secondly, a global linear readout neuron that receives synaptic inputs from a large set of neurons
142 in a recurrent neural network through trained synaptic weights is a too-powerful device that masks
143 the computational contribution of the recurrent neural network. Instead, the recurrent neural network
144 plays in this case just the role of a liquid or reservoir (Maass et al. 2002; Maass and Markram 2004).
145 Concretely, if one takes the V1 and LGN model as defined in (Billeh et al. 2020), without changing any of
146 its synaptic weights or other parameters, and just trains a linear readout from all of V1 neurons for the
147 visual-change-detection tasks, one gets already a very high average accuracy of 0.85 (see the rightmost
148 bar in Fig. 2F).

149 We trained all synaptic weights from the LGN to V1 and within the V1 model using stochastic gradient
150 descent for a suitable loss function that assumed a low value only when the readout neuron(s) fired within

151 a short response window 50 ms after the presentation of an image in case that this image differed from
152 the preceding one (Methods). We included regularization terms similarly as in (Chen et al. 2022) in the
153 loss function in order to keep the firing activity of the network in a biologically realistic sparse firing
154 regime. Synaptic connectivity was not changed through this training process. We also did not allow
155 synaptic weights to change their sign, thereby obeying Dale’s law. In particular, the average firing rate
156 after training was 3.86 Hz. Hence the model computed in an energy-efficient sparse firing regime. The
157 new values of synaptic weights after training remained in a biologically reasonable range, see Fig. S1
158 and S2.

159 The trained V1 model achieved high performance for the visual-change-detection task (Fig. 2F), lying
160 in the same range as the performance achieved by mice (Garrett et al. 2020). The model was also able
161 to generalize well, achieving almost the same performance for images that were not used during training
162 (Fig. 2F). Hence the model had acquired a network algorithm that generalized, i.e., was not constrained
163 to a particular set of previously seen images.

164 2.2 The data-based V1 model reproduces characteristic features of cortical 165 computations

166 Simultaneous recordings from large numbers of neurons in the awake brain show that neural networks of
167 the neocortex exhibit a peculiar type of network dynamics that is rarely seen in artificial neural networks:
168 Most neurons fire only during a rather short time window during a trial, see e.g. (Driscoll et al. 2017)
169 and Fig. 2 of (Koay et al. 2022). Furthermore, different neurons have different preferred firing times,
170 and the relative order in which they fire depends on the task condition and the sensory input. Hence,
171 similarly as in synfire chain models (Abeles et al. 2004), the network activity has a prominent sequential
172 character, but in contrast to synfire chain models only a very small fraction of neurons is active during
173 each segment of the sequence.

174 Since this type of network activity is hard to reproduce in neural network models, we wondered whether
175 the data-based V1 model would be able to do that. We considered trial-averaged neural activity as in
176 (Koay et al. 2022) and followed the same routine to order the neurons according to the time of the
177 peak activity. Furthermore, as in their data analysis we normalized the firing activity, averaged over 200
178 trials, of each neuron over the 300 ms of the computation on each image, consisting of a segment of the
179 continuous input- and processing stream that contained the 100 ms of the image presentation, a 50 ms
180 delay, the subsequent 50 ms of the network response window, and a 100 ms delay before the presentation
181 of the next image, marked at the top of Fig. 3 A-C. We found that the neural activity in the V1 model was
182 indeed very similar to that in the experimental data. Fig. 3 A-C show that most neurons of the V1 model
183 participated in the network computation, but focused their firing activity to a very short segment of the
184 300 ms time window. Furthermore, as in the experimental data, the relative order of these short segments
185 of high activity depended on whether the currently processed image was the same as the preceding one
186 or not, see Fig. 3. The neurons are plotted in panel C in the same order as in panel A, but in panel A for
187 the change condition and in panel C for the no-change condition. Panel B, where we have ordered the
188 neurons according to their preferred firing time for the no-change condition, shows that the neurons have
189 also for the no-change condition a clear order of preferred firing times. But the order is different from that
190 for the change condition, as panel C shows: The resulting sequence is blurred and lacks the characteristic
191 thin-line pattern observed in Fig. 3A and B. In other words, the network employed temporal coding
192 through the relative timing of the peak-firing activity of neurons for distinguishing between these two
193 experimental conditions. Furthermore, this temporal coding was simultaneously expressed in all layers
194 of the laminar cortical microcircuit model. Note that this type of temporal coding, through the order to
195 peak activity of different neurons, takes place in spite of substantial noise both in the brain and in our
196 V1 model, see (Chen et al. 2022) for details, that substantially affects the timing of individual spikes in

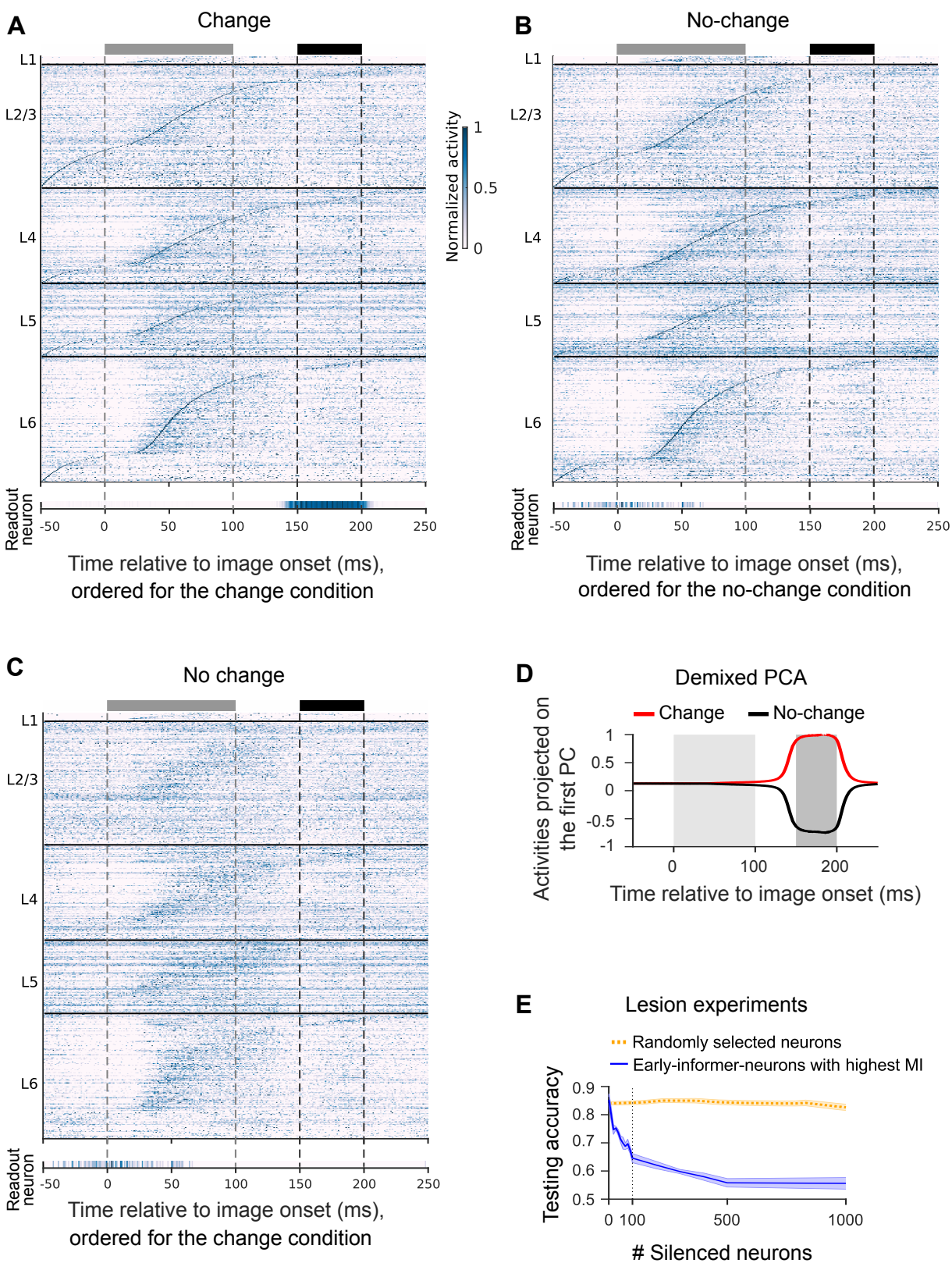
197 a single trial. The V1 model predicts that cortical microcircuits employ an even more refined type of
198 temporal coding: The order of peak activity also depends on the identity of the current image (Fig. S3).
199 This prediction needs to be tested through neurobiological experiments.

200 The experimental data of (Houweling and Brecht 2008; Doron et al. 2014; Marshel et al. 2019; Dalglish
201 et al. 2020; Doron et al. 2020) have elucidated another characteristic feature of cortical computations:
202 The network computation is surprisingly sensitive to the activity of a tiny fraction of neurons of these
203 very large networks, since stimulation of a few selected neurons could change the result of the network
204 computation, i.e., the behavioral response. This high sensitivity of the neocortex on the activity of
205 particular neurons is surprising insofar as it has to cope with a substantial amount of noise, see the
206 analysis and resulting noise model of (Chen et al. 2022) that we also used in this study. Since artificial
207 activation of a small set of neurons in a network is likely to activate also other neurons, a factor that
208 had not been quantified in these biological experiments, we used a slightly different paradigm for testing
209 the sensitivity of the network decision of the V1 model to the activity of a small set of neurons: We
210 silenced, rather than activated, selected neurons. Like in the experimental data on the neocortex, we
211 found that the result of the network computation is not equally sensitive to the activity of all of the
212 neurons, but that there are particular neurons that are pivotal for the network decision. In order to
213 avoid artifacts resulting from silencing of the readout neurons or neurons that directly activate them, we
214 focused on early-informer-neurons which produced the first information on whether the current image
215 was the same as the preceding one well before the response window, while the image was still processed.
216 Fig. 3D shows that the first information about this arises during the time window from 50 to 100 ms after
217 image onset, when the first information about the image reaches the V1 part of the model. Hence we
218 defined early-informer-neurons as neurons whose spike output contained already during this time window
219 substantial mutual information (MI, Methods) with the upcoming network decision. Fig. 3E shows that
220 while the network is not very sensitive to silencing of randomly selected neurons, silencing of just 100
221 early-informer-neurons has a drastic impact on the network performance. Also, the number of readout
222 neurons does not affect the conclusion of the lesion experiment (Fig. S4).

223 **2.3 RANN models cannot reproduce these characteristic features of cortical** 224 **computations**

225 Randomly connected recurrent networks of artificial neurons (RANNs) have commonly been used as
226 models for computations in cortical neural networks (Sussillo and Barak 2013; Sussillo et al. 2015; Yang
227 et al. 2019; Yang and X.-J. Wang 2020; Pollock and Jazayeri 2020). We show here that RANNs are not
228 able to reproduce the two previously discussed fingerprints of cortical computations: a short period of
229 peak activity for most neurons, a characteristic sequential order of this peak activity according to the
230 trial type, and the sensitivity of the network to the activity of small subsets of neurons. In order to
231 eliminate a possible impact of differences in the network size or training procedure, we trained through
232 stochastic gradient descent a randomly connected RANN with the same number of neurons and synapses
233 as the V1 model of (Billeh et al. 2020), for the same computational task. Furthermore, we used the same
234 preprocessor (the LGN model of (Billeh et al. 2020)) for transforming images into temporally dispersed
235 inputs to random subsets of neurons in the network. We used a standard neuron model from RANN
236 models for neural networks of the neocortex (Sussillo et al. 2015; Pollock and Jazayeri 2020): A non-
237 spiking neuron with tanh as activation function and a membrane time constant of 50 ms. For extracting
238 the network decision we used the same global linear readout from all neurons in the RANN as in these
239 paradigms.

240 This RANN model was after training able to perform the visual-change-detection task at a higher perfor-
241 mance level than the data-based V1 model and the subjects in the neurobiological experiments (Garrett
242 et al. 2020). But it could not reproduce the two previously discussed fingerprints of cortical computations.



(caption next page)

Figure 3 (previous page): Temporal organization of computations in the V1 model. (A) As in the experimental data, neural activity is sparse and exhibits a clear sequential organization with high temporal resolution. Shown are normalized average responses over 200 trials with the change condition but different images, with neurons ordered according to the time of their peak activity under the change condition. The gray and black bars at the top denote the image presentation and response windows, respectively. (B) Same as in (A), but for the no-change condition, with neurons ordered according to the time of their peak activity for the no-change condition. (C) The same data as in (B) for the no-change condition, but with neurons ordered as in (A). The resulting blurred sequence indicates that the order of peak activity of neurons is quite different for the change and no-change conditions. (D) Demixed principal components analysis. In order to visualize the formation of the network decision, we carried out demixed principal component analysis for trial-averaged network activity (Methods). Its projection onto the first principal component is shown. One sees that the network decision starts to emerge during the time window from 50 to 100 ms after image onset. The light and dark gray rectangles denote the window of image presentation and response, respectively. (E) Causal impact of specific neurons on the network decision: Task performance quickly decreases when early-informer-neurons are silenced (in the order of their MI with the network decision), see the blue curve. On the other hand, task performance is robust to silencing the same number of randomly selected neurons (dotted yellow curve). Both curves show average values for 10 V1 models where different sets of training data were used. The shaded area represents the SEM across 10 models.

243 Fig. 4 A and B show that neurons of the RANN do not have a similarly short time period of high activity
244 during the computation on the same task as the V1 model. Furthermore, the same analysis of the firing
245 order as in (Driscoll et al. 2017; Koay et al. 2022) does not reveal a substantial dependence of the order
246 of peak activity of neurons on the trial type, see Fig. 4C. We also tested a different way of plotting the
247 activity of the RANN where we did not normalize the activity of each neuron as in the data analyses of
248 (Driscoll et al. 2017; Koay et al. 2022). Resulting plots (Fig. S5) show that neurons in the RANNs have a
249 wide range of different activity levels. But the analysis of the role of temporal order as in (Driscoll et al.
250 2017), which was also employed in Fig. 3, did still not provide any indication that the temporal order of
251 peak activity in the RANN depended in a significant way on the trial type. Since the activity level of
252 neurons in the RANN depends on the weight of the regularization term in the loss function for stochastic
253 gradient descent, we repeated the training of the RANN with different weights of this regularization term
254 that controls the average activity of neurons in the network (Fig. S6); see Fig. S7 for the same results
255 without normalizing the activity of each neuron. Also in these controls, the neurons of the RANN do not
256 constrain their activity to a short time window like in the experimental data and the V1 model. The same
257 holds for a further control (Fig. S8) where we changed the threshold of the regularization term in the
258 loss function (Methods). These results suggest that it is not possible to reproduce in the RANN model
259 under common configurations the sparse neural activity with trial-type-dependent temporal sequences of
260 neural peak activity that has been found in the neocortex.

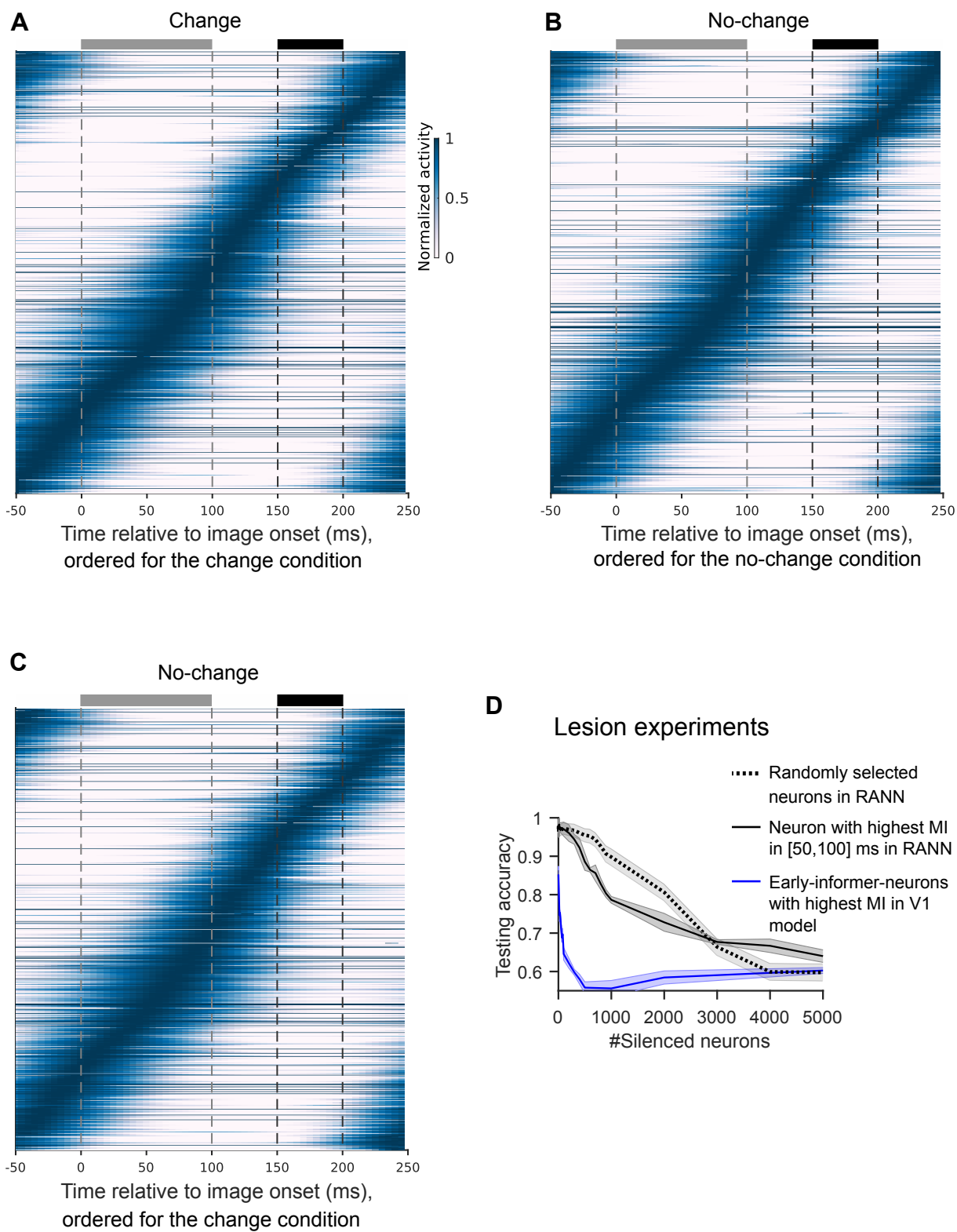
261 We also tested the sensitivity of the RANN to the activity of small sets of its neurons. We found that about
262 6,000 neurons need to be silenced in the RANN to reduce the accuracy of the network computation to
263 0.6 (Fig. 4D), a performance level which was reached in the V1 model according to Fig. 3D by silencing
264 just 200 neurons. Note that we silenced here RANN neurons in descending order of the MI of their
265 activity during the second half of an image presentation with the upcoming network decision, like in the
266 V1 model. The RANN has more high-MI neurons than the V1 model (Fig. S9). These results suggest
267 that the RANN is substantially less sensitive to the activity of small subsets of its neurons.

268 Altogether, these results suggest that the computation for the visual-change-detection task is organized
269 in the RANN quite differently than in the mouse brain and in the V1 model. In particular, the in-vivo
270 data and the V1 model suggest that neurons become “experts” for particular phases of a particular
271 computation and otherwise remain silent, reminiscent of mixture-of-experts models (Yuksel et al. 2012)
272 and hidden Markov models (Kappel et al. 2014). In contrast, information and impact on the network

273 output are distributed in the RANN over substantially larger subsets of neurons. This is likely to result
274 from the random connectivity of the RANN, which makes it harder to accumulate specific information in
275 specific parts of the network, the neuron models of the RANN, which do not induce its neurons to restrict
276 their activity to a particular phase of computation, and the use of a global readout from all neurons for
277 extracting the network decision.

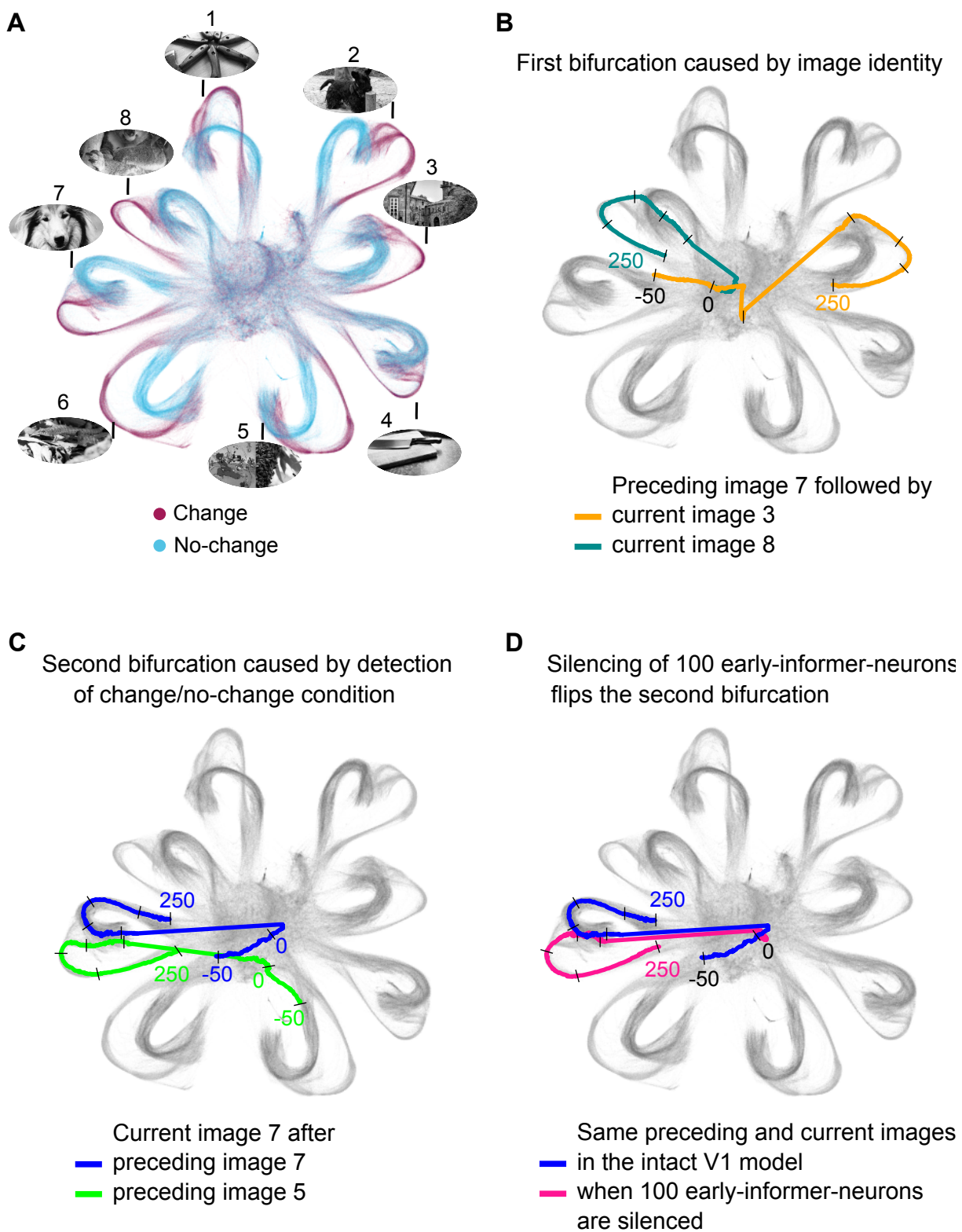
278 **2.4 Computational progress of the V1 model becomes visible as nested bi-** 279 **furcations of the network dynamics**

280 Analyses of results of simultaneous recordings from large numbers of neurons in the brain have shown that
281 low-dimensional projections of the high-dimensional network activity provide interesting links between the
282 network dynamics and the computations that it performs (Broome et al. 2006; Kato et al. 2015; Allen et al.
283 2017; Steinmetz et al. 2019). We wondered whether similar analyses could elucidate computations of the
284 V1 model. We embedded its activity vectors, defined for every ms by the low-pass filtered spiking activity
285 of its 51,978 neurons, into 2 dimensions with the help of PCA and a subsequent application of UMAP
286 (McInnes et al. 2018). The processing of each particular image gives rise to a bundle of trajectories, with
287 trial-to-trial variability resulting from the preceding images and noise within the network. Two nested
288 bifurcations in these bundles of trajectories mark the computational progress of the dynamical system, see
289 Fig. 5A. First, the trajectories of network states bifurcate during the first 50 ms of an image presentation
290 from the mid-region of the plotted state space and move into a region that is characteristic for the identity
291 of the currently presented image, see Fig. 5B. Afterward, between 50 and 100 ms after image onset, the
292 second bifurcation occurs in dependence on whether the current image is the same or different from
293 the previously presented image, see Fig. 5C. These modeling results provide concrete predictions for the
294 way how these computations are carried out by cortical microcircuits of the brain, viewed as dynamical
295 systems. Among various options on how these could compute, as discussed in (Rabinovich et al. 2006),
296 nested bifurcations of bundles of trajectories emerge as the clearest visible fingerprint of these network
297 computations.



(caption next page)

Figure 4 (previous page): Temporal organization of the same computations in the RANN model. (A) Neuronal activity in the equally large RANN model was plotted for the same network inputs and task condition as in Fig. 3A. Here, the activation regularization was not used in the RANN; in other words, the weight of activation regularization is 0 (Methods). (B) Same as in (A), but for the no-change condition as in Fig. 3B. (C) Same data as in panel (B), but with neurons ordered as in panel (A). In contrast to Fig. 3, little difference emerges between panels (B) and (C), indicating that the order of peak activity is less dependent on the task condition in the RANN. (D) Lesion experiments corresponding to those in the V1 model (Fig. 3E). The blue curve for the V1 model is the same as in Fig. 3E. One clearly sees that the network decision is substantially less sensitive to the activity of 100 neurons. The shaded areas represent SEM across 10 RANNs where different sets of training data are used.



(caption next page)

Figure 5 (previous page): Nested bifurcations of trajectories of network states provide links between network dynamics and network computations. Spiking activity of the 51,978 neurons was first filtered with an exponentially decaying kernel (time constant: 20 ms), and then projected onto its first 1,500 PCA dimensions (capturing 38% of the variance, which indicates that the network activity is quite high-dimensional). These data were then projected into 2 dimensions with the UMAP method. Each dot represents network activity at a particular ms of a 250-ms long fragment of the computation of the V1 model on a sequence of natural images. Short black bars mark 50-ms long subsections of network trajectories. **(A)** Network trajectories of the V1 model during the presentation of 8 (out of the 100) test images that had not been presented during training. Two colors indicate whether the current image was identical (blue) or different (brown) from the preceding image. One sees that the network state moves for each image into a different region of the state space, no matter whether it was the same or different from the preceding image. **(B)** The first bifurcation occurs according to the identity of the current image, highlighted here for the case where the trajectory starts in both cases from the same state, which largely results from the identity of the preceding image. This first bifurcation occurs within the first 50 ms after image onset. **(C)** The second bifurcation after image onset. It does not depend on the identity of the current image, but on whether it is the same or different from the preceding image. This 2nd bifurcation is more difficult to visualize since the trajectory arrives from two different regions that are characteristic of the identity of the preceding image. **(D)** A small set of early-informer-neurons is causal for the second bifurcation. Here two times the image 7 is shown, both in the intact model and when 100 early-informer neurons are silenced. One clearly sees that these neurons are causal for the second bifurcation occurring within 50 to 100 ms after image onset, since silencing them lets the trajectory flip to the bundle for the no-change condition. Silencing of these neurons also flips a trajectory for the no-change condition to the bundle for the change condition, see Fig. S10.

2.5 Causal relations between the activity of individual neurons and the network decision

The results of our lesion experiments in Fig. 3E indicate that the firing activity of a specific small set of neurons was causal for the network decision. These early-informer-neurons were distinguished by the fact that their firing activity during 50 to 100 ms after image onset contained substantial MI with the upcoming network decision. We show in Fig. 5D that their firing activity was also pivotal for the network dynamics. Silencing the 100 early-informer-neurons with the highest mutual information (MI) with the upcoming network decision flips the trajectory of the network dynamics at the 2nd bifurcation to another bundle, see the magenta curve in Fig. 5D. This suggests that these bundles of trajectories had a certain attractor quality for time-varying network states. The model received for both trajectories shown in this panel exactly the same network input, hence differences were only caused by the silencing of the 100 early-informer-neurons (with some further variance possibly caused by ongoing noise in the network). This was a no-change trial, as the blue curve in Fig. 5D indicates for the intact network (compare with the trajectory bundles in Fig. 5A). Silencing of these 100 neurons can also flip the trajectory of a change trial to the bundle of trajectories for no-change trials, see Fig. S10.

We then investigated what mechanism enabled these 100 early-informer-neurons to decide whether the currently presented image differed from the previously presented one. These neurons were primarily located in layers 2/3 and 4, see Fig. 6A. We selected 7 of them (Methods) from the 4 major neuron classes so that these had during the interval from 50 to 100 ms after image onset the largest MI with the subsequent network decision. Their firing activity is shown in Fig. 6B. One sees that they fired at different rates during the interval from 50 to 100 ms after image onset for the change and no-change conditions. These rate differences continued to be present during the subsequent delay and response window.

In order to determine the mechanism by which these neurons acquired their early information about the relationship between the identity of the current and preceding image we analyzed the dynamics of their internal variables. Fig. 6C depicts the time course of their internal variable with the largest time constant. One sees that this internal variable had in many trials for 5 of these neurons already at the beginning of the presentation of the current image (indicated by the beginning of the light-gray zone) a strongly negative value. In order to understand the computational role of these internal variables, we analyzed whether their strongly negative values at image onset provided information about the identity of the preceding image. The result of this analysis is plotted for the first 4 of these neurons, whose locations are marked in Fig. 7A, in Fig. 7B. One sees that for those 3 among these 4 neurons that had an internal variable with a large time constant, shown here in the 2nd to 4th column, this slow internal variable assumed its lowest values at image onset when a particular image (number 5) had been presented before. Note however that these neurons could not specialize in reporting this particular preceding image, because the network was tested on new test images that had not been shown during training. Therefore, in order to serve as early-informer-neurons, the image features that produced the lowest values of their internal variables at the beginning of the next image had to tile the image space. The neuron whose analysis is plotted in the first column has only short time constants, and the value of its internal variable is less characteristic for a particular preceding image. Hence it is likely to collect and transmit information that it receives from other early-informer-neurons. In order to further elucidate the causal relation between the network decision and preceding activity (or non-activity) of early-informer-neurons with and without long time constants, we separately silenced from each of these two neuron classes those which had the largest MI with the network decision. The result in Fig. 7C shows that the network performance is significantly reduced when we silence 100 neurons with long time constants, and more than 500 neurons with short time constants need to be silenced to produce a similar reduction of network performance.

The impact of early-informer-neurons on the firing or non-firing of the readout neuron was in general rather indirect, because they were in general not directly connected to the readout neuron (Fig. S11).

345 This had to be expected, because the response window came 100 ms after the critical period of early-
346 informer neurons (50 to 100 ms after image onset). Nevertheless, a direct impact of silencing the 100
347 early-informer-neurons on the membrane potential of the readout neuron can be detected. The result
348 is shown in Fig. 7D. One sees in the 2nd and 3rd row show that this silencing significantly moved the
349 membrane potential of the readout neuron during the response window, thereby explaining why their
350 silencing caused errors in the network decision.

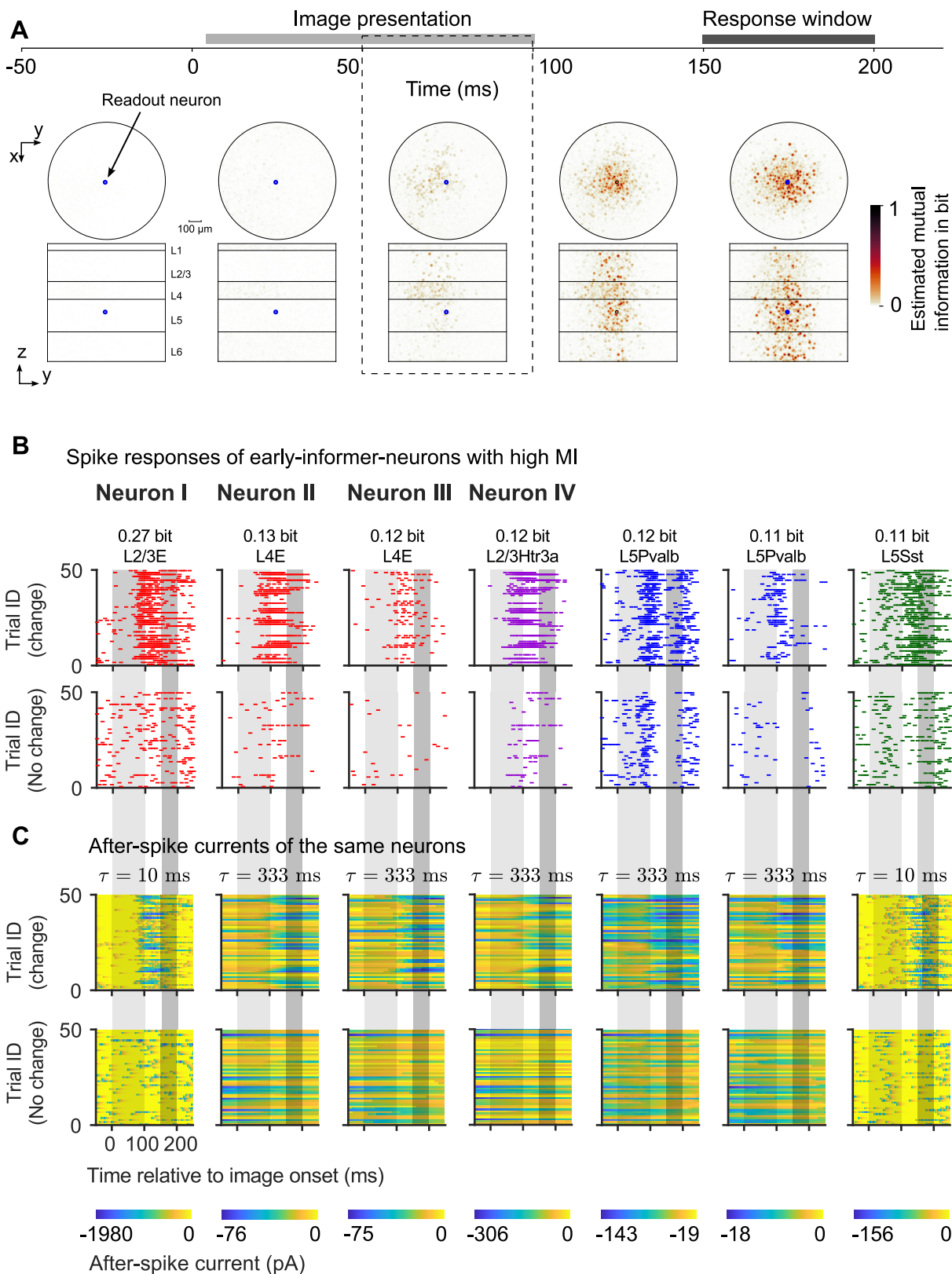
351 We have shown that the silencing of 100 neurons in the V1 model is able to change the result of a network
352 computation. This high sensitivity of the network output to the firing activity of individual neurons is
353 consistent with experimental data. More specifically, it has been shown that artificial activation of very
354 few neurons on layer 5 in the brain is able to switch a behavioral decision (Houweling and Brecht 2008;
355 Doron et al. 2014; Marshel et al. 2019; Dalglish et al. 2020; Doron et al. 2020). Note that it is difficult
356 to estimate how many further neurons were indirectly activated in these in-vivo experiments, hence it
357 remains open exactly how many neurons need to be manipulated in order to switch a network decision
358 of a cortical microcircuit.

359 Finally, we would like to emphasize that the computational analyses in Figs. 3, 5 and 6 were carried
360 out for a new set of images that had not been shown during training of the network. Hence our reverse
361 engineering has elucidated a generic computational network mechanism, rather than a mechanism that
362 only works for specific images. In particular, the diverse selectivity of the working memory of these
363 neurons shown in Fig. 6C in combination with their causal role for the network decision demonstrates
364 that their working memory specializations tile the image space, thereby enabling correct network decisions
365 for generic images.

366 3 Discussion

367 We have presented a new paradigm for modeling and analyzing computations in the neocortex. More
368 specifically, we have trained a detailed model for a patch of neocortex with a diameter of 800 μm to carry
369 out a demanding computational task that has often been used in mouse experiments: Deciding for a
370 sequence of natural images, interleaved with delays where no image is shown, whether the most recent
371 image was the same as the preceding one. The V1 model was able to solve this task after training on one
372 set of images, like the mouse, and also for never presented new images with high accuracy. This task is
373 a demanding computational task for any neural network since salient information needs to be extracted
374 from each image and stored, compared with information from the next image while simultaneously storing
375 information from this next image that will be needed to compare it with the subsequent image, and the
376 result of the comparison has to be reported at a time when none of the images are present. This task is
377 especially demanding for a neural network that has to cope with a substantial amount of internal noise,
378 which is the case both for neural networks in the brain and our model.

379 We have shown that a detailed model for a patch of V1 (Billeh et al. 2020), equipped with a data-based
380 and more challenging noise model from (Chen et al. 2022), can be trained with stochastic gradient descent
381 to solve this visual-change-detection task. We found that the resulting organization of its computation
382 differed in essential aspects from that of a RANN with the same number of neurons and synapses that has
383 been trained with stochastic gradient descent for solving the same task. In particular, the computations
384 of the data-based cortical microcircuit model revealed an exquisite temporal organization, where most
385 neurons focused their activity onto a particular phase of the computation, as seen in Ca-imaging data
386 from the awake brain (Driscoll et al. 2017; Koay et al. 2022). The V1 model also reproduced the finding
387 of (Driscoll et al. 2017) that this temporal order of peak activity depends on the trial type (in our case:
388 change or no-change), see Fig. 3. In addition, the rank order of peak activity depended in the data-based



(caption next page)

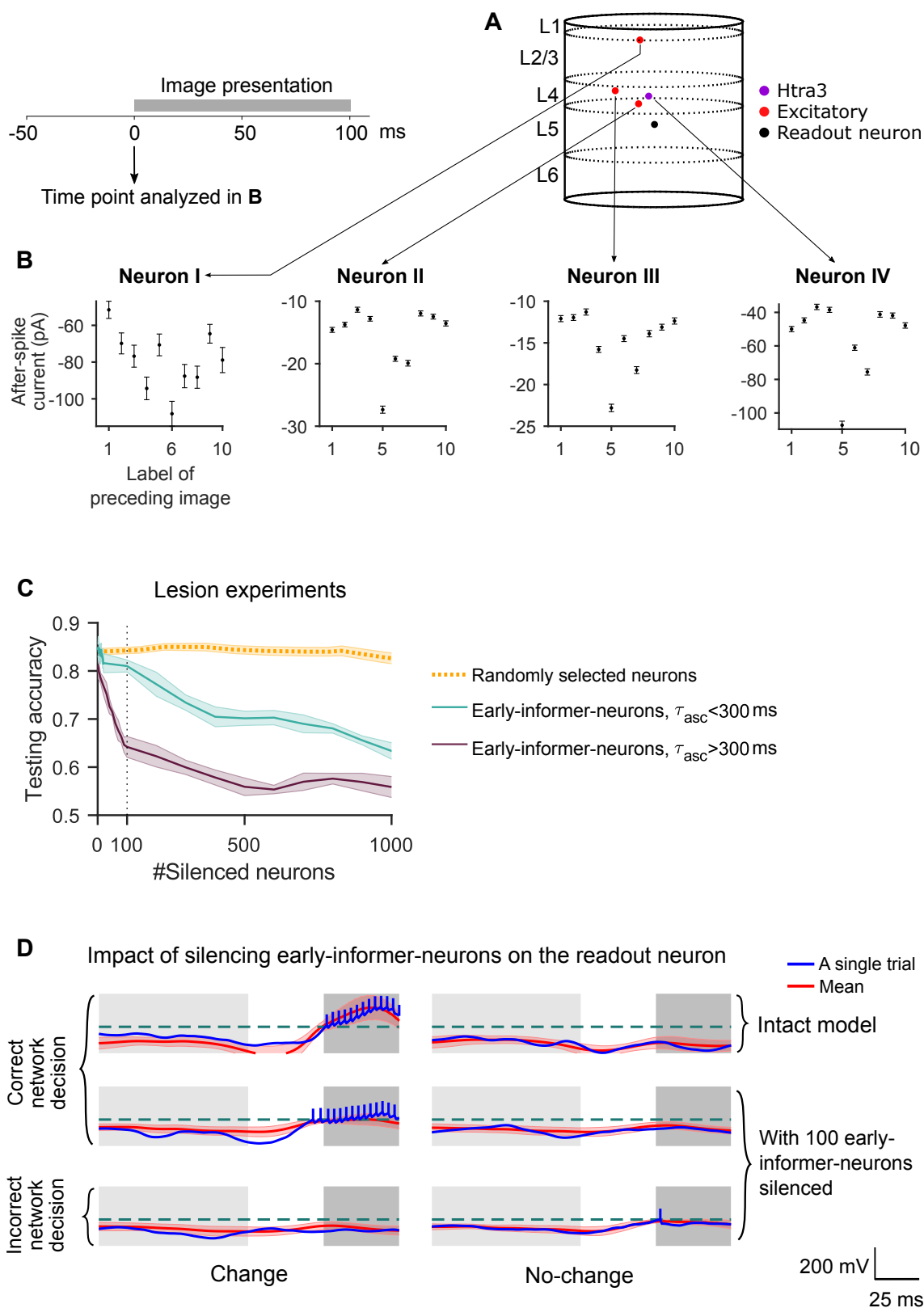
Figure 6 (previous page): Fine-scale analysis of the emergence of a network decision in the V1 model. (A) The mutual information (MI) between activities in 50 ms windows of single neurons and the change/no-change decision of the network (as emerging during the subsequent response window) is estimated for each neuron. For neurons that overlap in the projection from 3D to 2D, the maximum value is visualized to avoid dark points arising from accumulation of small contributions from several neurons. The critical time period from 50 to 100 ms after image onset is marked by dashed lines. (B) Spike responses of 7 early-informer-neurons, selected to cover all four basic neuron types: excitatory (E), Htra3, Sst, and Pvalb neurons (color-coded as in Fig. 1A). L2/3, L4, and L5 represent layer 2/3, 4, and 5, respectively. The value of the MI of the firing activity of each of these neurons during [50, 100] ms with the network decision is indicated in bits at the top of each panel. The periods of image presentation and subsequent response window are shaded in grey. Trials are separated according to the change/no-change condition. Condition-dependent differences in the firing responses of these neurons first appear at the start of the critical time window, 50 ms after image onset. (C) Time course of the values of the internal variables of these neurons that had the longest time constant (and modeled after-spike currents), for the same trials as in (B) (time constant shown on top). As GLIF₃ neurons have two after-spike currents that could have different time constants, we only show the after-spike currents with the largest time constant. Their values at image onset could potentially contain information about the identity of the preceding image, shown 200 ms before the current one. This information will be analyzed in the next figure.

389 model not only on the trial type, but also on the identity of the current image (Fig. S3). This additional
390 dependence of the temporal order of peak activity on the sensory input (for the same trial type) is a
391 prediction for future neurobiological experiments.

392 We found that the RANN was not able to reproduce the experimental data on the temporal organization
393 of cortical computations: Most of its neurons were active during rather long segments of a trial (see
394 Fig. 4A), and a fairly large fraction of neurons was simultaneously active at the same time. The temporal
395 order of their periods of high activity appeared to have a stereotypical character (see Fig. 4B, C) with
396 no apparent dependence on the trial type.

397 Apart from this different temporal organization of the network computation, the data-based model also
398 reproduced another characteristic feature of cortical computations: In spite of the large number of neurons
399 in the model (about 52,000), a rather small subset of neurons had a decisive impact on the time course
400 and outcome of the computation. Silencing of 100 neurons was sufficient to switch the trajectory of
401 network states during the computation (Fig. 5D), thereby drastically reducing the accuracy of the network
402 computation (Figs. 4D, 7C, D). This modeling result is consistent with experimental data which have
403 shown that activation of a small number of neurons in a cortical microcircuit is able to switch the
404 behavioral response of the mouse (Houweling and Brecht 2008; Doron et al. 2014; Marshel et al. 2019;
405 Dagleish et al. 2020; Doron et al. 2020). This high sensitivity of network decisions to the activity of
406 very small subsets of neurons is also of interest from a theoretical perspective: It suggests that the
407 V1 model operates, like the brain, in a critical regime. Importantly, this can be reproduced in the V1
408 model in spite of the substantial level of noise and trial-to-trial variability that we have placed, based
409 on experimental data from (Stringer et al. 2019), into the model, see (Chen et al. 2022) for details. In
410 contrast, we found that the RANN does not operate in a critical regime: Very large subsets of its neurons
411 have to be silenced in order to strongly reduce its task performance. Altogether, our results show that
412 the organization of computations is in RANNs substantially different from the brain, and that data-based
413 cortical microcircuit models provide a new family of models that can close this gap.

414 Reverse-engineering will be essential for understanding the organization of cortical computations. Doing
415 that in the living brain is still handicapped by limitations of current experimental techniques. But one
416 can spearhead such research by exploring and fine-tuning methods for reverse-engineering computations
417 in data-based models that employ a similar architecture and neuron types as the brain. We have shown
418 that low-dimensional projections provide substantial information about the organization of computations



(caption next page)

Figure 7 (previous page): Inner workings of 4 sample neurons that are pivotal for the network decision. (A) Spatial location and type of the 4 neurons labeled as Neuron I-IV in Fig. 6B and C whose firing rates have during [50, 100] ms after image presentation the largest MI with the network decision that is made 100 ms later during the response window. (B) Analysis of the information contained in the after-spike currents of the four neurons. The after-spike current with the largest time constant is analyzed at the onset of the current image (see time point marked in the scheme above), in dependence on the identity of the PRECEDING test image, marked on the horizontal axis. Its mean value is shown with error bars denoting the SEM across 1000 trials. One sees that this value depends strongly on the identity of the preceding image, especially for neurons II - IV. (C) Silencing early-informer-neurons with after-spike currents that have large time constants ($\tau_{asc} > 300$ ms) in the descending order of their MI reduces testing performance much faster than silencing those with shorter ones ($\tau_{asc} < 300$ ms). The shaded areas represent the SEM across 10 models with different training datasets. (D) The membrane potential of the readout neuron is significantly changed under both the change and the no-change condition when 100 early-informer-neurons are silenced. We used the same stimuli but different realizations of the noise model and initial states. The dashed green lines represent the firing threshold. Blue curves represent the membrane potentials of single trials; the extended vertical bars denote the spikes in these trials. Red curves represent the average membrane potentials across 100 trials. The shaded red area represents the standard deviation across 100 trials.

419 in the data-based model of Billeh et al., see Fig. 5. These results suggest that among the numerous
420 potentially relevant dynamical principles (Rabinovich et al. 2006) bifurcations of the network activity
421 turn out to be highly relevant for analyzing computational progress in these computations. While the
422 functional role of bifurcations of neural activity had previously focused on single neuron models, more
423 recent experimental data provide evidence that bifurcations are also essential for understanding the
424 functional role of populations of neurons in the living brain (Z. Wang et al. 2022). We have exhibited
425 in Fig. 5D a further-going prediction: These bifurcations are highly sensitive to the firing activity of
426 small sets of neurons: Silencing of just 100 neurons is able to flip a trajectory to a different bundle of
427 trajectories that produces a different network decision. This causal relationship between the activity of
428 individual neurons during a network computation and the resulting network decision needs to be tested
429 in future neurobiological experiments.

430 The data-based model allows us also to take a closer look at the inner workings of these pivotal neurons,
431 and to analyze how they can collect and transmit cues from two sequentially presented images which
432 indicate whether the second image is a different one. We found that the values of internal variables with
433 long time constants, which are abundantly present in the generalized leaky integrate-and-fire (GLIF₃)
434 neuron models of Billeh et al. that had been fitted to data of specific neurons from the Allen Brain
435 Atlas (Allen Institute 2018), assume values at the beginning of the processing of the current image that
436 contain salient information about the identity of the preceding image (Fig. 6C and 7B). The specific
437 causal impact of these neurons with long time constants on the network decision was verified through
438 further lesion experiments in our model (Fig. 7D).

439 Altogether, our work demonstrates the feasibility of a new methodology for understanding the organiza-
440 tion of computations in the neocortex: One can achieve a direct alignment of computational modeling
441 and neurobiological experiments by analyzing computations in detailed large-scale models, whose spatial
442 organization and neuron types match directly those found in the corresponding region of the neocortex,
443 and which solve the same computational task on the same ensemble of stimuli as the subjects in the neu-
444robiological experiments. This approach has now become feasible through recent advances in software
445 design, such as TensorFlow, and computing hardware, such as graphical processing units (GPUs), that
446 have been produced for the purpose of accelerating deep learning applications in AI.

447 An obvious next step in the direction of this work is an analysis of computations in detailed models
448 of cortical microcircuits in other cortical areas such as motor cortex, and of distributed computations

449 in interconnected cortical microcircuits from different cortical areas. Furthermore, at least simplified
 450 models for the dendritic arborization of selected classes of neurons need to be added in order to make
 451 the impact of top-down inputs more realistic. But this needs to be done in a way that still supports
 452 fast simulations of large-scale models, so that training of these models for specific computational tasks
 453 remains computationally feasible. In addition, the computational role of projections to and synaptic
 454 inputs from subcortical areas such as basal ganglia and the thalamus (see (Cruz et al. 2023) for a recent
 455 review) needs to be modeled and analyzed through integrated large-scale models. Altogether, this work is
 456 likely to complement, challenge, and enhance experimental work that aims at clarifying the organization
 457 of brain computations. In addition, it will provide paradigms for a new generation of artificial neural
 458 network models that can capture both the astounding functional capabilities and the energy efficiency of
 459 sparsely active neural networks in the brain.

460 4 Methods

461 4.1 Neuron models

462 As in (Chen et al. 2022) we are focusing on the “core” part of the point-neuron version of the realistic
 463 V1 model introduced by (Billeh et al. 2020). To make it gradient-friendly, we replaced the hard reset
 464 of membrane potential after a spike emerges with the reduction of membrane potential $z_j(t)(v_{\text{th}} - E_L)$,
 465 where $z_j(t) = 1$ when neuron j fires at time t and $z_j(t) = 0$ otherwise. v_{th} is the firing threshold of
 466 membrane potential. E_L the resting membrane potential. This causes no significant change in the neural
 467 response (Chen et al. 2022). We simulated each trial for 600 ms. The dynamics of the modified GLIF₃
 468 model was defined as

$$\begin{aligned}
 v_j(t + \delta t) &= \alpha v_j(t) + \frac{1 - \alpha\tau}{C} \left(I_j^e(t + 1) + \sum_m I_j^m(t + 1) + gE_L + I_j^{\text{syn}}(t) \right) - z_j(t)(v_{\text{th}} - E_L) \\
 z_j(t) &= H(v_j(t) - v_{\text{th}}) \\
 I_j^e(t) &= \sum_i W_{ji}^{\text{in}} x_i(t) + qK_j^{\text{quick}}(t) + sK_j^{\text{slow}},
 \end{aligned}
 \tag{1}$$

469 where C represents the neuron capacitance, I^e the external current, I^{syn} the synaptic current, g the
 470 membrane conductance, and v_{th} the spiking threshold. W_{ji}^{in} is the synaptic weight from LGN neuron i to
 471 V1 neuron j . The scales of the quick noise $K_j^{\text{quick}}(t)$ and the slow noise K_j^{slow} to neuron j are $q = 2$ and
 472 $s = 2$, respectively, unless otherwise stated. K_j was randomly drawn from the empirical noise distribution
 473 which will be elaborated on later. The decay factor α is given by $e^{-\delta t/\tau}$, where τ is the membrane time
 474 constant. δt denotes the discrete-time step size, which is set to 1 ms in our simulations. H denotes the
 475 Heaviside step function. To introduce a simple model of neuronal refractoriness, we further assumed that
 476 $z_j(t)$ is fixed to 0 after each spike of neuron j for a short refractory period depending on the neuron type.
 477 The after-spike current $I^m(t)$ was modeled as

$$I^m(t + \delta t) = f^m I^m(t) + z(t)\delta I^m; \quad m = 1, \dots, N_{\text{asc}},
 \tag{2}$$

478 where the multiplicative constant $f^m = \exp(-k^m\delta t)$ and an additive constant, δI^m . In our study, $m = 1$
 479 or 2. Neuron parameters have been fitted to experimental data from 111 selected neurons according to the
 480 cell database of the Allen Brain Atlas (Allen Institute 2018), see (Teeter et al. 2018; Billeh et al. 2020),
 481 including neuron capacity C , conductance g , resting potential E_L , the length of the refractory period, as
 482 well as amplitudes δI^m and decay time constants k^m of two types of after-spike currents, $m = 1, 2$.

483 4.2 Synaptic inputs

484 The V1 model utilizes experimental data to specify the connection probability between neurons. The
485 base connection probability for any pair of neurons from the 17 cell classes is provided in (Billeh et al.
486 2020) in a table (shown in Fig. 1B), where white cells denote unknown values. The values in this table
487 are derived from measured frequencies of synaptic connections for neurons at maximal 75 μm horizontal
488 inter-somatic distance. The base connection probability was then scaled by an exponentially decaying
489 factor based on the horizontal distance between the somata of the two neurons (Fig. 1C), also derived
490 from experimental data. The synaptic delay was spread in [1, 4] ms, as extracted from Fig. 4E of (Billeh
491 et al. 2020) and rounded to the nearest integer as the integration step is 1 ms.

The postsynaptic current of neuron j was defined by the following dynamics (Billeh et al. 2020):

$$I_j^{\text{syn}}(t + \delta t) = e^{-\frac{\delta t}{\tau_{\text{syn}}}} I_j^{\text{syn}}(t) + \delta t e^{-\frac{\delta t}{\tau_{\text{syn}}}} C_j^{\text{rise}}(t) \quad (3)$$

$$C_j^{\text{rise}}(t + \delta t) = e^{-\frac{\delta t}{\tau_{\text{syn}}}} C_j^{\text{rise}}(t) + \sum_i W_{ji}^{\text{rec}} z_i(t) \frac{e}{\tau_{\text{syn}}}, \quad (4)$$

492 where τ_{syn} is the synaptic time constant, W_{ji}^{rec} is the recurrent input connection weight from neuron i to
493 j , and z_i is the spike of presynaptic neuron i . The τ_{syn} constants depend on neuron types of pre- and
494 postsynaptic neurons (Billeh et al. 2020).

495 4.3 Initial conditions

496 The default initial conditions for spikes and membrane potentials were set to zero, unless otherwise
497 specified. The initial conditions for \mathbf{W}^{in} and \mathbf{W}^{rec} were taken from (Billeh et al. 2020), unless otherwise
498 stated.

499 4.4 Data-driven noise model

500 We used a noise model that was introduced in our previous study (Chen et al. 2022). The model was
501 based on an empirical noise distribution that was obtained from experimental data of mice responses to
502 2,800 nature images (Stringer et al. 2019). The noise currents $K_j^{\text{quick}}(t)$ and K_j^{slow} in Eq. 1 were drawn
503 independently for all neurons from this distribution. The quick noise $K_j^{\text{quick}}(t)$ was drawn every 1 ms
504 while the slow noise K_j^{slow} was drawn once every 600 ms. The empirical noise distribution was derived
505 from the variability (additive noise) collected from the experimental data. A detailed mathematical
506 analysis of this method is available in the methods and supplementary materials of (Stringer et al. 2019).

507 4.5 Readout neurons

508 We employed a readout population in the V1 model, whose firing activity during the response window
509 encoded the network decisions for the visual-change-detection task. Each population consisted of a certain
510 number (30 or 1) of randomly selected excitatory neurons in layer 5, located within a sphere of a radius
511 of 55 μm (Fig. 2E).

512 4.6 Visual-change-detection task

513 **LGN model.** The visual stimuli were processed by a qualitative retina and LGN model, as depicted in
 514 Fig. 2C and following (Billeh et al. 2020). Their full LGN model consists of 17,400 spatiotemporal filters
 515 that simulate the responses of LGN neurons in mice to visual stimuli (Durand et al. 2016). Each filter
 516 generates a positive output, which represents the firing rates of a corresponding LGN neuron. We used
 517 only a subset of 2,589 of these LGN filters that provide inputs from a smaller part of the visual field
 518 to the core part of the V1 model, on which we are focusing in this study. The input images were first
 519 converted to grayscale and scaled to fit in the interval $[-Int, Int]$, where $Int > 0$. The output of the
 520 LGN model was then used as an external current input in the V1 model as follows:

$$I_{sti} = \mathbf{W}^{in} \cdot \text{LGN}(G_{Int}), \quad (5)$$

521 where G_{Int} represents images scaled into $[-Int, Int]$ for $Int = 2$.

522 **Visual-change-detection task with natural images.** We designed the visual-change-detection task
 523 to be as close as possible to corresponding biological experiments while keeping them as simple as possible.
 524 In the mouse experiments of (Garrett et al. 2020; Joshua H. Siegle et al. 2021), mice were trained to
 525 perform a visual change detection task using static natural images presented in a sequence of 250 ms
 526 with short phases (500 ms) of gray screens in between. The mice had to report whether the most recently
 527 presented image was the same as the previously presented one. To replicate this task while taking into
 528 account GPU memory limitations, we presented natural images for 100 ms each with delays between them
 529 lasting 200 ms (Fig. 2A, B). The first image was presented after 50 ms, and all images were selected from
 530 a set of 40 randomly chosen images from the ImageNet dataset (Deng et al. 2009). The model had to
 531 report within a 50 ms time window starting 150 ms after image onset (response window) if the image had
 532 changed.

533 In the response window, we defined the mean firing rate of readout population as

$$r_{\text{readout}} = \frac{1}{T_{\text{resp}} \cdot N_{\text{readout}}} \sum_{t=1}^{T_{\text{resp}}} \sum_{j=1}^{N_{\text{readout}}} z_j(t), \quad (6)$$

534 where the sum over j is over the $N_{\text{readout}} = 30$ readout neurons and the sum over t is over the time
 535 length of response window $T_{\text{resp}} = 50$ ms. If $r > r_0 = 0.01$, the model reported a network decision that
 536 the image had changed. Otherwise, it reported no-change.

537 4.7 Loss function

538 The loss function was defined as

$$L = L_{\text{cross-entropy}} + \lambda_f L_{\text{rate reg.}} + \lambda_v L_{\text{v reg.}}, \quad (7)$$

539 where $L_{\text{cross-entropy}}$ represents the cross-entropy loss, λ_f and λ_v represent the weights of firing-rate regu-
 540 larization $L_{\text{rate reg.}}$ and voltage regularization $L_{\text{v reg.}}$, respectively. As an example, the cross-entropy loss
 541 of visual change detection tasks was given by

$$L_{\text{cross-entropy}} = - \sum_m \left[T^{(m)} \log \sigma \left(\theta \left(r_{\text{readout}}^{(m)} - r_0 \right) \right) + \left(1 - T^{(m)} \right) \log \sigma \left(\theta \left(r_0 - r_{\text{readout}}^{(m)} \right) \right) \right], \quad (8)$$

542 where the sum over m is organized into chunks of 50 ms and $r_{\text{readout}}^{(m)}$ denotes the mean readout population
 543 firing rate defined in Eq. 6. Similarly, $T^{(m)}$ denotes the target output in time window m , being 1 if a

544 change in image identity should be reported and otherwise 0. The baseline firing rate r_0 was 0.01. σ
545 represents the sigmoid function. θ is a trainable scale ($\theta > 0$) of firing rate.

546 We used regularization terms in the loss function to penalize very high firing rates as well as values of
547 membrane voltages that were not biologically realistic. The default values of their weights were $\lambda_f = 0.1$
548 and $\lambda_v = 10^{-5}$. The rate regularization is defined via the Huber loss (Huber 1992) between the target
549 firing rates, y , calculated from the model in (Billeh et al. 2020), and the firing rates, r , sampled the same
550 number of neurons from the network model:

$$L_{\text{rate reg.}} = \sum_j^N |\tau_j - \mathbb{I}\{\delta_j < 0\}| \frac{\mathcal{L}_\kappa(\delta_j)}{\kappa}, \quad \text{with} \quad (9)$$

$$\mathcal{L}_\kappa(\delta_j) = \begin{cases} \frac{1}{2}\delta_j^2, & \text{if } |\delta_j| \leq \kappa \\ \kappa(|\delta_j| - \frac{1}{2}\kappa), & \text{otherwise} \end{cases}$$

551 where j represents neuron j , N the number of neurons, $\tau_j = j/N$, $\delta = 0.002$, $\delta_j = r_j - r_j^{\text{target}}$. $\mathbb{I}(x) = 1$
552 when x is true; $\mathbb{I}(x) = 0$ when x is false.

553 The voltage regularization is defined through the term

$$L_{\text{v reg.}} = \frac{1}{N} \sum_{j=0}^{j=N} \left(\left[\frac{v_j - E_L}{v_{\text{th}} - E_L} - 1 \right]^+ \right)^2 + \left(\left[-\frac{v_j - E_L}{v_{\text{th}} - E_L} + 1 \right]^+ \right)^2, \quad (10)$$

554 where N represents the total number of neurons, v_j , the membrane potential of neuron j , $[\dots]^+$, the
555 rectifier function. v_{th} is the firing threshold of membrane potential. E_L the resting membrane potential.

556 4.8 Training and testing

557 We applied back-propagation through time (BPTT) (Chen et al. 2022) to minimize the loss function. The
558 non-existing derivative $\frac{\partial z_j}{\partial v_j}$ was replaced in simulations by a simple nonlinear function of the membrane
559 potential that is called the pseudo-derivative. Outside of the refractory period, we chose a pseudo-
560 derivative of the form

$$\psi^t = \frac{\gamma_{\text{pd}}}{v_{\text{th}} - E_L} \exp\left(-\frac{(v_{\text{sc}}^t)^2}{\sigma_p^2}\right), \quad (11)$$

$$v_{\text{sc}}^t = \frac{v^t - v_{\text{th}}}{v_{\text{th}} - E_L},$$

561 where the dampening factor $\gamma_{\text{pd}} = 0.5$, the Gaussian kernel width $\sigma_p = 0.28$. During the refractory
562 period, the pseudo derivative was set to 0.

563 To demonstrate how sensitive the performance is to the scale of the surrogate derivative, I trained the
564 model with $\gamma_{\text{pd}} = 0.25$ and 0.75 and kept all other hyperparameters the same. When $\gamma_{\text{pd}} = 0.25$, the
565 testing accuracy is 0.7; when $\gamma_{\text{pd}} = 0.75$, the testing accuracy is 0.75. Compared with the case of
566 $\gamma_{\text{pd}} = 0.5$ where the testing accuracy is 0.83, other values are worse. This demonstrates that the choice
567 of the derivative's scale can substantially affect gradient-based learning performance in spiking neural
568 networks (Zenke and Vogels 2021).

569 We drew a batch of visual stimuli (64) and calculated the gradient after every trial for each synaptic
570 weight whether an increase or decrease of it (but without changing its sign) would reduce the loss function.

571 Weights were then updated by the average gradient across the batch. This method had originally only
572 been applied to neuron networks with differentiable neuron models and was normally referred to as
573 stochastic gradient descent.

574 During the training, we added the sign constraint on the weights of the neural network to keep Dale’s
575 law. Specifically, if an excitatory weight was updated to a negative value, it would be set to 0; vice versa.
576 In every training run, we used a different random seed in order to draw fresh noise samples from the
577 empirical distribution, and to randomly generate/select training samples.

578 4.9 Other simulation details

579 The BPTT training algorithm was implemented in TensorFlow, which is optimized to run efficiently on
580 GPUs, allowing us to take advantage of their parallel computing capabilities. We distributed the visual-
581 change-detection task trials over batches, with each batch containing 64 trials, and performed independent
582 simulations in parallel. Each trial lasted for 600 ms of biological time, and computing gradients for each
583 batch took around 5 s on an NVIDIA A100 GPU. Once all batches had finished (one step), gradients were
584 calculated and averaged to update the weights by BPTT. We define an epoch as 500 iterations/steps.
585 This computation had to be iterated for 22 epochs to make sure the performance was saturated. This
586 took 12 h of wall clock time on 32 GPUs.

587 4.10 Recurrent artificial neural network models (RANNs)

588 Model

589 The dynamics of a RANN can be defined as

$$\tau \frac{d\mathbf{x}}{dt} = -\mathbf{x}(t) + \mathbf{W}_r \mathbf{r}(t) + \mathbf{W}_{ff} \mathbf{I}(t) + \mathbf{b}_r, \quad (12)$$

590 where the \mathbf{x} is the activation of the network units and the corresponding firing rate is defined as

$$\mathbf{r} = \tanh(\mathbf{x}) \quad (13)$$

591 τ is the single-unit timescale, I used 50 ms as in (Sussillo et al. 2015; Pollock and Jazayeri 2020). \mathbf{W}_r is
592 the recurrent synaptic weights; \mathbf{W}_{ff} is the feedforward synaptic weights; \mathbf{b}_r is bias. The initialization of
593 \mathbf{W}_r , \mathbf{W}_{ff} and \mathbf{b}_r are Gaussian noise $\mathcal{N}(0, \sigma^2)$. $\sigma = 1/\sqrt{\text{average \# fan-in}}$ for \mathbf{W}_r and \mathbf{W}_{ff} ; $\sigma = 0.1$ for
594 \mathbf{b}_r . $\mathbf{I}(t)$ is the input to the network at time t . It is the output of LGN model (Billeh et al. 2020).

595 The linear readout $\mathbf{y}(t)$ from activities of all neurons $\mathbf{r}(t)$,

$$\mathbf{y}(t) = \mathbf{W}_y^\top \mathbf{r}(t) + \mathbf{b}_y, \quad (14)$$

596 where \mathbf{b}_y is bias. \mathbf{W}_y is the readout weight. \mathbf{W}_y is a $N \times 2$ matrix; N is the number of recurrent neurons;
597 2 is the number of possible decisions.

598 The number of neurons and synapses in our model are the same as those in Billeh’s model (Billeh et al.
599 2020). We randomly shuffled the connectivity or kept it the same as in Billeh’s model.

600 Loss function

601 The loss function was defined as

$$L = L_{\text{cross-entropy}} + \lambda_{\text{reg}} L_{\text{reg}}, \quad (15)$$

602 where $L_{\text{cross-entropy}}$ represents the cross-entropy loss which was defined in Eq. 8, $\lambda_{\text{reg.}}$ represents the
603 weight of activation regularization $L_{\text{reg.}}$. The activation regularization was defined as

$$L_{\text{reg.}} = \left(\left[\frac{1}{N} \sum_i |r_i| - \theta \right]^+ \right)^2, \quad (16)$$

604 where r_i is the “firing rate” of neuron i . N is the number of neurons and θ is a threshold ($\theta = 0.01$),
605 unless otherwise stated. $[\cdot]^+$ is the rectifier linear unit function. The value of θ was determined as $\frac{1}{4}$
606 of the mean value of $|r|$ when training the RANN model without regularization.

607 **Training details**

608 We used the same training methods as in the V1 model but the learning rate is 10^{-4} . We also used
609 $dt = 5$ ms to alleviate the vanishing gradient problem, as the one used in (Sussillo et al. 2015).

610 **Details for calculating normalized averaged activity in Fig. 4A-C**

611 As the “firing rate” of neurons in RANN, \mathbf{r} is in $[-1, 1]$, we took the absolute value of \mathbf{r} to compare with
612 the neural activity in the V1 model.

613 **4.11 Demixed principal component analysis**

614 Demixed principal component analysis (Demixed PCA) is a statistical method that decomposes high-
615 dimensional neural data into a set of orthogonal latent variables, each of which captures a unique aspect
616 of the neural response Kobak et al. 2016. Briefly, let $\mathbf{X}n \times T$ be the neural data matrix, where n is
617 the number of neurons and T is the number of time points. Let $\mathbf{S}n \times C \times T$ be the tensor of stimulus
618 conditions, where C is the number of experimental conditions. The goal of Demixed PCA is to find a
619 low-dimensional latent space $\mathbf{Y}_{d \times T}$, where d is the number of latent variables, that captures the majority
620 of the variance in the neural data \mathbf{X} , while also separating out the variance that is specific to each
621 experimental condition in \mathbf{S} .

622 In Fig. 3D, the principal component used for the projection arises by analysis of the first eigenvector
623 of the covariance matrix that reflects variation through joint dependencies of the network decision and
624 relative timing, see marginalization procedure of (Kobak et al. 2016). Hence, this matrix does not reflect
625 variation that is caused only through the course of time within a trial or through the network decision
626 alone. Moreover, to emphasize the formation of the network decision, we include in the computation of
627 the aforementioned covariance matrix only data within $[-50, 50]$ ms of the image presentation.

628 **4.12 Mutual information**

629 To estimate the mutual information between single neuron activity and the network decision, we binned
630 the spike counts of each neuron into 10 uniformly distributed bins between the minimum and maximum
631 spike count observed for that neuron within 50 ms windows. We then established an empirical joint
632 distribution for the binned spike count and the network decision and computed the mutual information
633 using the below formula.

$$MI(X; Y) = \sum_{y \in Y} \sum_{x \in X} P(x, y) \log_2 \frac{P(x, y)}{P(x)P(y)}, \quad (17)$$

634 where X is the set of firing activities of each neuron within a 50 ms window, and Y is the set of network
635 decisions (either change or no-change). $P(x, y)$ is the joint probability distribution of spike count and

636 network decision, while $P(x)$ and $P(y)$ are the marginal probability distributions of spike count and
637 network decision, respectively.

638 To estimate $P(x, y)$, we calculated the spike count of each neuron within a 50 ms window and established
639 an empirical joint distribution by counting the number of occurrences of each possible combination of
640 spike counts and network decisions across 100,000 trials. We then normalized the joint distribution to
641 obtain a probability distribution. We estimated $P(x)$ and $P(y)$ in a similar way by counting the number
642 of occurrences of each possible value of spike count and network decision, respectively, across all trials.

643 4.13 UMAP

644 We applied an exponential filter with a time constant of 20 ms to the spike output of each neuron
645 for 8 new images that had not been used during training. We then discarded all but the 1,500 most
646 important principal components of these network states (explain 38% of variance; compromise to the
647 memory consumption), and embedded these into 2D space by UMAP. The number of neighbors is 200.
648 These projected network states were recorded for every ms, represented by a dot in Fig. 5.

649 UMAP (Uniform Manifold Approximation and Projection) is a nonlinear dimensionality reduction algo-
650 rithm that aims to preserve the local structure of high-dimensional data in low-dimensional space (McInnes
651 et al. 2018). In this study, UMAP was used to embed the network activity during task performance into
652 2D space.

653 First, we applied an exponential filter with a time constant of 20 ms to the spike output of each neuron
654 for 8 new images that had not been used during training. We then computed the principal components
655 of the resulting network states, discarding all but the 1,500 most important components that explained
656 38% of the variance.

657 Next, we used UMAP to embed the high-dimensional network states into 2D space, while preserving the
658 local structure of the data. Specifically, we used the UMAP implementation from the Python library
659 `umap`, with the following parameters: `n_neighbors = 200`, `min_dist = 0.1`, and `metric = euclidean`.

660 The resulting 2D embeddings represent the low-dimensional trajectories of the network states during
661 image processing, and were recorded for every ms. Each point in the 2D space represents a network state
662 at a given time point, and is displayed as a dot in Fig. 5. The trajectory of the network states can be
663 visualized by connecting these dots in chronological order.

664 4.14 Additional figure description

665 Normalized activity in Fig. 3, 4, and S6

666 Spiking activity at a specific relative time step, regarding image presentation, was averaged over 200
667 trials. These average activities per time step were then normalized with the maximum values of their
668 average activation.

669 Fig. 6B and C

670 7 neurons shown in Fig. 6B and C were selected from the 20 early-informer-neurons with the largest MI
671 that represented each of the 4 neuron classes, taking within each neuron class (excitatory, PV, Htra3, or
672 Sst neurons) the ones with the largest MI.

673 Acknowledgements

674 We would like to thank Yuqing Zhu for helpful discussions. We also thank Sandra Diaz for advice and
675 help in using supercomputers. This research was partially supported by the Human Brain Project (Grant
676 Agreement number 785907) of the European Union and a grant from Intel. Computations were carried
677 out on the Human Brain Project PCP Pilot Systems at the Jülich Supercomputing Centre, which received
678 co-funding from the European Union (Grant Agreement number 604102).

679 References

- 680 Abeles, Moshe, Gaby Hayon, and Daniel Lehmann (2004). “Modeling compositionality by dynamic bind-
681 ing of synfire chains”. In: *Journal of computational neuroscience* 17, pp. 179–201.
- 682 Allen, William E, Isaac V Kauvar, Michael Z Chen, Ethan B Richman, Samuel J Yang, Ken Chan, Viviana
683 Gradinaru, Benjamin E Deverman, Liqun Luo, and Karl Deisseroth (2017). “Global representations
684 of goal-directed behavior in distinct cell types of mouse neocortex”. In: *Neuron* 94.4, pp. 891–907.
- 685 Allen Institute (2018). “© 2018 Allen Institute for Brain Science. Allen Cell Types Database, cell feature
686 search. Available from: celltypes.brain-map.org/data”. In:
- 687 Bellec, Guillaume, Darjan Salaj, Anand Subramoney, Robert Legenstein, and Wolfgang Maass (2018).
688 “Long short-term memory and Learning-to-learn in networks of spiking neurons”. In: *Advances in
689 Neural Information Processing Systems*. Ed. by S. Bengio, H. Wallach, H. Larochelle, K. Grauman,
690 N. Cesa-Bianchi, and R. Garnett. Vol. 31. Curran Associates, Inc., pp. 787–797. URL: [https://
691 proceedings.neurips.cc/paper/2018/file/c203d8a151612acf12457e4d67635a95-Paper.pdf](https://proceedings.neurips.cc/paper/2018/file/c203d8a151612acf12457e4d67635a95-Paper.pdf).
- 692 Billeh, Yazan N, Binghuang Cai, Sergey L Gratiy, Kael Dai, Ramakrishnan Iyer, Nathan W Gouwens,
693 Reza Abbasi-Asl, Xiaoxuan Jia, Joshua H Siegle, Shawn R Olsen, et al. (2020). “Systematic integration
694 of structural and functional data into multi-scale models of mouse primary visual cortex”. In: *Neuron*.
- 695 Broome, Bede M, Vivek Jayaraman, and Gilles Laurent (2006). “Encoding and decoding of overlapping
696 odor sequences”. In: *Neuron* 51.4, pp. 467–482.
- 697 Chen, Guozhang, Franz Scherr, and Wolfgang Maass (2022). “A data-based large-scale model for primary
698 visual cortex enables brain-like robust and versatile visual processing”. In: *Science Advances* 8.44,
699 eabq7592.
- 700 Cruz, K Guadalupe, Yi Ning Leow, Nhat Minh Le, Elie Adam, Rafiq Huda, and Mriganka Sur (2023).
701 “Cortical-subcortical interactions in goal-directed behavior”. In: *Physiological reviews* 103.1, pp. 347–
702 389.
- 703 Dagleish, Henry WP, Lloyd E Russell, Adam M Packer, Arnd Roth, Oliver M Gauld, Francesca Green-
704 street, Emmett J Thompson, and Michael Häusser (2020). “How many neurons are sufficient for
705 perception of cortical activity?” In: *Elife* 9, e58889.
- 706 Deng, Jia, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei (2009). “Imagenet: A large-scale
707 hierarchical image database”. In: *2009 IEEE conference on computer vision and pattern recognition*.
708 Ieee, pp. 248–255.
- 709 Doron, Guy, Jiyun N Shin, Naoya Takahashi, Moritz Drüke, Christina Bocklisch, Salina Skenderi, Lisa de
710 Mont, Maria Toumazou, Julia Ledderose, Michael Brecht, et al. (2020). “Perirhinal input to neocortical
711 layer 1 controls learning”. In: *Science* 370.6523, eaaz3136.
- 712 Doron, Guy, Moritz Von Heimendahl, Peter Schlattmann, Arthur R Houweling, and Michael Brecht
713 (2014). “Spiking irregularity and frequency modulate the behavioral report of single-neuron stimula-
714 tion”. In: *Neuron* 81.3, pp. 653–663.
- 715 Douglas, Rodney J and Kevan AC Martin (2004). “Neuronal circuits of the neocortex”. In: *Annu. Rev.*
716 *Neurosci.* 27, pp. 419–451.

- 717 Driscoll, Laura N, Noah L Pettit, Matthias Minderer, Selmaan N Chettih, and Christopher D Harvey
718 (2017). “Dynamic reorganization of neuronal activity patterns in parietal cortex”. In: *Cell* 170.5,
719 pp. 986–999.
- 720 Durand, Séverine, Ramakrishnan Iyer, Kenji Mizuseki, Saskia de Vries, Stefan Mihalas, and R Clay Reid
721 (2016). “A comparison of visual response properties in the lateral geniculate nucleus and primary
722 visual cortex of awake and anesthetized mice”. In: *Journal of Neuroscience* 36.48, pp. 12144–12156.
- 723 Garrett, Marina, Sahar Manavi, Kate Roll, Douglas R Ollerenshaw, Peter A Groblewski, Nicholas D
724 Ponvert, Justin T Kiggins, Linzy Casal, Kyla Mace, Ali Williford, Arielle Leon, Xiaoxuan Jia, Peter
725 Ledochowitsch, Michael A Buice, Wayne Wakeman, Stefan Mihalas, and Shawn R Olsen (Feb. 2020).
726 “Experience shapes activity dynamics and stimulus coding of VIP inhibitory cells”. In: *eLife* 9. Ed. by
727 Brice Bathellier, Joshua I Gold, Brice Bathellier, and Georg B Keller, e50340. ISSN: 2050-084X. DOI:
728 [10.7554/eLife.50340](https://doi.org/10.7554/eLife.50340). URL: <https://doi.org/10.7554/eLife.50340>.
- 729 Harris, Kenneth D. and Gordon M. G. Shepherd (Feb. 2015). “The neocortical circuit: themes and
730 variations”. In: *Nature Neuroscience* 18.2, pp. 170–181. ISSN: 1546-1726. DOI: [10.1038/nn.3917](https://doi.org/10.1038/nn.3917). URL:
731 <https://doi.org/10.1038/nn.3917>.
- 732 Houweling, Arthur R and Michael Brecht (2008). “Behavioural report of single neuron stimulation in
733 somatosensory cortex”. In: *Nature* 451.7174, pp. 65–68.
- 734 Huber, Peter J (1992). “Robust estimation of a location parameter”. In: *Breakthroughs in statistics*.
735 Springer, pp. 492–518.
- 736 Kappel, David, Bernhard Nessler, and Wolfgang Maass (2014). “STDP installs in winner-take-all circuits
737 an online approximation to hidden Markov model learning”. In: *PLoS computational biology* 10.3,
738 e1003511.
- 739 Kato, Saul, Harris S Kaplan, Tina Schrödel, Susanne Skora, Theodore H Lindsay, Eviatar Yemini, Shawn
740 Lockery, and Manuel Zimmer (2015). “Global brain dynamics embed the motor command sequence
741 of *Caenorhabditis elegans*”. In: *Cell* 163.3, pp. 656–669.
- 742 Koay, Sue Ann, Adam S Charles, Stephan Y Thiberge, Carlos D Brody, and David W Tank (2022).
743 “Sequential and efficient neural-population coding of complex task information”. In: *Neuron* 110.2,
744 pp. 328–349.
- 745 Kobak, Dmitry, Wieland Brendel, Christos Constantinidis, Claudia E Feierstein, Adam Kepecs, Zachary
746 F Mainen, Xue-Lian Qi, Ranulfo Romo, Naoshige Uchida, and Christian K Machens (Apr. 2016).
747 “Demixed principal component analysis of neural population data”. In: *eLife* 5. Ed. by Mark CW van
748 Rossum, e10989. ISSN: 2050-084X. DOI: [10.7554/eLife.10989](https://doi.org/10.7554/eLife.10989). URL: <https://doi.org/10.7554/eLife.10989>.
- 749
- 750 Maass, Wolfgang and Henry Markram (2004). “On the computational power of circuits of spiking neu-
751 rons”. In: *Journal of computer and system sciences* 69.4, pp. 593–616.
- 752 Maass, Wolfgang, Thomas Natschläger, and Henry Markram (2002). “Real-time computing without stable
753 states: A new framework for neural computation based on perturbations”. In: *Neural computation*
754 14.11, pp. 2531–2560.
- 755 Markram, Henry, Eilif Muller, Srikanth Ramaswamy, Michael W Reimann, Marwan Abdellah, Carlos
756 Aguado Sanchez, Anastasia Ailamaki, Lidia Alonso-Nanclares, Nicolas Antille, Selim Arsever, et al.
757 (2015). “Reconstruction and simulation of neocortical microcircuitry”. In: *Cell* 163.2, pp. 456–492.
- 758 Marshel, James H, Yoon Seok Kim, Timothy A Machado, Sean Quirin, Brandon Benson, Jonathan
759 Kadmon, Cephra Raja, Adelaida Chibukhchyan, Charu Ramakrishnan, Masatoshi Inoue, et al. (2019).
760 “Cortical layer-specific critical dynamics triggering perception”. In: *Science* 365.6453.
- 761 Martin Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S. Cor-
762 rado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew Harp,
763 Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser, Manjunath Kudlur,
764 Josh Levenberg, Dandelion Mané, Rajat Monga, Sherry Moore, Derek Murray, Chris Olah, Mike
765 Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent Van-
766 houcke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin

- 767 Wicke, Yuan Yu, and Xiaoqiang Zheng (2015). *TensorFlow: Large-Scale Machine Learning on Het-*
768 *erogeneous Systems*. Software available from tensorflow.org. URL: <https://www.tensorflow.org/>.
- 769 McInnes, L., J. Healy, and J. Melville (Feb. 2018). “UMAP: Uniform Manifold Approximation and
770 Projection for Dimension Reduction”. In: *ArXiv e-prints*. arXiv: [1802.03426](https://arxiv.org/abs/1802.03426) [stat.ML].
- 771 Mountcastle, Vernon B (1998). *Perceptual neuroscience: The cerebral cortex*. Harvard University Press.
- 772 Pollock, Eli and Mehrdad Jazayeri (2020). “Engineering recurrent neural networks from task-relevant
773 manifolds and dynamics”. In: *PLoS computational biology* 16.8, e1008128.
- 774 Rabinovich, Mikhail I, Pablo Varona, Allen I Selverston, and Henry DI Abarbanel (2006). “Dynamical
775 principles in neuroscience”. In: *Reviews of modern physics* 78.4, p. 1213.
- 776 Russakovsky, Olga, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang,
777 Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei (2015). “Im-
778 ageNet Large Scale Visual Recognition Challenge”. In: *International Journal of Computer Vision*
779 *(IJCV)* 115.3, pp. 211–252. DOI: [10.1007/s11263-015-0816-y](https://doi.org/10.1007/s11263-015-0816-y).
- 780 Siegle, Joshua H., Xiaoxuan Jia, Séverine Durand, Sam Gale, Corbett Bennett, Nile Graddis, Gregory
781 Heller, Tamina K. Ramirez, Hannah Choi, Jennifer A. Luviano, Peter A. Groblewski, Ruweida Ahmed,
782 Anton Arkhipov, Amy Bernard, Yazan N. Billeh, Dillan Brown, Michael A. Buice, Nicolas Cain,
783 Shiella Caldejon, Linzy Casal, Andrew Cho, Maggie Chvilicek, Timothy C. Cox, Kael Dai, Daniel
784 J. Denman, Saskia E. J. de Vries, Roald Dietzman, Luke Esposito, Colin Farrell, David Feng, John
785 Galbraith, Marina Garrett, Emily C. Gelfand, Nicole Hancock, Julie A. Harris, Robert Howard, Brian
786 Hu, Ross Hytnen, Ramakrishnan Iyer, Erika Jessett, Katelyn Johnson, India Kato, Justin Kiggins,
787 Sophie Lambert, Jerome Lecoq, Peter Ledochowitsch, Jung Hoon Lee, Arielle Leon, Yang Li, Elizabeth
788 Liang, Fuhui Long, Kyla Mace, Jose Melchior, Daniel Millman, Tyler Mollenkopf, Chelsea Nayan,
789 Lydia Ng, Kiet Ngo, Thuyahn Nguyen, Philip R. Nicovich, Kat North, Gabriel Koch Ocker, Doug
790 Ollerenshaw, Michael Oliver, Marius Pachitariu, Jed Perkins, Melissa Reding, David Reid, Miranda
791 Robertson, Kara Ronellenfitch, Sam Seid, Cliff Slaughterbeck, Michelle Stoecklin, David Sullivan, Ben
792 Sutton, Jackie Swapp, Carol Thompson, Kristen Turner, Wayne Wakeman, Jennifer D. Whitesell,
793 Derric Williams, Ali Williford, Rob Young, Hongkui Zeng, Sarah Naylor, John W. Phillips, R. Clay
794 Reid, Stefan Mihalas, Shawn R. Olsen, and Christof Koch (Jan. 2021). “Survey of spiking in the mouse
795 visual system reveals functional hierarchy”. In: *Nature*. ISSN: 1476-4687. DOI: [10.1038/s41586-020-](https://doi.org/10.1038/s41586-020-03171-x)
796 [03171-x](https://doi.org/10.1038/s41586-020-03171-x). URL: <https://doi.org/10.1038/s41586-020-03171-x>.
- 797 Steinmetz, Nicholas A, Peter Zatzka-Haas, Matteo Carandini, and Kenneth D Harris (2019). “Distributed
798 coding of choice, action and engagement across the mouse brain”. In: *Nature* 576.7786, pp. 266–273.
- 799 Stringer, Carsen, Marius Pachitariu, Nicholas Steinmetz, Matteo Carandini, and Kenneth D Harris (2019).
800 “High-dimensional geometry of population responses in visual cortex”. In: *Nature* 571, pp. 361–365.
- 801 Sussillo, David and Omri Barak (Mar. 2013). “Opening the Black Box: Low-Dimensional Dynamics in
802 High-Dimensional Recurrent Neural Networks”. In: *Neural Computation* 25.3, pp. 626–649. ISSN: 0899-
803 7667. DOI: [10.1162/NECO_a_00409](https://doi.org/10.1162/NECO_a_00409). eprint: [https://direct.mit.edu/neco/article-pdf/25/3/
804 626/881886/neco_a_00409.pdf](https://direct.mit.edu/neco/article-pdf/25/3/626/881886/neco_a_00409.pdf). URL: https://doi.org/10.1162/NECO%5C_a%5C_00409.
- 805 Sussillo, David, Mark M Churchland, Matthew T Kaufman, and Krishna V Shenoy (2015). “A neural net-
806 work that finds a naturalistic solution for the production of muscle activity”. In: *Nature neuroscience*
807 18.7, pp. 1025–1033.
- 808 Teeter, Corinne, Ramakrishnan Iyer, Vilas Menon, Nathan Gouwens, David Feng, Jim Berg, Aaron Szafer,
809 Nicholas Cain, Hongkui Zeng, Michael Hawrylycz, et al. (2018). “Generalized leaky integrate-and-fire
810 models classify multiple neuron types”. In: *Nature communications* 9.1, pp. 1–15.
- 811 Thomson, Alex M and Christophe Lamy (2007). “Functional maps of neocortical local circuitry”. In:
812 *Frontiers in neuroscience* 1, p. 2.
- 813 Wang, Zhaoxiang, Zhouyan Feng, Yue Yuan, Gangsheng Yang, Yifan Hu, and Lvpioa Zheng (2022).
814 “Bifurcations in the firing of neuronal population caused by a small difference in pulse parameters
815 during sustained stimulations in rat hippocampus in vivo”. In: *IEEE Transactions on Biomedical*
816 *Engineering* 69.9, pp. 2893–2904.

- 817 Yang, Guangyu Robert, Madhura R Joglekar, H Francis Song, William T Newsome, and Xiao-Jing Wang
818 (2019). “Task representations in neural networks trained to perform many cognitive tasks”. In: *Nature*
819 *neuroscience* 22.2, pp. 297–306.
- 820 Yang, Guangyu Robert and Xiao-Jing Wang (2020). “Artificial neural networks for neuroscientists: A
821 primer”. In: *Neuron* 107.6, pp. 1048–1070.
- 822 Yuksel, Seniha Esen, Joseph N Wilson, and Paul D Gader (2012). “Twenty years of mixture of experts”.
823 In: *IEEE transactions on neural networks and learning systems* 23.8, pp. 1177–1193.
- 824 Zenke, Friedemann and Tim P Vogels (2021). “The remarkable robustness of surrogate gradient learning
825 for instilling complex function in spiking neural networks”. In: *Neural Computation* 33.4, pp. 899–925.

826 **Supplementary Information**

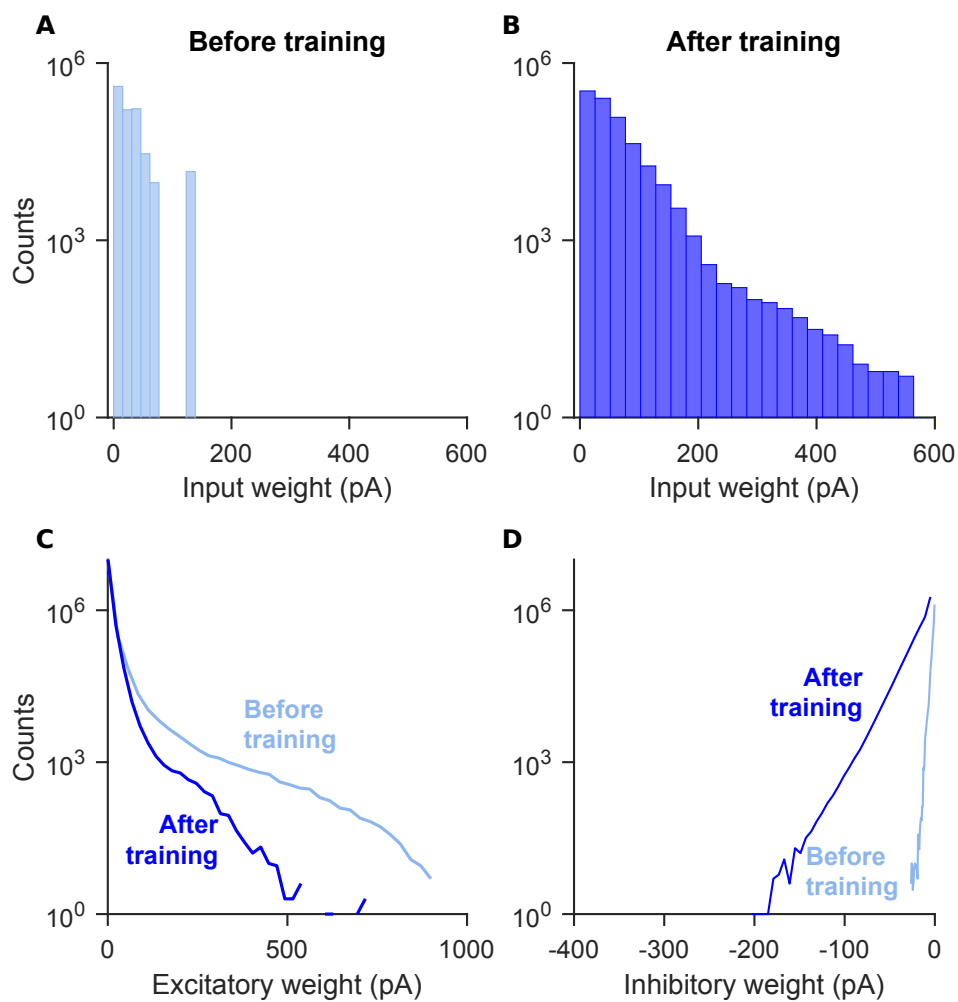


Figure S1: Changes in the distribution of synaptic weights through training. (A) Distributions of input weights before training. (B) Same as (A), but after training. The mean of input weights increases from 24.1 to 25.2 through training. (C) Distributions of excitatory weights in V1 model before and after training. (D) Distributions of inhibitory weights in V1 model before and after training. Note that the weights before training were given by (Billeh et al. 2020).

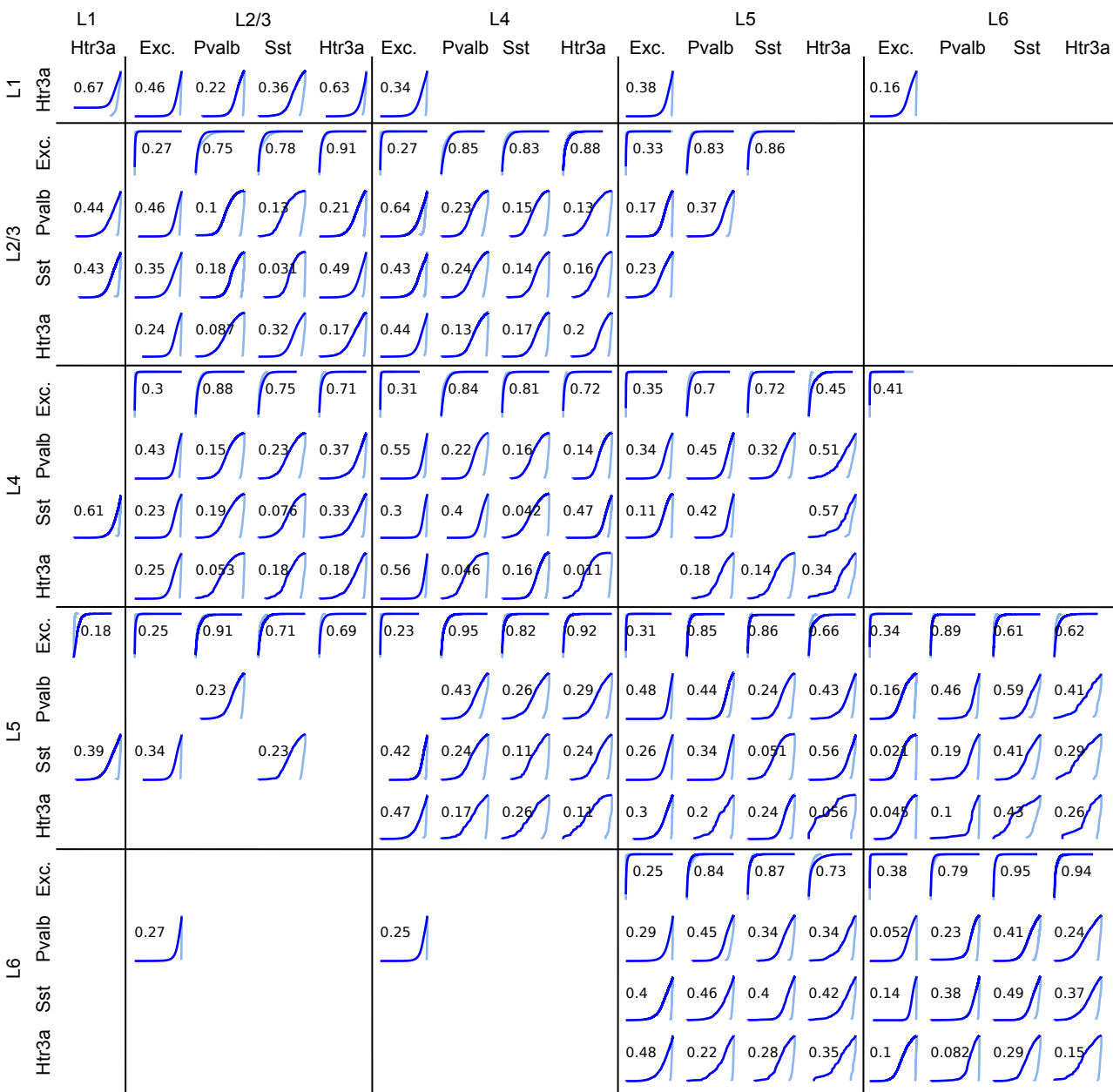
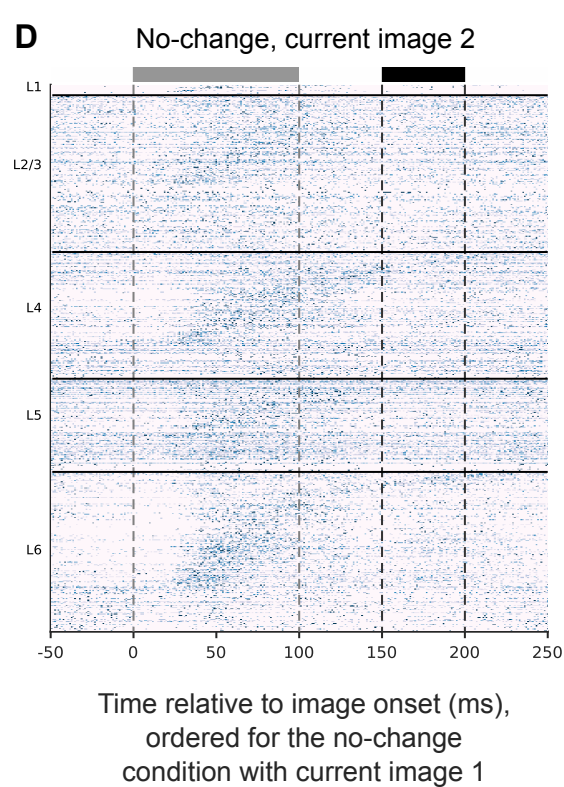
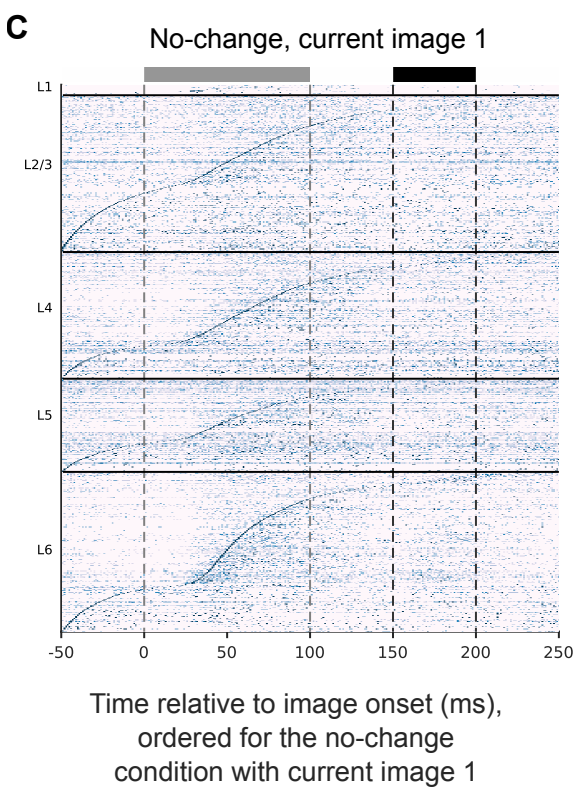
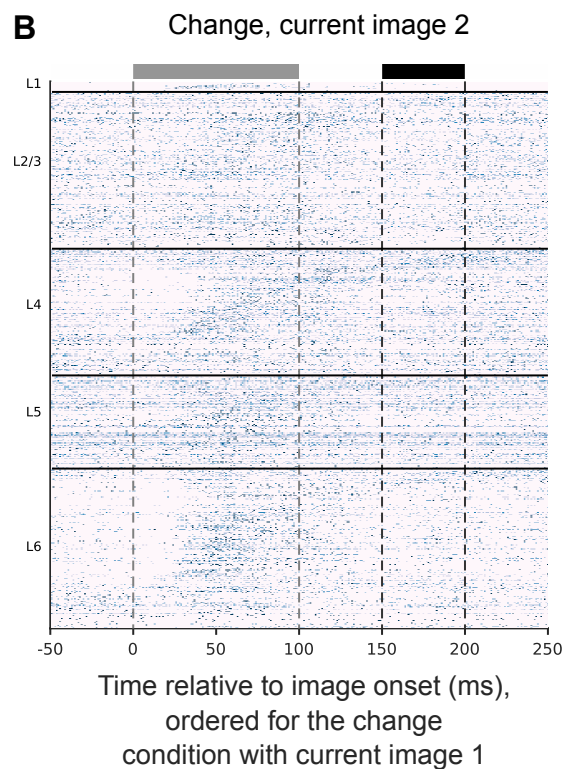
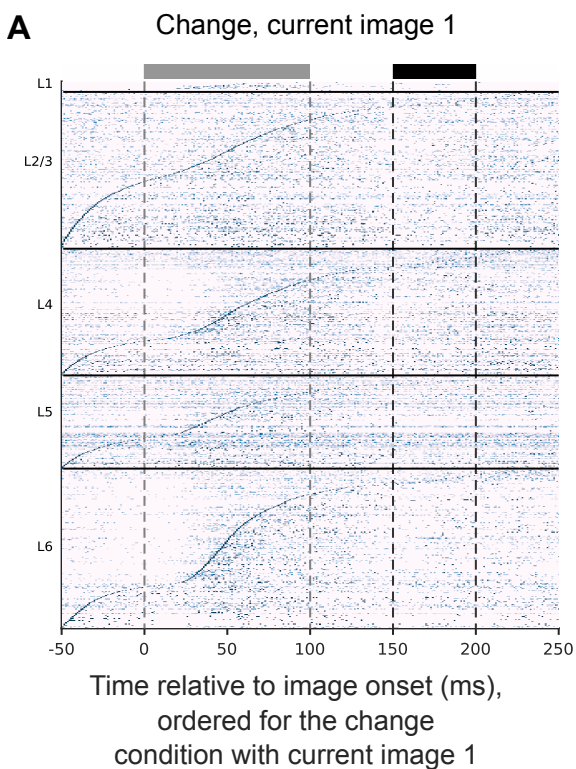


Figure S2: Distribution of recurrent synaptic weights between each pair of populations before (light blue) and after learning (dark blue). Each row represents a pre-synaptic neuron population, and each column represents a post-synaptic neuron population. The histogram represents the distribution of synaptic weights of all synaptic connections that share the same pre-synaptic and post-synaptic neuron population. Vertical axis in each panel is log-scale. Horizontal axis is linear scale and horizontal range is from the smallest value to the largest value of each population. The number is $1 - D$ where D is from the Kolmogorov-Smirnov test, quantifying the similarity between distributions (Billeh et al. 2020). Exc., excitatory neurons.



(caption next page)

Figure S3 (previous page): The temporal order of peak activity encodes in the V1 model also information about the identity of the current image. We have shown in Fig. 3 that the temporal order of peak activity in the V1 model has characteristic differences for different trial types (change or no-change). We show here that this temporal or rank-order coding is even more refined: the order contains in addition information about the identity of the current image. **(A)** Normalized average responses over 200 trials with the change condition and the same current image (image 1) but different preceding images, with neurons ordered according to the time of their peak activity. The gray and black bars at the top denote the image presentation and response windows, respectively. **(B)** Same as in **(A)**, but all trials have the same current image 2 and neurons were ordered as in **(A)**. The resulting blurred sequence indicates that the order of peak activity of neurons is different for images 1 and 2, also within the same trial type (change condition). **(C)** and **(D)** Same as in **(A)** and **(B)**, respectively, but for the no-change condition.

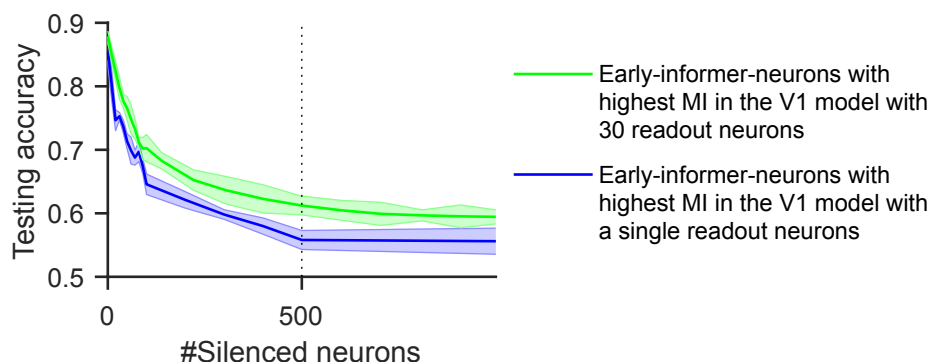
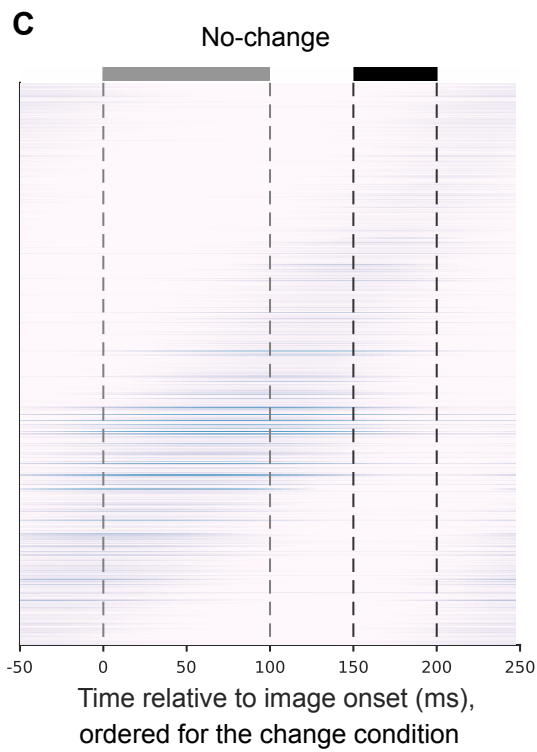
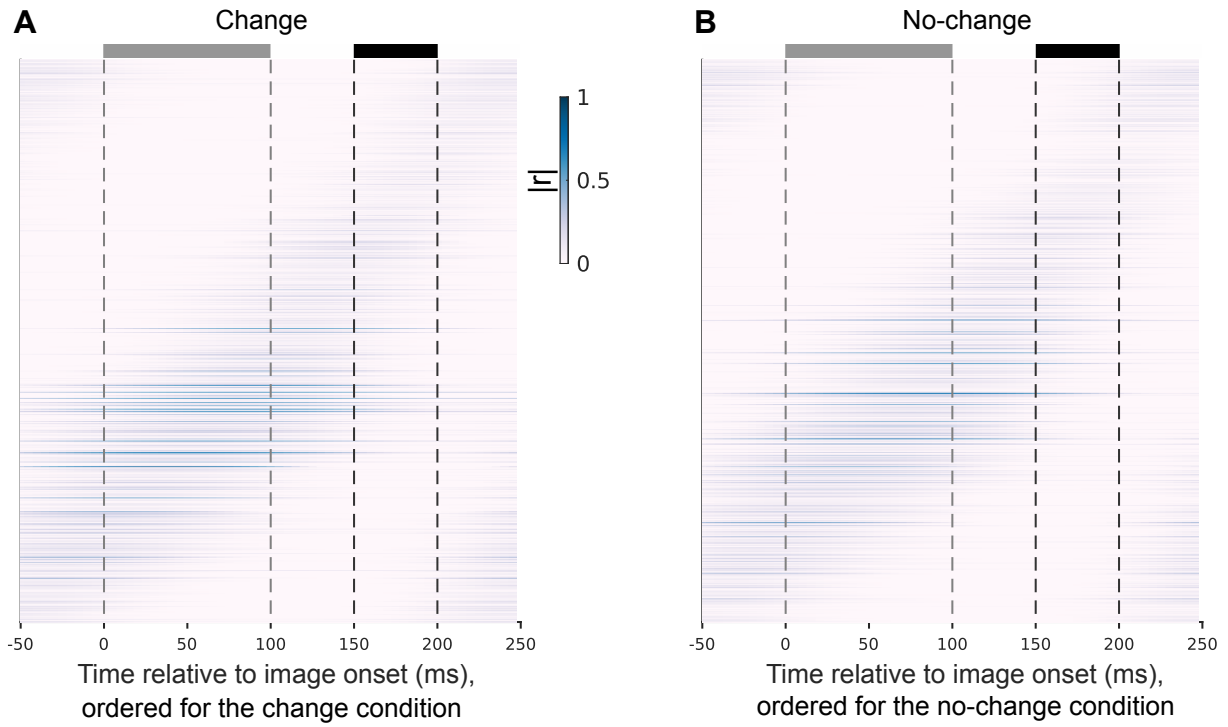
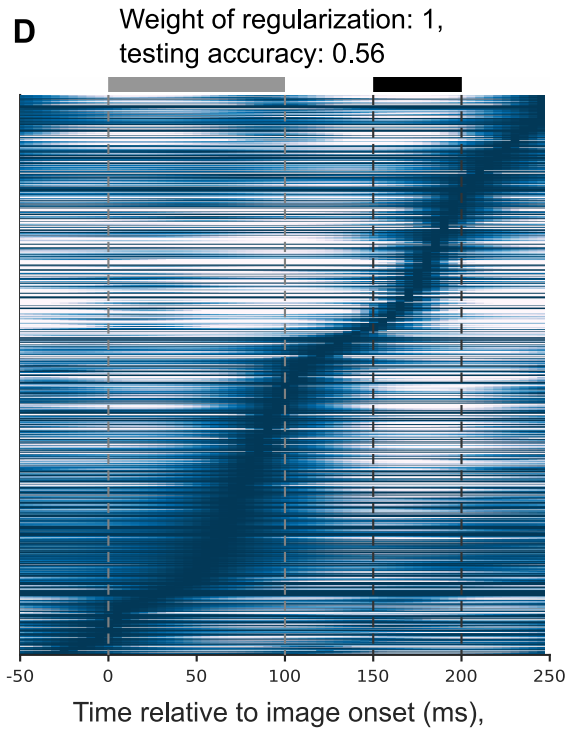
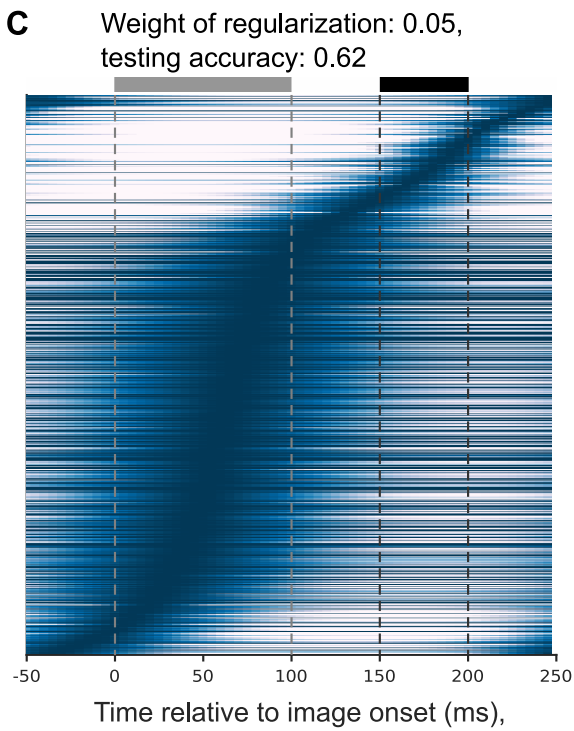
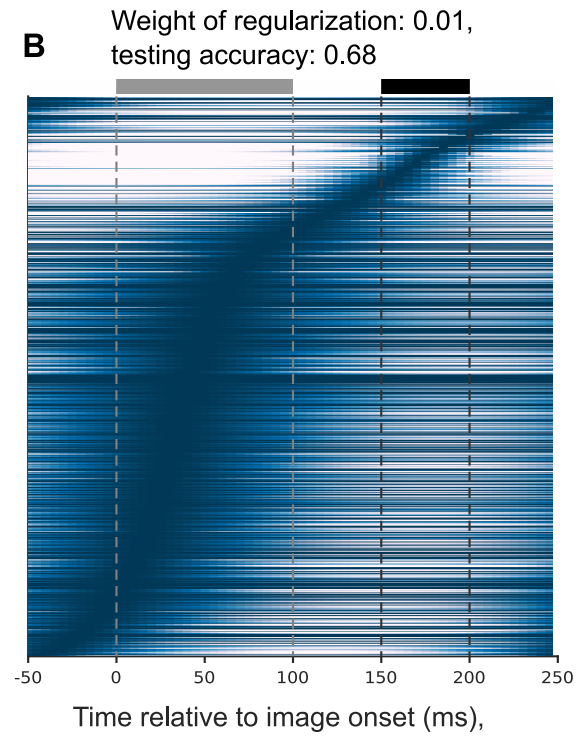
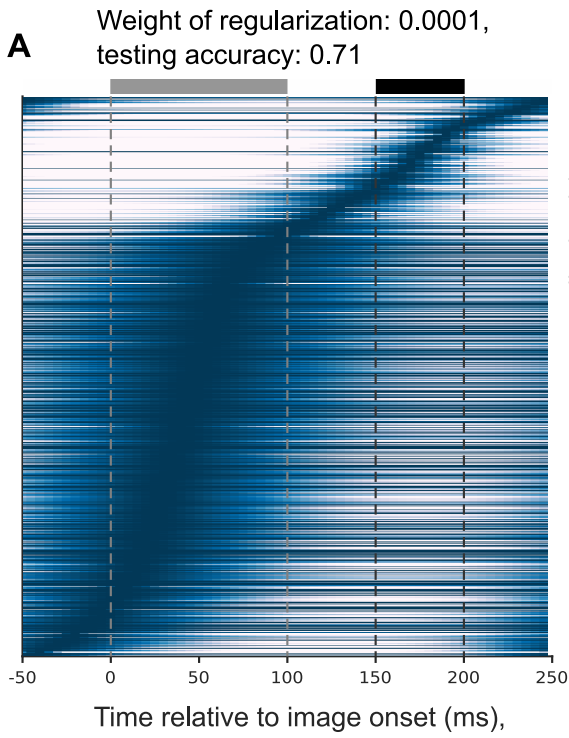


Figure S4: Lesion experiments for two versions of the trained V1 model with 30 and 1 readout neurons that are trained to report an image chance. Task performance quickly decreases when early-informer-neurons are silenced (in the order of their MI with the network decision) in both versions. The 30-readout-neuron model has also after silencing a given number of neurons slightly higher accuracy on test images than the single-readout-neuron model, consistent with Fig. 2F. Both curves show average values for 10 V1 models where different sets of training data were used. The shaded area represents the SEM across 10 models.



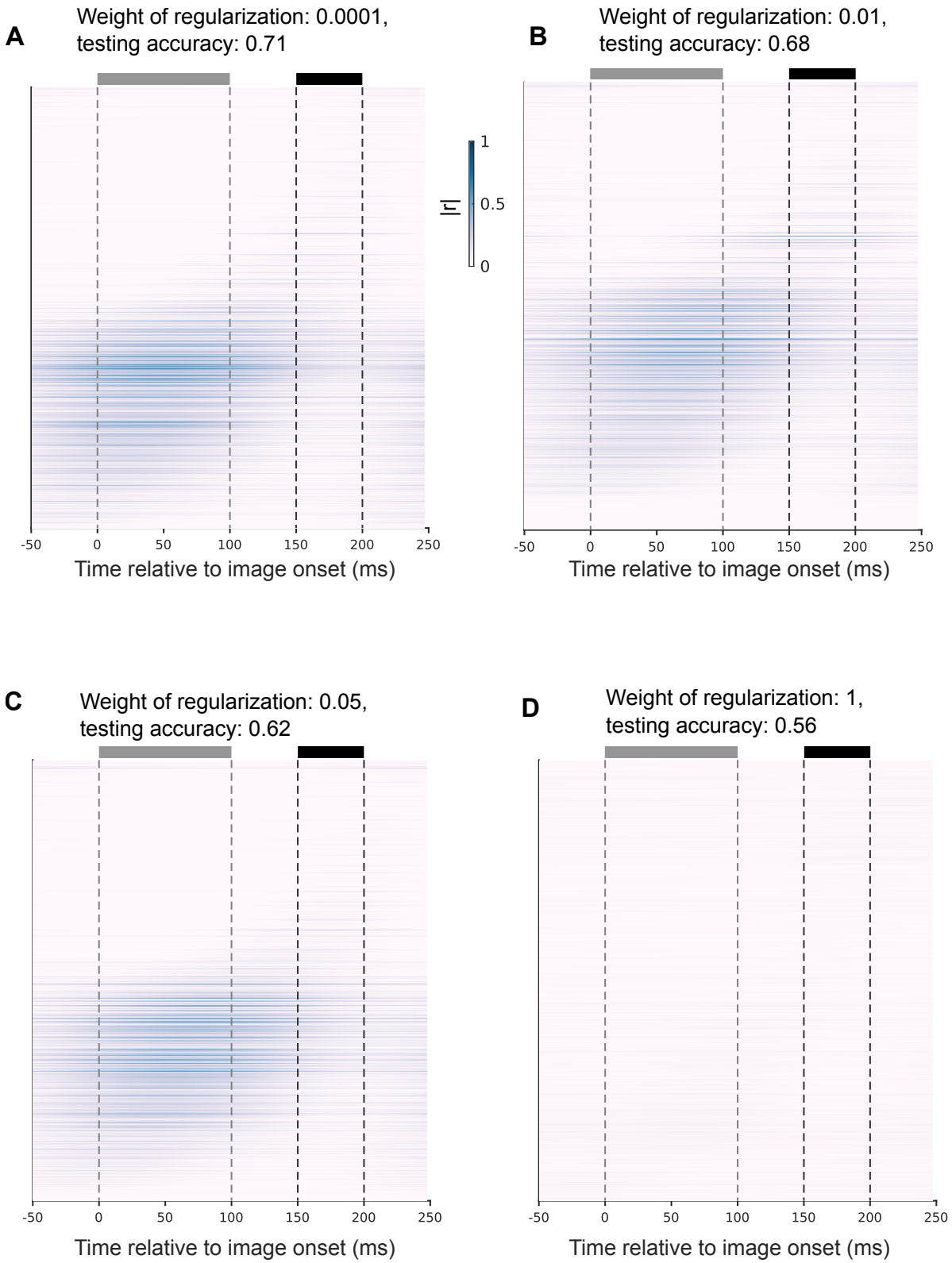
(caption next page)

Figure S5 (previous page): Another attempt to identify trial-type dependent sequential activity in the RANN, similarly as shown in Fig. 3 for the V1 model (A-C) Same RANN and conditions as in Fig. 4A-C but visualized without neuron-wise normalization. The fuzzy sequential order does not get sharper in panel (B) than in panel (C), although the neurons are ordered in (B) for this particular trial type (no-change).



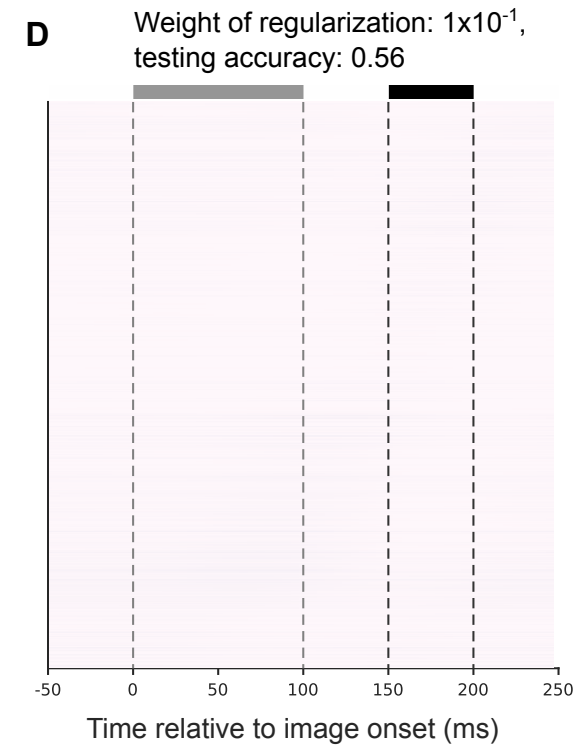
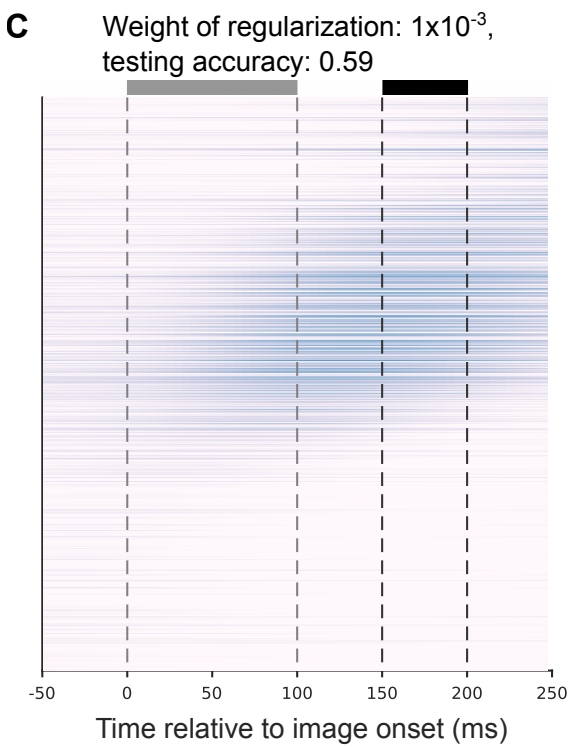
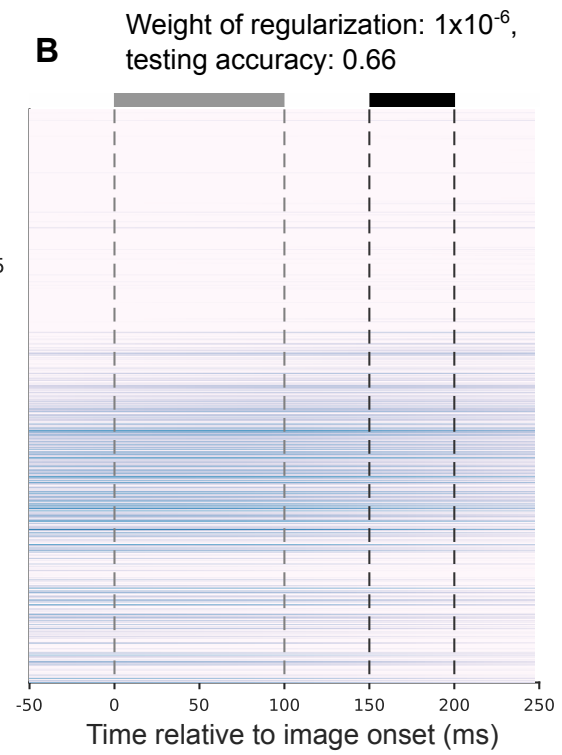
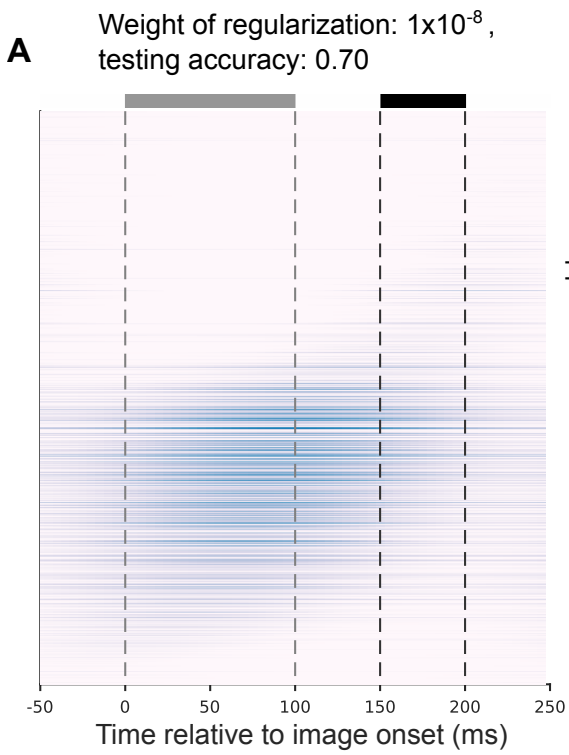
(caption next page)

Figure S6 (previous page): Different weights of the regularization term cannot produce sparse sequential neural activity while maintaining high task performance. (A-D) The normalized average responses of 200 change-condition trials in the RANN with different weights of regularization are plotted over time, with neurons ordered based on the time of their peak activities under the change condition. One sees that training the RANN for the same task with different weights of the regularization term (Methods) does not produce a sparse sequential neural activity as in the experimental data and the V1 model. Furthermore, more aggressive regularization strongly reduces task performance.



(caption next page)

Figure S7 (previous page): Plotting the same results as in Fig. S6 without neuron-wise normalization still does not indicate a RANN regime with sparse sequential neural activity. (A-D) Here the average absolute responses $|r|$ of 200 change-condition trials in RANN with different weights of regularization are plotted over time, with neurons ordered based on the time of their peak activities under the change condition.



(caption next page)

Figure S8 (previous page): Setting the value of the threshold in the activation regularization term to 0 also does not produce sparse sequential activity in the RANN. The threshold, θ in Eq. 16 was changed here from its default value 0.01 to 0. Average responses of 200 change-condition trials in RANN with different weights of the regularization term are plotted over time, with neurons ordered based on the time of their peak activities under the change condition. Similarly as in Fig. S7, none of the weights of the regularization term that we tried produces sparse sequential activity in the RANN while maintaining high task performance.

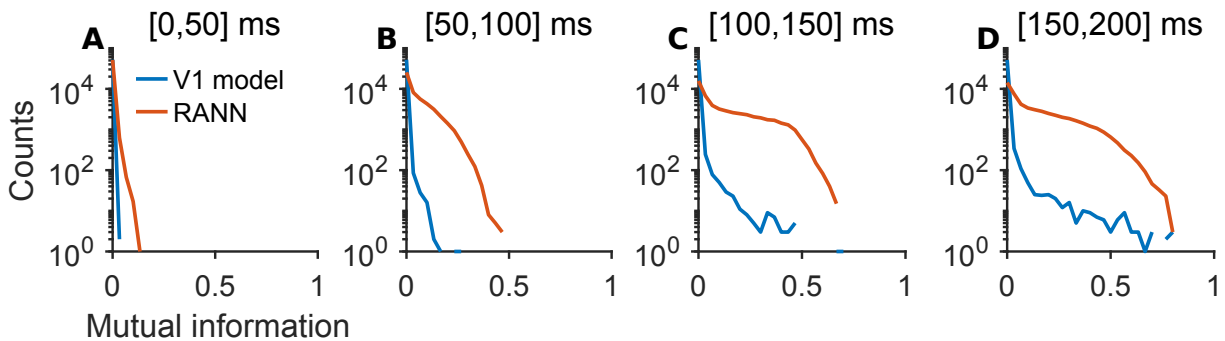
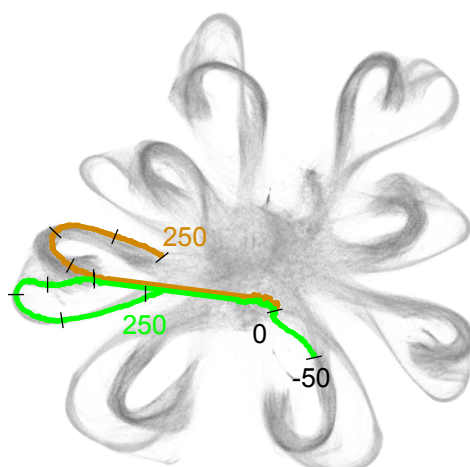


Figure S9: Comparison of mutual-information distributions among neurons in the V1 Model and the RANN. (A-D) The number of neurons that has a given level of MI with the network output is shown for each 50-ms window in the V1 model and the RANN (image onset was at time 0). One clearly sees that for each of these time windows that are by 1 or 2 orders of magnitude more neurons that have high MI with the network output. This explains also why so much more neurons need to be silenced in the RANN in order to reduce the task performance of the network to a given level.

Silencing of 100 early-informer-neurons also flips the decision from change to no-change



Preceding image 5 and current image 7

- in the intact V1 model
- when 100 early-informer-neurons are silenced

Figure S10: Silencing of 100 neurons can also flip a trajectory from the bundle for change to the no-change bundle of trajectories. In Fig. 5D, we demonstrated that silencing 100 early-informer-neurons can cause the trajectory of network states to flip from the bundle for no-change to the bundle for change trials. We show here that silencing of the same 100 neurons can also flip the network bifurcation in the other direction.

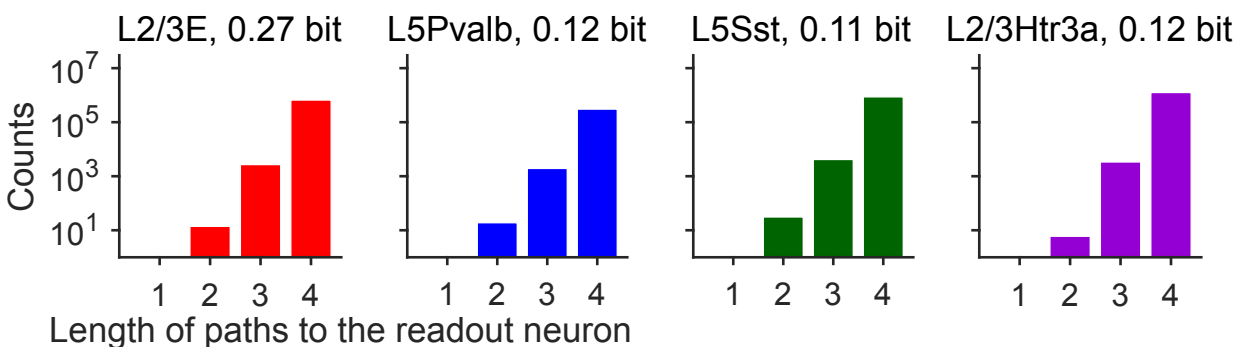


Figure S11: Histogram of the lengths (= number of synapses) of paths from 4 sampled early-informer-neurons with high MI to the readout neuron. We selected early-informer-neurons with the highest MI in four basic neuron types (excitatory, Pvalb, Sst, and Htra3 neurons). The paths from early-informer-neurons to the readout neuron were found by the MATLAB function “allpaths” in the directed graph. One can also find arbitrarily long paths; here we only demonstrate the short paths (length < 5) to the readout neuron for each of the early-informer-neurons. This distribution of path lengths suggests that the firing activity of an early-informer-neuron affects the membrane voltage of the readout neuron in multiple and diverse indirect ways.