

# The neurocellular implementation of representational geometry in primate prefrontal cortex

Xiao-Xiong Lin<sup>1,2</sup>, Andreas Nieder<sup>3</sup>, Simon N. Jacob<sup>1\*</sup>

<sup>1</sup> Translational Neurotechnology Laboratory, Department of Neurosurgery, Klinikum rechts der Isar, Technical University of Munich, Germany

<sup>2</sup> Graduate School of Systemic Neurosciences, Ludwig-Maximilians-University Munich, Germany

<sup>3</sup> Animal Physiology, University of Tübingen, Germany

\* Correspondence: [simon.jacob@tum.de](mailto:simon.jacob@tum.de)

## Summary

Modern neuroscience has seen the rise of a population-doctrine that represents cognitive variables using geometrical structures in activity space. Representational geometry does not, however, account for how individual neurons implement these representations. Here, leveraging the principle of sparse coding, we present a framework to dissect representational geometry into biologically interpretable components that retain links to single neurons. Applied to extracellular recordings from the primate prefrontal cortex in a working memory task with interference, the identified components revealed disentangled and sequential memory representations including the recovery of memory content after distraction, signals hidden to conventional analyses. Each component was contributed by small subpopulations of neurons with distinct electrophysiological properties and response dynamics. Modelling showed that such sparse implementations are supported by recurrently connected circuits as in prefrontal cortex. The perspective of neuronal implementation links representational geometries to their cellular constituents, providing mechanistic insights into how neural systems encode and process information.

## Keywords

Representational geometry; neuronal implementation; sparsity; working memory; prefrontal cortex; non-human primate

35 **Acknowledgements**

36 This study was supported by grants from the German Research Foundation (DFG JA 1999/1-  
37 1, JA 1999/5-1, JA 1999/6-1) and the European Research Council (ERC StG MEMCIRCUIT,  
38 GA 758032) to S.N.J and grants NI 618/10-1 and NI 618/13-1 to A.N.

39

40 **Author contributions**

41 X.-X.L. conceived the study and performed the analyses with contributions from S.N.J. A.N.  
42 and S.N.J. designed the experiments and collected the data. X.-X.L. and S.N.J. wrote the  
43 manuscript and prepared the figures. All authors edited the manuscript.

44

45 **Declaration of interests**

46 The authors declare no competing interests.

## 47 Introduction

48 For decades, the dominant approach to understanding neural systems has been to  
49 characterize the role and contributions of individual neurons. In a recent paradigm shift, the  
50 concept of high-dimensional activity spaces that represent cognitive and other variables at the  
51 level of neuronal populations has taken the center stage and sidelined the single-neuron  
52 perspective (Barack & Krakauer, 2021; Saxena & Cunningham, 2019). These population  
53 representations capture multi-neuron activity in different behavioral task conditions in the form  
54 of geometrical structures (Bernardi et al., 2020; Okazawa et al., 2021). Representational  
55 geometry provides a complete description of the information encoded by and processed in a  
56 neuronal population. It does not, however, account for how individual neurons – the nuts and  
57 bolts of brain processing – give rise to the representations and the operations performed on  
58 them (Kriegeskorte & Wei, 2021) because there is no direct connection between informational  
59 representation and biological implementation at the cellular and circuit level.

60 In constructing representational geometries, the choice of coordinate system, that is the set  
61 of components that capture the population activity, is arbitrary. The question then arises what  
62 the most meaningful coordinate system is to represent the data. In principal component  
63 analysis (PCA), a widely used method for dimensionality reduction, the principal components  
64 (PCs) capture the neuronal activity's variance, but they are not designed to yield biologically  
65 interpretable aspects of the representational geometry. Identifying coordinate systems that  
66 are rooted in biology is particularly relevant in association cortices where neurons often have  
67 mixed-selective responses that are not easily interpreted as the representation of any single  
68 stimulus or task variable alone (Bernardi et al., 2020; Rigotti et al., 2013). Neuronal signals in  
69 association cortices also show complex temporal dynamics and task-dependent modulations  
70 that reflect distinct sensory and memory processing stages (Cavanagh et al., 2018; Jacob et  
71 al., 2018; Jacob & Nieder, 2014). During working memory, for example, behaviorally relevant  
72 target items are maintained in online storage and must be protected against interfering  
73 distractors (Jacob et al., 2018; Jacob & Nieder, 2014). However, depending on which  
74 coordinate system is used to express the representational geometry, the same task-related  
75 neuronal activity could be interpreted in one of two ways: either as components representing  
76 the target in each task epoch individually, suggesting a memory mechanism built on sequential  
77 relay of target information among components (Parthasarathy et al., 2019), or, alternatively,  
78 as components that represent the target across task epochs, suggesting a memory  
79 mechanism of continuous representation of target information by the same components (Tang  
80 et al., 2020).

81 The biological implementation of representations points to how components are accessed and  
82 information is communicated. Unlike the units in neuronal network models, *in vivo* neurons  
83 are subject to anatomical and physiological constraints. There are approximately  $10^{10}$  neurons  
84 in the human brain and  $10^9$  in a hypothetical functional module such as the dorsolateral  
85 prefrontal cortex (PFC) (Courchesne et al., 2011; Herculano-Houzel et al., 2015). A pyramidal  
86 cortical neuron has on the order of  $10^4$  dendritic spines (Eyal et al., 2018). Thus, given the  
87 disproportion between the low number of possible connections and the large number of  
88 potentially informative neurons, a neuron downstream of the PFC can only 'read out' from a  
89 small fraction of neurons in this region. That is, it cannot access arbitrary components of the  
90 representational geometry. Instead, it would be more efficient and biologically plausible to read  
91 out components that a few neurons predominantly contribute to, that is the components with  
92 a sparse neuronal implementation.

93 Here, we present a framework that exploits the structure in the representational geometry's  
94 neuronal implementation. We show that this approach yields unbiased components of  
95 population activity that retain links to individual neurons. We performed data dimensionality  
96 reduction on extracellular multi-channel recordings from the non-human primate PFC by  
97 leveraging sparsity constraints in order to identify components that are contributed mainly by  
98 small subpopulations of strongly coding neurons (sparse component analysis, SCA; Georgiev  
99 et al., 2007; Olshausen & Field, 1996). We found that the activities on these components  
100 nontrivially matched the working memory task sequence performed by the animals, revealing  
101 separate sensory and memory components including a previously hidden component, namely  
102 the recovery of memory content after distraction. Notably, each component was made up of  
103 non-overlapping subpopulations of neurons with distinct electrophysiological properties and  
104 temporal dynamics. Finally, neuronal network modelling showed that recurrent connectivity as  
105 in the PFC favors such sparse implementations over non-structured Gaussian  
106 implementations. The framework and findings presented here bridge the gap between the  
107 single-neuron doctrine and the neuronal population doctrine (Barack & Krakauer, 2021;  
108 Saxena & Cunningham, 2019) and establish the perspective of neuronal implementation as  
109 an important complement to representational geometry.

## 110 Results

### 111 Different neuronal implementations may underlie the same representational geometry

112 Representational geometry abstracts the information coded by a population of neurons from  
113 their individual tuning profiles (Kriegeskorte & Wei, 2021). It specifies the pairwise distances  
114 between task-related collective neuronal responses, but no longer reflects the exact pattern  
115 of firing rates. This approach defines a stimulus-representing subspace. To illustrate, the  
116 representations for two stimuli A and B in PC space separate, rotate and collapse back to the  
117 origin (**Fig. 1a**).

118 The same stimulus-representing subspace can be defined with arbitrary sets of components.  
119 Components can be chosen to capture specific aspects of the representation, e.g., to  
120 continuously distinguish between stimuli (**Fig. 1b**), or to distinguish between stimuli at different  
121 time points (**Fig. 1c**). Note that in the former example, the components align with the PCs,  
122 while in the latter they do not. Various studies have followed this approach, selecting the  
123 components e.g. such that they express representations sequentially (Aoi et al., 2020) or such  
124 that they each correspond to a particular task variable of interest (Libby & Buschman, 2021;  
125 Mante et al., 2013).

126 Neuronal activity can be reconstructed by the weighted sum of components. Every neuron  
127 has a set of weights quantifying its relation to the different components, i.e. its loadings on the  
128 components. The loadings of neurons on the PCs visualize their positions in implementation  
129 space (**Fig. 1d-f**), where the loadings along any axis correspond to a component in  
130 representation space with the same orientation (**Fig. 1a-c**). The structure in the  
131 implementation space, i.e., the distribution of loadings across neurons, can be exploited to  
132 identify a unique, non-arbitrary set of components that emphasizes biological plausibility of  
133 stimulus coding over enforcing possibly unjustified priors.

134 Representational geometry is invariant to the rotation of neuronal coordinates (Kornblith et al.,  
135 2019). Different neuronal implementations may therefore underlie the same representational  
136 geometry. We first consider the scenario of a Gaussian (dense) distribution of loadings  
137 (**Fig. 1d**), where the standardized moments (e.g., skewness and kurtosis) are constant,  
138 meaning there are no differences in these distributional statistics across axis orientations. We  
139 define the sparsity index (SI; **Fig. 1d**, top inset) to denote the sparsity of the implementation  
140 along a given axis. SI is proportional to a distribution's kurtosis. If SI is constant across axis  
141 orientations, neurons do not preferentially align to any axes.

142 Next, we consider a sparse distribution (**Fig. 1e**). Most neurons lie around the origin of the  
143 coordinate system. However, because SI is not constant (**Fig. 1e**, top inset), we can find the  
144 sparse components that strongly coding neurons align to. In the present case, these sparse  
145 axes correspond to the components in representational space that code the difference  
146 between stimulus A and B continuously (with one of the components reversing between  
147 epochs; compare **Fig. 1e** with **Fig. 1b**). Importantly, sparse distributions can exist for arbitrary  
148 axis orientations. For example, strongly coding neurons could align to the components that  
149 sequentially represent the stimulus information at time point 1 and time point 2 (compare  
150 **Fig. 1f** with **Fig. 1c**).

151 Although both scenarios are characterized by sparse neuronal implementations, we note that  
152 they have fundamentally different implications for readout, lending particular importance to the  
153 positioning of sparse axes orientations. Continuous readout (**Fig. 1b** and **e**, component 1) is

154 stable, but not optimized for either time point 1 or time point 2, whereas sequential readouts  
155 (**Fig. 1c** and **1f**) are more precise at the respective time points, but not stable across time  
156 points.

157 In summary, the perspective of neuronal implementation offers a way to connect  
158 representational geometries to their cellular constituents, revealing mechanistic insights into  
159 how a neural system encodes, processes and relays information.

## 160 **The neuronal implementation of working memory**

161 With this framework, we now examine neuronal implementation of working memory, a core  
162 cognitive function for online maintenance and manipulation of information in the absence of  
163 sensory inputs. Extracellular multi-channel recordings were performed in the lateral PFC of  
164 two monkeys trained on a delayed-match-to-numerosity task, requiring them to memorize the  
165 number of dots (i.e., numerosity) in a visually presented sample and resist an interfering  
166 distracting numerosity (Jacob and Nieder, 2014) (**Fig. 2a**). A total of 467 single units recorded  
167 across 78 sessions were included in the analysis. Spike rates were binned, averaged across  
168 conditions of the same type and demixed into their constituent parts (**Fig. 2b**) (Kobak et al.,  
169 2016). Because the task design was balanced (i.e., all sample-distractor combinations were  
170 included), the different task variables were statistically independent of each other. Demixing  
171 therefore allowed to isolate and analyze signal components that would otherwise be  
172 overshadowed by signals that dominate the raw firing rates. Across neurons, the neuronal  
173 activities coding for trial time, sample numerosity, distractor numerosity and the sample-  
174 distractor interaction accounted for 72.7 %, 8.7 %, 5.8 % and 12.9 % of the total variance,  
175 respectively (**Fig. 2b**).

176 We first focused on the representation of the sample numerosity throughout the trial, the  
177 crucial function for completing the task (**Fig. 2c**). In PC space, the representations of different  
178 numerosities (1 and 4 visualized here) started to separate, marking an increase of the  
179 information during sample presentation. Then the representations rotated and returned to the  
180 origin. Similar representational changes have been reported previously (Elsayed &  
181 Cunningham, 2017; Murray et al., 2017; Parthasarathy et al., 2019).

182 The distribution of loadings of individual neurons onto the first three PCs was highly non-  
183 Gaussian ( $p < 0.001$ ; Henze-Zirkler multivariate normality test; **Fig. 2d**). Accordingly, the  
184 sparsity index (SI) was not uniform across all axis orientations (**Fig. 2d**). Using sparse  
185 component analysis (SCA) that identifies components with sparse distributions of neuronal  
186 loadings (sparse components, SCs), we found three SCs that optimally decomposed the  
187 sample numerosities' representational geometry. The SCs displayed temporally well-defined  
188 active periods that matched the task structure and tiled the duration of a trial (**Fig. 2e**).  
189 Intuitively, they correspond to components for sensory encoding, memory maintenance and  
190 memory recovery following distraction, in accord with the scenario of sequential  
191 representations (cp. to **Fig. 1c** and **f**).

192 To control for the possibility that noise in non-sparse implementations is mistaken for structure  
193 by SCA, we created substitute datasets with random Gaussian implementations (i.e.,  
194 Gaussian distributions of neuronal loadings) while keeping the representational geometry  
195 intact and then systematically compared the original SCs with the substitute SCs (example  
196 substitute SCs in **Fig. 2f**). First, the sparsity parameter  $\beta$  (fit to the distribution of loadings on  
197 the SCs) was smaller for all three original SCs than for the substitutes ( $p < 0.001$  for all three  
198 SCs; permutation test with  $n = 3 \times 1000$  permutations; **Fig. 2g**), confirming the presence of

199 structure in the implementation. Second, the activities on the SCs showed temporally  
200 restricted sample representations with shorter spread ( $p < 0.002$ ; permutation test with  
201  $n = 1000$  permutations; same as for Fig. 2i-k; **Fig. 2h**), less temporal overlap with other SCs  
202 ( $p < 0.003$ ; **Fig. 2i**), and less reversal of sample numerosity tuning ( $p < 0.030$ ; **Fig. 2j**) than  
203 the substitutes, suggesting that the observed SC activity was more sequential than to be  
204 expected with a random implementation. Third and finally, the SCs were closer to orthogonal  
205 than the substitutes ( $p < 0.019$ ; **Fig. 2k**), demonstrating that the observed implementation is  
206 more efficient than a random implementation.

207 In summary, the neuronal implementation of the sample numerosities' representational  
208 geometry was structured and sparse. The activities on the sparse components demonstrated  
209 sequential rather than continuous coding of working memory content, indicating that the  
210 change of behavioral demands in the course of the trial triggers a switching of informative  
211 subpopulations.

### 212 **The effect of distraction on sample numerosity representations**

213 The lack of a component that continuously represented the behaviorally relevant sample  
214 numerosity throughout the trial was unexpected. We therefore investigated the influence of  
215 distraction on sample number coding.

216 First, we applied SCA to the demixed distractor coding part of the data (**Fig. 3a**, top). Two SCs  
217 were obtained that were sequentially active during presentation and maintenance of the  
218 distractor numerosity, respectively (**Fig. 3a**, bottom). These components resembled the  
219 sensory and memory sample coding SCs (cp. to **Fig. 2e**), suggesting that target and  
220 distracting information initially occupied similar resources despite their distinct behavioral  
221 relevance. Supporting this hypothesis, we found strongly overlapping neuronal loadings  
222 between sample SCs and distractor SCs (cosine similarity; 0.69 and 0.57 for the sensory and  
223 memory components, respectively; **Fig. 3b**) with displacement of sample information by  
224 distractor information as the trial evolved (**Fig. S1a**, top and middle). However, in contrast to  
225 the sample sensory and memory components, the sample recovery SC was unique and did  
226 not share loadings with any other SC (**Fig. 3b**). Furthermore, the sample recovery SC was not  
227 influenced by distractor information and carried sample information until test numerosity  
228 presentation (**Fig. S1a**, bottom). To correctly complete a trial, more activity in the sample  
229 sensory and recovery SCs was required when the trial contained a distractor than when a trial  
230 without a distractor was presented (**Fig. S1b**). Conversely, distractors led to reduced sample  
231 activity in the memory component.

232 Second, we applied SCA to the sample-distractor interaction part of the data. One SC was  
233 identified. Its activity was most pronounced when the sample and distractor numerosity were  
234 the same (**Fig. S2**). The neuronal loadings on this SC did not overlap with the loadings on  
235 sample or distractor SCs (**Fig. 3b**), suggesting that the boost in numerosity information was  
236 generated by a dedicated subpopulation responding to a repeated presentation of the same  
237 number, instead of changing the activity of the sample representing neurons.

238 Together, these results indicate a (partially) shared capacity for sample and distractor  
239 representations during the sensory input and subsequent memory delay stages. The invasion  
240 of distractor information forced the recruitment of an extra component, the recovery  
241 component, to maintain sample information in working memory.

242 So far, all analyses were performed on separated (demixed) representations. We next  
243 investigated whether sample and distractor information could be equally disentangled using  
244 SCA alone without demixing the numerosity coding signal (**Fig. 3c**). SCA performed on firing  
245 rates averaged across the second memory delay recovered two sparse components that each  
246 selectively captured sample and distractor information (**Fig. 3d**). The corresponding  
247 representational geometry was grid-like with clearly factorized sample and distractor  
248 information that each aligned well to one SC (**Fig. 3e**). Notably, this alignment was non-trivial  
249 and not enforced by our analytical method, arguing that the PFC spontaneously disentangles  
250 target and distractor representations in working memory. The underlying implementation  
251 showed clear sparse structure in the neuronal loadings onto these components (**Fig. 3f**).

252 For comparison, PCA, which is insensitive to the neuronal implementation, was unable to  
253 recover factorized components (**Fig. 3g**). The grid-like geometry was still largely preserved,  
254 but it did not align with the PCs (**Fig. 3h**). In contrast to SCA, PCA did not identify the  
255 components with the sparsest loadings (**Fig. 3i**).

### 256 **Subpopulations of neurons dominate working memory representations**

257 Next, we investigated whether the implementation was sparse enough to be able to reliably  
258 reconstruct the population-level sample representation using only a small fraction of neurons.  
259 We performed cross-temporal linear discriminant analysis (LDA) to decode sample numerosity  
260 at a given time point in the trial using training data from a different time point (**Fig. 4**). Decoding  
261 accuracy therefore quantifies the degree to which the representation is transferable. With four  
262 numerosities, chance level accuracy is 25 %. Using the entire population of 467 recorded  
263 neurons, we found a highly dynamic code with good within-epoch transfer, but very little  
264 generalization across epochs, in particular from the first to the second memory delay (**Fig. 4a**).  
265 In line with our previous results, this finding suggests that working memory representations  
266 are non-uniform and that distinct, complementary processes are required to protect  
267 behaviorally relevant information from interference.

268 We selected the neurons that contributed most to the previously identified SCs (loading on the  
269 SC larger than two standard deviations; **Fig. 4b**). 36, 28 and 28 single neurons passed the  
270 criterion for the sensory, memory and recovery SC, respectively. Although each subpopulation  
271 comprised only 6 to 8 % of the entire recorded population, these 'dominant neurons' explained  
272 88 %, 82 % and 87 % of their respective component's variance (sum of squares of dominant  
273 neurons' loadings over sum of squares of all neurons' loadings). Overlapping membership in  
274 two subpopulations was very rare (no more than three neurons in any SC pair; **Fig. 4b**).

275 Cross-temporal LDA using only the dominant neurons showed a very similar sample  
276 numerosity decoding pattern as with the entire population (**Fig. 4c**, cp. with **Fig. 4a**),  
277 confirming that the decoder previously relied mainly on this small subset of neurons. The  
278 sensory subpopulation contributed to decoding in particular during the sample and test  
279 numerosity presentation, but showed very little activity in the memory epochs (**Fig. 4d**, top).  
280 The memory subpopulation dominated in the first delay, but surprisingly was not involved in  
281 sample coding during the second delay (**Fig. 4d**, middle). Instead, after distraction, the  
282 recovery subpopulation was exclusively responsible for carrying sample information (**Fig. 4d**,  
283 bottom). This suggests that these neurons crucially contribute to shielding working memory  
284 information from interference (see also **Fig. S1**).

### 285 **Subpopulation-specific electrophysiological properties**



286 Above, we identified dominant neurons based on their stimulus selectivity. We now  
287 investigated whether their different roles in representing sample information were possibly  
288 mirrored by distinct electrophysiological properties.

289 First, we calculated the across-trial similarity (Pearson correlation) between each neuron's  
290 activity at different time points in the fixation period in order to derive the intrinsic time scale,  
291 a measure considered to index a neuron's ability to maintain memory traces (Murray et al.,  
292 2014). Representative neurons from all three subpopulations are shown (**Fig. 5a**). The  
293 example recovery neuron had a significantly larger spread from the diagonal than the sensory  
294 and memory neuron, i.e., its activity in distant time points was more strongly correlated, thus  
295 signifying a longer time constant (**Fig. 5a**, bottom panel). For each subpopulation, an  
296 exponential decay was fitted to the mean correlation coefficient across neurons (**Fig. 5b**). The  
297 recovery subpopulation had the largest time constant  $\tau$  (165 ms, 127 ms, and 338 ms for  
298 sensory, memory and recovery neurons, respectively). The distribution of  $\tau$  values in the  
299 recovery population also stood out from the distributions observed in subsampled  
300 subpopulations of PFC neurons, whereas the sensory and memory neurons' distributions  
301 were not significantly different ( $p = 0.874$ ,  $p = 0.455$ ,  $p = 0.002$  for sensory, memory and  
302 recovery subpopulations, respectively; KL-divergence with bootstraps; **Fig. 5c**).

303 Next, we investigated spike train statistics using the inter-spike intervals (ISI) measured during  
304 the neurons' entire recording lifetime. The coefficient of variation (CV) measures the  
305 irregularity of a spike train (**Fig. 5d**). CVs of all recorded neurons were larger than 1 (i.e., more  
306 irregular than a Poisson process) with a gradual increase of spiking irregularity across the  
307 sensory, memory and recovery subpopulations. CVs in the recovery neuron population were  
308 significantly larger than in the sensory subpopulation ( $p = 0.030$ , two-tailed  $t$ -Test; **Fig. 5d**).  
309 The local variation (LV) measures local ISI differences and complements CV, which is a global  
310 measure. LVs in all dominant neurons were smaller than 1 (i.e., less local variation than a  
311 Poisson process) and significantly lower than in the non-coding PFC population ( $p < 0.001$ ,  
312 two-tailed  $t$ -Tests; **Fig. 5e**).

313 Notably, these distinct electrophysiological properties were not involved in the original  
314 selection of subpopulations and therefore lend support to the notion that the implementation  
315 structure carries biological meaning.

### 316 **Subpopulation-specific temporal dynamics and representation of context**

317 There was no perceptual cue in the working memory task specifying the difference between  
318 sample and distractor. This forced the animals to internally keep track of a trial's temporal  
319 evolution. To investigate whether temporal dynamics and context played a role in supporting  
320 the subpopulation-specific stimulus representations, we next analyzed the temporal part of the  
321 demixed signal and visualized condition-averaged activity trajectories in each of the dominant  
322 subpopulations (**Fig. 6a**).

323 In the sensory subpopulation, the trajectory followed a periodic, quasi-circular course (**Fig. 6a**,  
324 top panel). The first and second memory epochs overlapped almost entirely. This indicates  
325 that the sensory neurons did not distinguish between the time periods after sample and after  
326 distractor presentation. The trajectory of the memory subpopulation was less periodic, but  
327 intertwined in the first and second memory epochs (**Fig. 6a**, middle panel). In contrast, the  
328 trajectory of the recovery subpopulation was less intertwined, with most time points  
329 distinguishable from each other, especially the first and second memory epochs, signifying a

330 better representation of the contextual difference following sample and distractor presentation  
331 (**Fig. 6a**, bottom panel).

332 Overlap of the memory epochs in the sensory and memory subpopulations could be due to  
333 the limitations of a linear projection and the emphasis of PCA on global structure. We therefore  
334 performed non-linear embedding using t-SNE (**Fig. 6b**). This analysis revealed comparable  
335 structures as the linear projection, with the first and second memory epochs separated only in  
336 the recovery neuron subpopulation.

337 To further investigate the temporal evolution of neuronal activity, we measured the Euclidean  
338 distances between individual time points in each subpopulation (full space; **Fig. 6c**). All  
339 distance matrices displayed a strong diagonal, reflecting the fact that close-by time points  
340 were represented similarly. Notably, there were also strong offset diagonals in the sensory  
341 subpopulation, meaning that activity in these neurons repeated with a cycle of about 1.5 s.  
342 Furthermore, activity in the sensory and memory epochs differed most in this subpopulation.  
343 These patterns were present, albeit weaker, in the memory subpopulation, but absent in the  
344 recovery neurons. We quantified periodicity for each neuron by computing the relative power  
345 of 1/1.5 s (0.67 Hz) activity and its harmonics normalized to the power of the full frequency  
346 spectrum (**Fig. 6d**). Compared to randomly sampled subpopulations of PFC neurons, the  
347 sensory subpopulation and the recovery subpopulation showed significantly different (higher  
348 and lower, respectively) periodicity ( $p < 0.001$ ,  $p = 0.051$ ,  $p = 0.043$  for sensory, memory and  
349 recovery subpopulations, respectively; KL-divergence with bootstraps; **Fig. 6d** inset).

350 Neuronal activity is not static and temporally independent. Instead, firing rates at every time  
351 point depend on previous time points. To characterize the dynamical properties of the  
352 recorded PFC population in more detail, we used the measure of tangling (Russo et al., 2018).  
353 Tangling measures the extent to which the velocity (direction and speed) of a given state on  
354 a trajectory diverges from the velocity of its neighboring states (**Fig. 6e**), reflecting the level of  
355 unpredictability and instability (chaos) in the system. High tangling means a small disturbance  
356 in the current state would lead to large changes in the next state (difference of derivatives of  
357 neighboring points). The instability or inability to determine the next state from the current state  
358 (i.e., high tangling) indicates that other neuronal populations or external stimuli may drive the  
359 trajectory. Consequently, tangling was increased following the onset and offset of sensory  
360 input in all three subpopulations. Tangling was highest, however, in the sensory subpopulation  
361 and lowest in the recovery subpopulation (sensory vs. memory,  $p < 0.001$ ; memory vs.  
362 recovery,  $p = 0.013$ ; two-tailed  $t$ -Test across all trial time points; **Fig. 6f**).

363 In summary, these results suggest that the subpopulation of recovery neurons keeps a record  
364 of time and temporal context, which could contribute to these neurons' ability to separate  
365 sample and distracting information. In contrast, the sensory subpopulation - and the memory  
366 subpopulation to a lesser degree - is characterized by its strong input-driven temporal  
367 dynamics, which is consistent with these neurons' passive representation of numerosity  
368 regardless of it being behaviorally relevant (sample) or irrelevant (distractor).

### 369 **Recurrent connectivity favors sparse implementations**

370 The implementation underlying the temporal evolution of neuronal representations is not  
371 arbitrary, but must be derived from the dynamical system of constituent neurons and their  
372 anatomical connectivity pattern. The PFC is a highly recurrent, rather than purely feed-forward,  
373 brain region (Harris et al., 2019). If biological structure and resource efficiency indeed favor

374 sparse implementations, these should be better captured by recurrently connected networks  
375 than non-structured Gaussian implementations.

376 To address this hypothesis, we constructed a recurrent neural network model (RNN) to  
377 reproduce the target (to-be-fitted) firing rate sequences of each sample-distractor combination  
378 (**Fig. 7a**). The model consists of 467 neurons (to match the recorded population) receiving  
379 inputs of stimulus information according to the task structure. The model learns the recurrent  
380 connectivity  $W$  among the neurons.  $W$  summarizes the influence of the current time point's  
381 firing rates  $r$  on the firing rates of the next time point. An indicator vector  $n$  (one non-zero entry)  
382 represents the sample and distractor numerosity, activating the numerosity-specific input in  $I$   
383 to the entire neuronal population. To reflect the absence of an explicit visual cue that  
384 differentiates between sample and distractor in the task design, sample and distractor  
385 numerosity share the same input channel ( $I, n$ ). The contextual difference is left for the model  
386 to resolve. The intercept term  $b$  captures the baseline activity of each neuron.

387 We first trained the model on the original dataset and visualized the trajectory of the output  
388 averaged across all conditions (**Fig. 7b**). The model reproduced the original dataset well,  
389 capturing 85.7 % of total variance. Next, we created substitute datasets with altered  
390 implementations of numerosity representations ( $x_{\text{sample}} + x_{\text{distractor}} + x_{\text{SD interaction}}$ ) for the model to  
391 fit. The temporal part of the demixed data was unchanged. Three different implementations  
392 were created: first, a non-structured Gaussian distribution of neuronal loadings and no  
393 alignment to any components (cp. **Fig. 1d**); second, a distribution with the same degree of  
394 sparsity as the original data, but with sparse axes randomly rotated to align to other  
395 components (cp. **Fig. 1e**); third, a substitute with the same sparse distribution of neuronal  
396 loadings as in the original data (cp. **Fig. 1f**).

397 The model captured an increasing proportion of variance of the full signal across the three  
398 substitutes ( $p < 0.001$ ; one-way ANOVA; **Fig. 7c**). The absolute differences in explained  
399 variance were comparatively small (left axis), but remarkable in relation to the variance of the  
400 manipulated signal (right axis) and given that the representational geometry was unchanged  
401 and identical for all substitutes (cp. **Fig. 1**). A comparable result was obtained for the explained  
402 variance of the numerosity coding part ( $p < 0.001$ ; one-way ANOVA; **Fig. 7d**).

403 Taken together, these results demonstrate that sparse implementations of working memory  
404 representations are favored by recurrent circuits, the characteristic wiring motif of association  
405 cortices such as the PFC.

## 406 **Discussion**

407 We presented a framework to examine the contributions of individual neurons to population-  
408 level responses in representation space and to utilize its implementation structure. We  
409 identified heavy-tailed, i.e., sparse distributions of neuronal loadings on components that  
410 captured disentangled and sequential memory representations including the recovery of  
411 memory content after distraction. The switching of working memory components circumvented  
412 interference. These components could be traced to small subpopulations of neurons with  
413 distinct electrophysiological properties and temporal dynamics. Modelling showed that such  
414 sparse implementations with sequentially active components are supported by recurrently  
415 connected networks.

## 416 **Bridging population activity and neuronal implementation**

417 Population-level activity and representational geometry were previously studied without  
418 forming direct links to individual neurons (Bernardi et al., 2020; Chung & Abbott, 2021;  
419 Kriegeskorte & Wei, 2021; Okazawa et al., 2021). However, while single-neuron selectivity  
420 measures have the advantage of being more easily connected to biological properties such  
421 as cell type, receptor expression and axonal projection targets, they are typically chosen  
422 based on intuition and past experience and only partially or indirectly reflect the full  
423 representational space (Hirokawa et al., 2019; Jacob & Nieder, 2014).

424 Our sparse component analysis (SCA) framework (**Fig.1**) combines the advantages of both  
425 perspectives. It builds on representational geometry for a comprehensive account of the data  
426 and then links the relevant coding dimensions in the activity space to populations of strongly  
427 contributing neurons, which allows relating the population-wide activity patterns to tangible  
428 physiological measures.

## 429 **Implementation reveals biologically relevant dimensions in activity space**

430 Without respecting implementation, selecting components in activity space for further analysis  
431 is arbitrary. It is often done post-hoc after visualizing the top PCs, or by relying on the heuristics  
432 of 'what should be coded' in the system (Aoi et al., 2020; Bernardi et al., 2020; Libby &  
433 Buschman, 2021). This approach becomes problematic when the dimensionality is too high  
434 or when too many variables are involved.

435 By exploiting neuronal implementation, SCA identifies activity components in an un-biased  
436 and non-arbitrary way. SCA can therefore capture a more complete set of stimulus-associated  
437 variables (dimensions), most notably the temporal modulation of stimulus coding. This  
438 reduces bias otherwise introduced by selecting specific time windows, across which neuronal  
439 activity is averaged, and acknowledges the role of different response dynamics for information  
440 coding (Bondanelli & Ostojic, 2020; Mante et al., 2013). Furthermore, incorporating temporal  
441 modulation renders analyses more robust to noise (Johnstone & Lu, 2009), which is usually  
442 Gaussian and could hide the structure in implementation.

443 The implementation's sparse structure is a result of biological constraints regarding the  
444 connections among individual neurons. The approximately  $10^4$  dendritic spines on each  
445 cortical neuron (Eyal et al., 2018) define an upper limit for the number of neurons it could read  
446 out from. The  $10^9$  neurons in a cortical region such as human PFC (Courchesne et al., 2011;  
447 Herculano-Houzel et al., 2015), and even sub-modules with one to two magnitudes fewer  
448 neurons, therefore cannot be reached directly. The addition of one connection step would  
449 allow reaching the majority of PFC neurons, but at the cost of producing a layer of  $10^4$  to  $10^5$

450 neurons that are dedicated exclusively to feeding the single hypothetical downstream neuron.  
451 This is prohibitively inefficient. In such polysynaptic chains, it is more likely that meaningful  
452 representations have already emerged in intermediate layers as a result of direct connections  
453 from the source region. This notion is also in line with the high dimensionality and non-linear  
454 mixed selectivity characteristic of PFC, which allow for direct linear readout of complex  
455 representations without further computations (Rigotti et al., 2013).

456 Neurons share inputs and have local recurrent connections, which are particularly pronounced  
457 in association cortices such as the PFC (Harris et al., 2019), resulting in more similar firing  
458 patterns among neurons within cortical regions. Consequently, neurons might display activity  
459 that is weakly correlated to some components of the representational geometry even though  
460 they do not participate in the readout. This emphasizes the importance of truncating neurons  
461 with weak loadings and enforcing sparsity constraints for estimating potential readout  
462 connections (**Fig. 4**) and motivates the use of dynamical systems modelling to validate  
463 correlative measures (**Fig. 7**).

#### 464 **Working memory persistence without neuronal persistence**

465 Applied to working memory maintenance in the face of distraction, our framework uncovered  
466 an unexpected sequential representation of numerosity information across multiple task  
467 epochs (**Fig. 2**). This result was neither encouraged nor guaranteed by SCA. This suggests  
468 that the readout of memory content from the PFC is optimized for accuracy in each behavioral  
469 context rather than optimized for stability across time periods. The distractor occupied the  
470 same resources as the sample numerosity with regard to the sensory and memory component  
471 (**Fig. 3**), forcing behaviorally relevant information to be shifted to the recovery component  
472 following distraction. Thus, working memory content was maintained by distinct mechanisms  
473 before and after interference (**Fig. 4**).

474 The subpopulation of recovery neurons was characterized by electrophysiological properties  
475 that set these neurons apart from the other populations and could render them particularly  
476 suited to working memory storage. Their longer intrinsic timescales (**Fig. 5**) suggest more  
477 stable memory retention (Kim & Sejnowski, 2021; Murray et al., 2014). These neurons also  
478 distinguished between sample and distractor contexts, which is crucial for determining what  
479 information to keep and what information to discard (**Fig. 6**). The contextual signal was  
480 additively mixed with the numerosity coding signal in these neurons, but might still act as gain  
481 modulation for numerosity information given the neuronal input-output non-linearity (Dubreuil  
482 et al., 2020).

483 Representing memory content by sequentially active subpopulations is advantageous. With  
484 relay of information, a result of locally feed-forward connectivity, a network can maintain  
485 multiple inputs from previous time points and show more resistance to noise (Orhan & Pitkow,  
486 2020). Furthermore, the PFC might be non-linearly mixing context and memory  
487 representations in all possible ways, expanding dimensionality to enable flexible readout  
488 (Rigotti et al., 2013). Extensive training could have strengthened the non-linear mixture of  
489 second memory epoch context and sample numerosity representations that was most  
490 important in the current task, with the PFC retaining other mixtures (e.g. the component coding  
491 for sample numerosity in the first memory epoch) for other behavioral demands. In this view,  
492 the subpopulation of memory neurons could function as a more passive short-term memory  
493 storage oblivious to the behavioral relevance of the memorized information.

494 Introducing distraction into the memory delay unmasked the crucial role of recovery neurons  
495 for working memory maintenance, which would have been hidden in simpler tasks. This  
496 highlights the importance of including richer temporal structure, multiple processing stages  
497 and behavioral perturbation into cognitive task designs to enable dissection of higher-order  
498 brain functions in finer detail and sampling from the full spectrum of underlying mechanisms.

#### 499 **Alternative implementation structures**

500 We focused here on detecting sparse structure in the representational geometry's neuronal  
501 implementation, which is linked to the standardized moment of kurtosis. Consequently, the  
502 loading distributions have both positive and negative heavy tails. Reading out a given sparse  
503 component thus requires both excitatory and inhibitory connections. However, long-range  
504 corticocortical projections are mainly excitatory. This means that other selection criteria that  
505 capture non-symmetrical structure such as the standardized moment of skewness should also  
506 be explored (Koren et al., 2020; Román Rosón et al., 2019).

507 Structure could be in the form of disjointed cell clusters (Hirokawa et al., 2019) or a mixture of  
508 Gaussians (Dubreuil et al., 2020). However, if present, these structures would not dissect the  
509 representational geometry, as they do not have a one-to-one relation to the dimensions in the  
510 activity space. Our neuronal implementation followed a unimodal Laplace distribution (Fig. 2g)  
511 instead of a multimodal distribution.

512 Structure can also be investigated when there are no prior assumptions about the underlying  
513 distributions of neuronal loadings. For example, given that neuronal firing is energy-consuming  
514 and non-negative, possibly encouraging neurons to align to the dimensions of the  
515 representational geometry that have shorter ranges of variation, non-uniform distributions of  
516 the number of selective neurons across different dimensions can arise (Whittington et al.,  
517 2022). However, because all neurons are counted equally, structure probed non-  
518 parametrically could potentially be clouded by the large number of weakly coding (non-  
519 dominant) neurons and thus difficult to detect, in particular in PFC (Bernardi et al., 2020).

#### 520 **Relation of SCA to other linear dimensionality reduction methods**

521 Different linear dimensionality reduction methods based on L2 reconstruction loss will yield  
522 comparable representational geometries, but they will not find the same projections of the  
523 representational geometry, i.e., the same components or the same coordinate system in which  
524 the data is expressed. The principle components of PCA are conveniently orthogonal and  
525 ranked by variance (Vu & Lei, 2013), but usually neither correspond to task-related  
526 components nor align to the activity of individual neurons (Higgins et al., 2021). Truncating the  
527 smaller PCs provides denoised signal as a preprocessing step for independent component  
528 analysis (ICA) that can infer the independent sources in the signal space (Hyvärinen & Oja,  
529 2000). Its most common form, fastICA, enforces sparsity constraints on the activity of the  
530 components, reflecting an assumption about the activity (Hyvarinen, 1999). In contrast, in SCA  
531 the sparsity constraint is on the neuronal implementation, i.e., the potential readout weights  
532 corresponding to the mixing matrix in ICA, reflecting an assumption about the connectivity.

533 Neuronal representations must be communicated. Information that cannot be accessed by  
534 other neurons does not exist. In order to understand complex neural systems such as the PFC  
535 where we lack clear priors about the signal sources, it is paramount to exploit the circuit and  
536 wiring motifs that underlie the observed activity patterns.

537 **References**

- 538 Aoi, M. C., Mante, V., & Pillow, J. W. (2020). Prefrontal cortex exhibits multidimensional  
 539 dynamic encoding during decision-making. *Nature Neuroscience*, 23(11), 1410–1420.  
 540 <https://doi.org/10.1038/s41593-020-0696-5>
- 541 Barack, D. L., & Krakauer, J. W. (2021). Two views on the cognitive brain. *Nature Reviews*  
 542 *Neuroscience*, 22(6), 359–371. <https://doi.org/10.1038/s41583-021-00448-6>
- 543 Bernardi, S., Benna, M. K., Rigotti, M., Munuera, J., Fusi, S., & Salzman, C. D. (2020). The  
 544 Geometry of Abstraction in the Hippocampus and Prefrontal Cortex. *Cell*, 183(4), 954-  
 545 967.e21. <https://doi.org/10.1016/j.cell.2020.09.031>
- 546 Bondanelli, G., & Ostojic, S. (2020). Coding with transient trajectories in recurrent neural  
 547 networks. *PLOS Computational Biology*, 16(2), e1007655.  
 548 <https://doi.org/10.1371/journal.pcbi.1007655>
- 549 Cavanagh, S. E., Towers, J. P., Wallis, J. D., Hunt, L. T., & Kennerley, S. W. (2018).  
 550 Reconciling persistent and dynamic hypotheses of working memory coding in  
 551 prefrontal cortex. *Nature Communications*, 9(1), 3498. <https://doi.org/10.1038/s41467-018-05873-3>
- 553 Chung, S., & Abbott, L. F. (2021). Neural population geometry: An approach for understanding  
 554 biological and artificial neural networks. *Current Opinion in Neurobiology*, 70, 137–144.  
 555 <https://doi.org/10.1016/j.conb.2021.10.010>
- 556 Courchesne, E., Mouton, P. R., Calhoun, M. E., Semendeferi, K., Ahrens-Barbeau, C., Hallet,  
 557 M. J., Barnes, C. C., & Pierce, K. (2011). Neuron Number and Size in Prefrontal Cortex  
 558 of Children With Autism. *JAMA*, 306(18), 2001–2010.  
 559 <https://doi.org/10.1001/jama.2011.1638>
- 560 Dubreuil, A., Valente, A., Beiran, M., Mastrogiuseppe, F., & Ostojic, S. (2020). *Complementary*  
 561 *roles of dimensionality and population structure in neural computations* [Preprint].  
 562 bioRxiv. <https://doi.org/10.1101/2020.07.03.185942>
- 563 Elsayed, G. F., & Cunningham, J. P. (2017). Structure in neural population recordings: An  
 564 expected byproduct of simpler phenomena? *Nature Neuroscience*, 20(9), 1310–1318.  
 565 <https://doi.org/10.1038/nn.4617>
- 566 Eyal, G., Verhoog, M. B., Testa-Silva, G., Deitcher, Y., Benavides-Piccione, R., DeFelipe, J.,  
 567 de Kock, C. P. J., Mansvelder, H. D., & Segev, I. (2018). Human Cortical Pyramidal  
 568 Neurons: From Spines to Spikes via Models. *Frontiers in Cellular Neuroscience*, 12,  
 569 181. <https://doi.org/10.3389/fncel.2018.00181>
- 570 Georgiev, P., Theis, F., Cichocki, A., & Bakardjian, H. (2007). Sparse component analysis: A  
 571 new tool for data mining. *Data Mining in Biomedicine*, 7(Part 1), 91–116.
- 572 Harris, J. A., Mihalas, S., Hirokawa, K. E., Whitesell, J. D., Choi, H., Bernard, A., Bohn, P.,  
 573 Caldejon, S., Casal, L., Cho, A., Feiner, A., Feng, D., Gaudreault, N., Gerfen, C. R.,  
 574 Graddis, N., Groblewski, P. A., Henry, A. M., Ho, A., Howard, R., ... Zeng, H. (2019).  
 575 Hierarchical organization of cortical and thalamic connectivity. *Nature*, 575(7781),  
 576 195–202. <https://doi.org/10.1038/s41586-019-1716-z>
- 577 Herculano-Houzel, S., Catania, K., Manger, P. R., & Kaas, J. H. (2015). Mammalian Brains  
 578 Are Made of These: A Dataset of the Numbers and Densities of Neuronal and  
 579 Nonneuronal Cells in the Brain of Glires, Primates, Scandentia, Eulipotyphlans,  
 580 Afrotherians and Artiodactyls, and Their Relationship with Body Mass. *Brain, Behavior*  
 581 *and Evolution*, 86(3–4), 145–163. <https://doi.org/10.1159/000437413>
- 582 Higgins, I., Chang, L., Langston, V., Hassabis, D., Summerfield, C., Tsao, D., & Botvinick, M.  
 583 (2021). Unsupervised deep learning identifies semantic disentanglement in single

584           inferotemporal face patch neurons. *Nature Communications*, 12(1), 6456.  
585           <https://doi.org/10.1038/s41467-021-26751-5>

586 Hirokawa, J., Vaughan, A., Masset, P., Ott, T., & Kepecs, A. (2019). Frontal cortex neuron  
587 types categorically encode single decision variables. *Nature*, 576(7787), 446–451.  
588           <https://doi.org/10.1038/s41586-019-1816-9>

589 Hyvarinen, A. (1999). *Fast ICA for noisy data using Gaussian moments*. 5, 57–61.

590 Hyvärinen, A., & Oja, E. (2000). Independent component analysis: Algorithms and applications.  
591           *Neural Networks*, 13(4–5), 411–430. [https://doi.org/10.1016/S0893-6080\(00\)00026-5](https://doi.org/10.1016/S0893-6080(00)00026-5)

592 Jacob, S. N., Hähnke, D., & Nieder, A. (2018). Structuring of Abstract Working Memory  
593 Content by Fronto-parietal Synchrony in Primate Cortex. *Neuron*, 99(3), 588-597.e5.  
594           <https://doi.org/10.1016/j.neuron.2018.07.025>

595 Jacob, S. N., & Nieder, A. (2014). Complementary Roles for Primate Frontal and Parietal  
596 Cortex in Guarding Working Memory from Distractor Stimuli. *Neuron*, 83(1), 226–237.  
597           <https://doi.org/10.1016/j.neuron.2014.05.009>

598 Johnstone, I. M., & Lu, A. Y. (2009). On Consistency and Sparsity for Principal Components  
599 Analysis in High Dimensions. *Journal of the American Statistical Association*, 104(486),  
600 682–693. <https://doi.org/10.1198/jasa.2009.0121>

601 Kim, R., & Sejnowski, T. J. (2021). Strong inhibitory signaling underlies stable temporal  
602 dynamics and working memory in spiking neural networks. *Nature Neuroscience*, 24(1),  
603 129–139. <https://doi.org/10.1038/s41593-020-00753-w>

604 Kobak, D., Brendel, W., Constantinidis, C., Feierstein, C. E., Kepecs, A., Mainen, Z. F., Qi, X.-  
605 L., Romo, R., Uchida, N., & Machens, C. K. (2016). Demixed principal component  
606 analysis of neural population data. *ELife*, 5, e10989.  
607           <https://doi.org/10.7554/eLife.10989>

608 Koren, V., Andrei, A. R., Hu, M., Dragoi, V., & Obermayer, K. (2020). Pairwise Synchrony and  
609 Correlations Depend on the Structure of the Population Code in Visual Cortex. *Cell*  
610 *Reports*, 33(6), 108367. <https://doi.org/10.1016/j.celrep.2020.108367>

611 Kornblith, S., Norouzi, M., Lee, H., & Hinton, G. (2019). *Similarity of neural network*  
612 *representations revisited*. 3519–3529.

613 Kriegeskorte, N., & Wei, X.-X. (2021). Neural tuning and representational geometry. *Nature*  
614 *Reviews Neuroscience*, 22(11), 703–718.

615 Ledoit, O., & Wolf, M. (2004). A well-conditioned estimator for large-dimensional covariance  
616 matrices. *Journal of Multivariate Analysis*, 88(2), 365–411.  
617           [https://doi.org/10.1016/S0047-259X\(03\)00096-4](https://doi.org/10.1016/S0047-259X(03)00096-4)

618 Lee, H., Battle, A., Raina, R., & Ng, A. Y. (2007). Efficient sparse coding algorithms. *Advances*  
619 *in Neural Information Processing Systems*, 801–808.

620 Libby, A., & Buschman, T. J. (2021). Rotational dynamics reduce interference between  
621 sensory and memory representations. *Nature Neuroscience*, 24(5), 715–726.  
622           <https://doi.org/10.1038/s41593-021-00821-9>

623 Mante, V., Sussillo, D., Shenoy, K. V., & Newsome, W. T. (2013). Context-dependent  
624 computation by recurrent dynamics in prefrontal cortex. *Nature*, 503(7474), 78–84.  
625           <https://doi.org/10.1038/nature12742>

626 Murray, J. D., Bernacchia, A., Freedman, D. J., Romo, R., Wallis, J. D., Cai, X., Padoa-  
627 Schioppa, C., Pasternak, T., Seo, H., Lee, D., & Wang, X.-J. (2014). A hierarchy of  
628 intrinsic timescales across primate cortex. *Nature Neuroscience*, 17(12), 1661–1663.  
629           <https://doi.org/10.1038/nn.3862>

630 Murray, J. D., Bernacchia, A., Roy, N. A., Constantinidis, C., Romo, R., & Wang, X.-J. (2017).  
631 Stable population coding for working memory coexists with heterogeneous neural



632 dynamics in prefrontal cortex. *Proceedings of the National Academy of Sciences*,  
633 114(2), 394–399. <https://doi.org/10.1073/pnas.1619449114>

634 Nieder, A., Freedman, D. J., & Miller, E. K. (2002). Representation of the Quantity of Visual  
635 Items in the Primate Prefrontal Cortex. *Science*, 297(5587), 1708–1711.  
636 <https://doi.org/10.1126/science.1072493>

637 Okazawa, G., Hatch, C. E., Mancoo, A., Machens, C. K., & Kiani, R. (2021). Representational  
638 geometry of perceptual decisions in the monkey parietal cortex. *Cell*, 184(14), 3748-  
639 3761.e18. <https://doi.org/10.1016/j.cell.2021.05.022>

640 Olshausen, B. A., & Field, D. J. (1996). Emergence of simple-cell receptive field properties by  
641 learning a sparse code for natural images. *Nature*, 381(6583), 607–609.

642 Orhan, A. E., & Pitkow, X. (2020). Improved memory in recurrent neural networks with  
643 sequential non-normal dynamics. *ArXiv:1905.13715 [Cs, Stat]*.  
644 <http://arxiv.org/abs/1905.13715>

645 Parthasarathy, A., Tang, C., Herikstad, R., Cheong, L. F., Yen, S.-C., & Libedinsky, C. (2019).  
646 *Time-Invariant Working Memory Representations in the Presence of Code-Morphing*  
647 *in the Lateral Prefrontal Cortex* [Preprint]. Neuroscience.  
648 <https://doi.org/10.1101/563668>

649 Rigotti, M., Barak, O., Warden, M. R., Wang, X.-J., Daw, N. D., Miller, E. K., & Fusi, S. (2013).  
650 The importance of mixed selectivity in complex cognitive tasks. *Nature*, 497(7451),  
651 585–590. <https://doi.org/10.1038/nature12160>

652 Román Rosón, M., Bauer, Y., Kotkat, A. H., Berens, P., Euler, T., & Busse, L. (2019). Mouse  
653 dLGN Receives Functional Input from a Diverse Population of Retinal Ganglion Cells  
654 with Limited Convergence. *Neuron*, 102(2), 462-476.e8.  
655 <https://doi.org/10.1016/j.neuron.2019.01.040>

656 Russo, A. A., Bittner, S. R., Perkins, S. M., Seely, J. S., London, B. M., Lara, A. H., Miri, A.,  
657 Marshall, N. J., Kohn, A., Jessell, T. M., Abbott, L. F., Cunningham, J. P., & Churchland,  
658 M. M. (2018). Motor Cortex Embeds Muscle-like Commands in an Untangled  
659 Population Response. *Neuron*, 97(4), 953-966.e8.  
660 <https://doi.org/10.1016/j.neuron.2018.01.004>

661 Saxena, S., & Cunningham, J. P. (2019). Towards the neural population doctrine. *Current*  
662 *Opinion in Neurobiology*, 55, 103–111. <https://doi.org/10.1016/j.conb.2019.02.002>

663 Shinomoto, S., Kim, H., Shimokawa, T., Matsuno, N., Funahashi, S., Shima, K., Fujita, I.,  
664 Tamura, H., Doi, T., Kawano, K., Inaba, N., Fukushima, K., Kurkin, S., Kurata, K., Taira,  
665 M., Tsutsui, K.-I., Komatsu, H., Ogawa, T., Koida, K., ... Toyama, K. (2009). Relating  
666 Neuronal Firing Patterns to Functional Differentiation of Cerebral Cortex. *PLoS*  
667 *Computational Biology*, 5(7), e1000433. <https://doi.org/10.1371/journal.pcbi.1000433>

668 Tang, C., Herikstad, R., Parthasarathy, A., Libedinsky, C., & Yen, S.-C. (2020). Minimally  
669 dependent activity subspaces for working memory and motor preparation in the lateral  
670 prefrontal cortex. *Elife*, 9, e58154.

671 Vu, V. Q., & Lei, J. (2013). Minimax sparse principal subspace estimation in high dimensions.  
672 *The Annals of Statistics*, 41(6). <https://doi.org/10.1214/13-AOS1151>

673 Whittington, J. C., Dorrell, W., Ganguli, S., & Behrens, T. E. (2022). Disentangling with  
674 Biological Constraints: A Theory of Functional Cell Types. *ArXiv Preprint*  
675 *ArXiv:2210.01768*.

## 676 **Methods**

677 Two adult male rhesus monkeys (*Macaca mulatta*, 12 and 13 years old) were used for this  
678 study. All experimental procedures were in accordance with the guidelines for animal  
679 experimentation approved by the national authority, the Regierungspräsidium Tübingen. A  
680 detailed description is provided elsewhere (Jacob et al., 2018; Jacob & Nieder, 2014).

## 681 **Surgical procedures**

682 Monkeys were implanted with two right-hemispheric recording chambers centered over the  
683 principal sulcus of the lateral prefrontal cortex (PFC) and the ventral intraparietal area (VIP) in  
684 the fundus of the intraparietal sulcus. This study reports on the PFC data.

## 685 **Task and stimuli**

686 The animals grabbed a bar to initiate a trial and maintained eye fixation (ISCAN, Woburn, MA)  
687 within 1.75° of visual angle of a central white dot. Stimuli were presented on a centrally placed  
688 gray circular background subtending 5.4° of visual angle. Following a 500 ms pre-sample  
689 (pure fixation) period, a 500 ms sample stimulus containing 1 to 4 dots was shown. The  
690 monkeys had to memorize the sample numerosity for 2,500 ms and compare it to the number  
691 of dots (1 to 4) presented in a 1,000 ms test stimulus. Test stimuli were marked by a red ring  
692 surrounding the background circle. If the numerosities matched (50 % of trials), the animals  
693 released the bar (correct Match trial). If the numerosities were different (50 % of trials), the  
694 animals continued to hold the bar until the matching number was presented in the subsequent  
695 image (correct Non-match trial). Match and non-match trials were pseudo-randomly  
696 intermixed. Correct trials were rewarded with a drop of water. In 80 % of trials, a 500 ms  
697 interfering numerosity of equal numerical range was presented between the sample and test  
698 stimulus. The interfering numerosity was independent from either the sample or test  
699 numerosity and therefore not useful for solving the task. In 20 % of trials, a 500 ms gray  
700 background circle without dots was presented instead of an interfering stimulus, i.e., trial  
701 length remained constant (control condition, blank). Trials with and without interfering  
702 numerosities were pseudo-randomly intermixed. Stimulus presentation was balanced: a given  
703 sample was followed by all interfering numerosities with equal frequency, and vice versa.  
704 Throughout the monkeys' training on the distractor task, there was never a condition where a  
705 stimulus appearing at the time of the distractor was task-relevant.

706 Low-level, non-numerical visual features could not systematically influence task performance  
707 (Jacob & Nieder, 2014; Nieder et al., 2002): in half of the trials, dot diameters were selected at  
708 random. In the other half, dot density and total occupied area were equated across stimuli.  
709 CORTEX software (NIMH, Bethesda, MD) was used for experimental control and behavioral  
710 data acquisition. New stimuli were generated before each recording session to ensure that the  
711 animals did not memorize stimulus sequences.

## 712 **Electrophysiology**

713 Up to eight 1 MΩ glass-insulated tungsten electrodes (Alpha Omega, Israel) per chamber and  
714 session were acutely inserted through an intact dura with 1 mm spacing. Single units were  
715 recorded at random; no attempt was made to preselect for particular response properties  
716 (Jacob & Nieder, 2014). Signal amplification, filtering, and digitalization were accomplished  
717 with the MAP system (Plexon, Dallas, TX). Waveform separation was performed offline  
718 (Plexon Offline Sorter).

## 719 **Data analysis**

720 Data analysis was performed with Python using custom scripts based on packages NumPy,  
721 SciPy, sci-kit learn, TensorFlow2, PyTorch, Matplotlib and Plotly.

## 722 **Preprocessing**

723 Single units were included in the analysis if they were recorded in at least 4 correct trials of  
724 each task condition (meaning each unique sample and distractor numerosity combination).  
725 This resulted in 467 neurons across 78 sessions recorded in the PFC. Trials without distractors  
726 were not included in the analyses unless specified otherwise.

727 Unless specified otherwise, the firing rates were binned in a Gaussian window with sigma of  
728 50 ms and step of 100 ms, aligned to the start of the fixation period. The data were then  
729 organized into a neuron-by-condition-by-timepoint tensor. Each tensor entry was normalized  
730 by the standard deviation across trials (within each condition).

## 731 **Demixing**

732 Given the independence of the task variables sample numerosity (s), distractor numerosity (d)  
733 and trial time (t), the neuronal activity can be directly factorized into parts for each variable  
734 and their interaction:

$$735 \quad x = \bar{x} + \bar{x}_t + \bar{x}_s + \bar{x}_d + \bar{x}_{st} + \bar{x}_{dt} + \bar{x}_{sd} + \bar{x}_{sdt}$$

736 Because the stimulus response is also modulated by time, each part was grouped together  
737 with its interaction with time (Kobak et al., 2016):

$$738 \quad x_{time} = \bar{x}_t$$

$$739 \quad x_{sample} = \bar{x}_s + \bar{x}_{st}$$

$$740 \quad x_{distractor} = \bar{x}_d + \bar{x}_{dt}$$

$$741 \quad x_{sd \text{ interaction}} = \bar{x}_{sd} + \bar{x}_{sdt}$$

## 742 **Visualization of representation and implementation space**

743 For a data matrix  $X$  where each column vector  $x$  is the demixed activity of a neuron, the  
744 singular value decomposition was taken:

$$745 \quad X = U\Sigma V^T$$

746 where  $U$  and  $V$  are unitary matrices and  $\Sigma$  is a diagonal matrix with ordered singular values.  
747 The first  $n$  columns of  $U\Sigma$  are the PCs that were used to visualize the representational  
748 geometry. The first  $n$  columns of  $V\Sigma$  are loadings on the PCs that were used to visualize the  
749 implementation space.

750 Within this subspace an arbitrary component can be specified with  $U\Sigma P_{:,1}$  ( $P_{:,1}$  being a column  
751 vector from a unitary matrix  $P$ ), with the orientation of this component given by  $P_{:,1}$ . The  
752 loadings on this component will be the first row of  $(U\Sigma P)^+ X = P^T V^T$ , that is  $P_{:,1}^T V^T$ . This  
753 way, the loadings are visualized with the same orientation  $P_{:,1}$ . in implementation space as  
754 their corresponding component in representation space. The sparsity index of the neuronal  
755 loadings on component  $U\Sigma P_{:,1}$  is then:

756 
$$SI(P_{:,1}) = \text{kurtosis}(P_{:,1}^T V^T) / 3$$

757 
$$\text{kurtosis}(x) = \langle (x - \bar{x})^4 \rangle / \langle (x - \bar{x})^2 \rangle^2$$

758 **Sparse component analysis**

759 Following the formulation of sparse coding (Georgiev et al., 2007; Lee et al., 2007; Olshausen  
760 & Field, 1996), sparse component analysis (SCA) reduces the dimensionality of the dataset  
761 and extracts the unique components by enforcing a sparse penalty on neuronal loadings:

762 
$$\text{Loss} = \left\| X - \sum_{i=1}^k \vec{u}_i \vec{v}_i^T \right\|_{\text{frobienius}} + \alpha \sum_{i=1}^k \|\vec{v}_i\|_1 + \beta \sum_{i=1}^k \|\vec{v}_i\|_2^2$$

763 
$$\|\vec{u}_i\| = 1$$

764 The loss function is defined as the sum of the reconstruction loss and the regularizations. Data  
765  $X$  is organized as a  $n$  firing instances by  $p$  neurons matrix.  $X$  is then approximated by  $k$  firing  
766 activity vectors  $\vec{u}$  and their corresponding neuronal loadings  $\vec{v}$ . Parameter  $\alpha$  controls the  
767 strength of L1-regularization that encourages sparsity of the loadings. Parameters  $\alpha$  and  $k$   
768 were determined by a cross-validated grid search.  $\beta$  was set at 0.01 to smooth the loss  
769 landscape and make the result stable across random initializations.

770 **Substitute data for SCA**

771 Substitute data were created for the demixed sample coding part  $X$  of the data (Fig. 2). For  
772 the singular value decomposition  $X = U\Sigma V^T$ ,  $U\Sigma$  specifies the representational geometry  
773 (see above). Operations were performed on  $V$  only.

774 A random unitary matrix  $R$  with the size of the number of neurons was drawn from a Haar  
775 distribution. The original matrix  $V$  was replaced with  $V' = VR$ .  $V'$  is also a unitary matrix,  
776 meaning that this manipulation will not change the geometries but will rotate them to random  
777 axes. In other words, it will linearly combine the loadings including those on the components  
778 with very low variance, which will render the substitute distribution of loadings on the sample  
779 numerosity components close to Gaussian. The substitute data is then  $X' = U\Sigma V'^T = XR$

780 **Measures of sparse component activity**

781  $\vec{u}_i$  in SCA specifies the activity of the sparse component  $i$ . The following measures of the set  
782 of  $\vec{u}_i$  were compared between the original dataset and its substitutes ( $n = 1000$ ).

783 *Spread of representation.* The standard deviation of  $\vec{u}_i$  across different numerosity conditions  
784  $k$  at each time point was used to define the relative (normalized) information at that time point.  
785 Specifically, each  $\vec{u}_i$  was first reshaped into a condition-by-timepoint matrix  $Y^i$ . Then the  
786 information in component  $i$  at time point  $t$  is given by:

787 
$$Z_{i,t} = \sqrt{\langle (Y_{k,t}^i - \langle Y_{k,t}^i \rangle_k)^2 \rangle_k}$$

788 The skewness of the information across time points was calculated for each component and  
789 averaged across components as follows:

790 
$$Skew_i = \langle (Z_{i,t} - \overline{Z_{i,t}})^3 \rangle_t / \langle (Z_{i,t} - \overline{Z_{i,t}})^2 \rangle_t^{3/2}$$

791 Positively skewed  $Z$  indicates a long tail in the distribution of information across time points,  
 792 corresponding to few time points having high information. Conversely, a smaller or even  
 793 negative skewness implies there are more high information timepoints than low information  
 794 time points, making the high information more spread out across time points. We define the  
 795 spread of representation as the negative skewness:

796 
$$Spread = -\langle Skew_i \rangle_i$$

797 *Overlap of active periods.* The dot product of the information of every pair of components  $i$  and  
 798  $j$  was taken and averaged across pairs:

799 
$$Overlap = \langle Z_{i,t} Z_{j,t}^T \rangle$$

800 *Maximum tuning reversal.* A given component  $i$  may show changes of tuning to sample  
 801 numerosities during the course of a trial. Its tuning at time  $t$  is specified by  $Y_{:,t}^i$ . For each  
 802 component  $i$ , the dot product similarity of tunings between timepoint pairs was specified in the  
 803 non-diagonal entries in  $C^i = Y^{iT} Y^i$ , where the diagonal entries are the strength of the tuning  
 804 at each time point.  $C^i$  was then normalized to the strongest tuning:  $C^{i'} = C^i / \max(C^i)$ . The  
 805 most negative entry in  $C^{i'}$  was then the degree of reversal in this component.  $Reversal_i =$   
 806  $-\min(C^{i'})$ . It would reach the maximum of 1 when tuning at a given time point is the  
 807 complete reversal of the strongest tuning. It would be close to 0 when the tuning does not  
 808 reverse. The maximum tuning reversal is then the largest reversal in a set of SCs:

809 
$$Max\ tuning\ reversal = \max_i Reversal_i = \max_i \left[ -\min \left( \frac{Y^{iT} Y^i}{\max(Y^{iT} Y^i)} \right) \right]$$

810 *Component similarity.* Let  $U_{sca}$  be the concatenation of activity  $\vec{u}_i$  and  $V_{sca}$  the concatenation  
 811 of loadings  $\vec{v}_i$  of the sparse component  $i$ . The data matrix can be expressed as  $X =$   
 812  $U_{sca} V_{sca}^T + \epsilon$ .  $\epsilon$  denotes the noise term. Then it follows  $U_{sca}^+ (X - \epsilon) = V_{sca}^T$ . The  
 813 pseudoinverse  $U_{sca}^+$  can be viewed as a linear transform of the original data. Since all the  
 814 activities  $\vec{u}$  have unit length, larger loadings would be required to express an arbitrary  
 815 geometry when the activities are correlated, meaning lower efficiency. The component  
 816 similarity is measured by the product of the singular values of  $U_{sca}$ . Formally, if the singular  
 817 value decomposition gives  $U_{sca} = U \Sigma V^T$ , then

818 
$$Similarity = \prod_i \Sigma_{i,i}$$

819 The similarity can also be viewed as the determinant of the transformation matrix from arbitrary  
 820 orthogonal bases to the bases of  $U_{sca}$ .

## 821 Numerosity information in different components

822 The standard deviation  $Z_{i,t}$  for all time points  $t$  specifies the evolution of normalized  
 823 information within this component. But since  $\vec{u}_i$  in component  $i$  has unit length, this measure

824 does not allow for direct comparisons between components (see above). To allow for such  
 825 comparisons (Fig. S1), the norm of  $\vec{v}_i$  is therefore applied to  $Z_{i,t}$  as a scaling factor:

$$826 \quad \text{Information} = \|\vec{v}_i\| Z_{i,t}$$

### 827 **Linear discriminant analysis decoding**

828 Neurons recorded in different sessions were stitched together. To account for the different  
 829 number of trials recorded per neuron, a criterion was set to ensure there were at least 1.5  
 830 times more trials than neurons. This resulted in 228 neurons with at least 385 trials each.  
 831 Removing incorrect trials and selecting the minimum number of trials recorded per condition  
 832 and neuron left 118 trials per neuron. Trials of the same condition were then randomly selected  
 833 for each repetition of the analysis.

834 Multi-class linear discriminant analysis (LDA; sci-kit learn package) was used for decoding  
 835 because of its advantageous property of accounting for data covariance. LDA assumes the  
 836 same covariance in every class. It finds the projection that preserves the Mahalanobis  
 837 distance between classes and predicts the label of a new data point by its Mahalanobis  
 838 distance to the class centroid. Shrinkage of the measured covariance matrix was performed  
 839 by averaging with a diagonal matrix. The strength of shrinkage was determined following the  
 840 Ledoit-Wolf lemma (Ledoit & Wolf, 2004).

841 Decoding accuracy, i.e., the ratio of correctly predicted trials, was averaged across 7  
 842 repetitions of 7-fold cross-validation.

### 843 **Spike train statistics**

844 Firing rates were binned in a Gaussian window with sigma of 12.5 ms and step of 25 ms.

845 Correlation, autocorrelation and intrinsic timescales were determined as described elsewhere  
 846 (Murray et al., 2014). The firing rate of each neuron  $n$  at timepoint  $t$  of trial  $i$  is expressed as  
 847  $x_{n,i,t}$ . The Pearson correlation between timepoints  $t1$  and  $t2$  is then:

$$848 \quad r_n(t1, t2) = \frac{\left\langle \left( x_{n,i,t1} - \langle x_{n,i,t1} \rangle_i \right) \left( x_{n,i,t2} - \langle x_{n,i,t2} \rangle_i \right) \right\rangle_i}{\left\langle \left( x_{n,i,t1} - \langle x_{n,i,t1} \rangle_i \right)^2 \right\rangle_i^{1/2} \left\langle \left( x_{n,i,t2} - \langle x_{n,i,t2} \rangle_i \right)^2 \right\rangle_i^{1/2}}$$

849 Autocorrelation is defined as:

$$850 \quad AC_n(\Delta t) = \langle r_n(t0, t0 + \Delta t) \rangle_{t0}$$

851 To account for the refractoriness and adaptation at small time lags, fitting started at the time  
 852 lag where the autocorrelation function had dropped most strongly. Neurons with the strongest  
 853 drop after 400 ms were discarded (6 neurons). The autocorrelation was then fitted with an  
 854 exponential decay:

$$855 \quad AC(\Delta t) = A[\exp(-\Delta t/\tau) + B]$$

856 Parameters  $A$  and  $B$  were constrained in  $[0,1]$  and  $\tau$  was constrained from 10 ms to 2000 ms.  
 857 The autocorrelation function of 8 neurons could not be fitted. The neurons with  $\tau$  fitted below  
 858 20 ms (20 neurons) or above 1600 ms (25 neurons) were excluded because of the biologically  
 859 unrealistic fit. This left 408 neurons. Very few neurons were excluded in the dominant

860 subpopulations (2, 2, and 1 neurons for the sensory, memory and recovery subpopulation,  
861 respectively).

862 The inter-spike intervals (ISI) were determined for the entire session. The coefficient of  
863 variation (CV) measures the global variation of a neuron's ISI and is defined as:

$$864 \quad CV = s. d. (ISI) / \langle ISI \rangle$$

865 In contrast to CV, local variation (LV) measures the local ISI change (Shinomoto et al., 2009).  
866 It is defined as:

$$867 \quad LV = \frac{3}{n-1} \sum_{i=1}^{n-1} (ISI_i - ISI_{i+1})^2 / (ISI_i + ISI_{i+1})^2$$

868 CV and LV are both expected to be 1 for spiking activity following a Poisson process. CV and  
869 LV would be 0 for perfectly regular firing and larger than 1 for more irregular firing than by a  
870 Poisson process.

### 871 **Kullback-Leibler divergence**

872 KL divergence measures the difference between two distributions. For the analyses of intrinsic  
873 time scales and periodicity, KL divergence was calculated between the distribution of statistic  
874  $x$  for the entire population  $P$  and that of sub-samples  $Q$  (either dominant subpopulations or  
875 bootstrap subsamples). It is given by:

$$876 \quad D_{KL}(P \parallel Q) = - \sum_x P(x) \cdot \log Q(x)/P(x)$$

877 To create the null distribution of  $D_{KL}$ , 27 neurons (comparable to the number of neurons in the  
878 dominant subpopulations after exclusion of neurons in which no autocorrelation function could  
879 be fitted) were randomly sampled from the PFC population 1000 times.

### 880 **Temporal dynamics**

881 *Periodicity.* The Fourier transform of the demixed temporal part of the firing rate of each neuron  
882 is given by:

$$883 \quad PSD(f) = DFT(x_{time}(t))$$

884 Then, the periodicity was defined as the ratio between the power of the harmonics of 1/1.5 Hz  
885 (reflecting the onset of visual input at regular spacing of 1.5 s) and the power of all frequencies:

$$886 \quad Periodicity = \sum_{i \in \mathbb{Z}^+} PSD(i \frac{2}{3}) / \sum_f PSD(f)$$

887 *Tangling.* Tangling reflects the smoothness and stability of the flow field around the vicinity of  
888 state  $x_t$  on a trajectory (Russo et al., 2018). It is given by:

$$889 \quad Q(t) = \max_{t'} \frac{\|\dot{x}_t - \dot{x}_{t'}\|^2}{\|x_t - x_{t'}\|^2 + \epsilon}$$

890 It specifies the maximum difference between the derivative at state  $x_t$  and the derivative at  
 891 other states  $x_{t'}$ , normalized by their Euclidean distance. A small constant  $\epsilon$  was added to  
 892 avoid numerical error when the two states were too close.

### 893 Recurrent neural network

894 A recurrent neural network (RNN) model was implemented using the PyTorch neural network  
 895 module. The model has the formulation:

$$896 \quad \mathbf{r}(s, d, t + 1) = \phi(W\mathbf{r}(s, d, t) + I\mathbf{n}(s, d, t) + \mathbf{b})$$

897  $\mathbf{r}$  is the firing rate of units in the condition of sample numerosity  $s$  and distractor numerosity  $d$   
 898 at time point  $t$ .  $\phi$  is the non-linear activation function, chosen to be a rectified linear unit (ReLU)  
 899 to respect the biological characteristics of non-negative firing rates with high upper limits.  $W$   
 900 is the within-population connectivity matrix.  $I$  is the input matrix with the dimensions of 467  
 901 (total number of units) by 4 (number of numerosities). A column  $I_{:,a}$  is the input to the units  
 902 when numerosity  $a$  is being presented.  $\mathbf{n}$  is an indicator vector with the entry  $n_a$   
 903 corresponding to the presented numerosity being 1 and all other entries being 0.  $\mathbf{b}$  is the  
 904 intercept.  $W$ ,  $I$  and  $\mathbf{b}$  are the parameters to be trained. Formally,  $\mathbf{n}$  as a function of trial type  
 905 specified by  $s$  and  $d$  and time point  $t$  is defined by:

$$906 \quad \mathbf{n}(s, d, t) = \mathbf{m}(s) \cdot \text{mask}_{[0.5,1)}(t) + \mathbf{m}(d) \cdot \text{mask}_{[2,2.5)}(t)$$

$$907 \quad \mathbf{m}(x) = [\mathbf{1}_{\{1\}}(x), \mathbf{1}_{\{2\}}(x), \mathbf{1}_{\{3\}}(x), \mathbf{1}_{\{4\}}(x)]^T$$

$$908 \quad \text{mask}_A(t) = \mathbf{1}_A(t * 0.1)$$

$$909 \quad \mathbf{1}_A(x) := \begin{cases} 1, & x \in A \\ 0, & x \notin A \end{cases}$$

910  $\mathbf{m}$  maps a numerosity to the corresponding one-hot vector.  $\text{mask}_A(t)$  indicates the time  
 911 (0.1 s steps) when the corresponding stimulus is presented.  $\mathbf{1}_A(x)$  is an ancillary indicator  
 912 function to define  $\mathbf{m}$  and  $\text{mask}$ .

913 The model was trained to produce the whole sequence of firing rates  $\mathbf{r}(s, d, t)$  in order to  
 914 match the target data  $\mathbf{x}_{s,d,t}$ , given the initial firing rate in the fixation period  $\mathbf{r}(s, d, 0)$  and the  
 915 input  $\mathbf{n}(s, d, t)$ . The loss function is defined as:

$$916 \quad \text{Loss}(W, I, \mathbf{b}) = \sum_{s,d,t} [\mathbf{r}(s, d, t) - \mathbf{x}_{s,d,t}]^2 + \lambda \|W\|_1 + \lambda \|I\|_1$$

$$917 \quad \mathbf{r}(s, d, t_0) = \mathbf{x}_{s,d,t_0}$$

918 The coefficient  $\lambda$  controls the strength of regularization and was determined by a grid search  
 919 with cross validation.

920 The prediction of the later timepoints relies on the quality of the prediction of the early  
 921 timepoints. If the training was done only by giving the first timepoint, convergence would be  
 922 difficult to achieve and learning heavily biased towards reproducing early timepoints in the  
 923 data. To overcome this possible instability, the model was trained in a recursive fashion by  
 924 first using every timepoint as the initial firing rate, training the model to predict the following  
 925 timepoints and gradually increasing the number of timepoints the model needs to predict. As



926 such, at each iteration  $i$ , the temporal sequence  $x_{s,d,t}$  was reorganized into  $T - i$  chunks of  
 927 length  $i + 1$ ,  $\langle x_{s,d,t_0}, \dots, x_{s,d,t_0+i} \rangle$ ,  $t_0 \in \langle 1, \dots, T - i \rangle$ , with the first firing rate in each chunk as  
 928 initial firing rate and the rest as target to be fit by the model.

### 929 Variance explained by RNN

930 The variance explained by the model was determined by the difference between the model's  
 931 predicted trajectory and the trajectory of the original data normalized to the difference between  
 932 a reference trajectory (constant activity set to the first entry of the fixation period) and the  
 933 trajectory of the original data:

$$934 \quad EV = 1 - \frac{\sum_{s,d,t} [r(s, d, t) - x_{s,d,t}]^2}{\sum_{s,d,t} [x_{s,d,t_0} - x_{s,d,t}]^2}$$

935 The normalized EV (Fig. 7c, right axis) was defined as the difference between a substitute's  
 936 EV and the original data's EV, divided by the percentage of the manipulated variance  
 937 (numerosity coding signal, 27.4 %; cp. Fig. 2b). EV for the numerosity signal (Fig. 7d) was  
 938 calculated by replacing both  $r(s, d, t)$  and  $x_{s,d,t}$  with their demixed numerosity representing  
 939 parts.

### 940 Substitute data for RNN

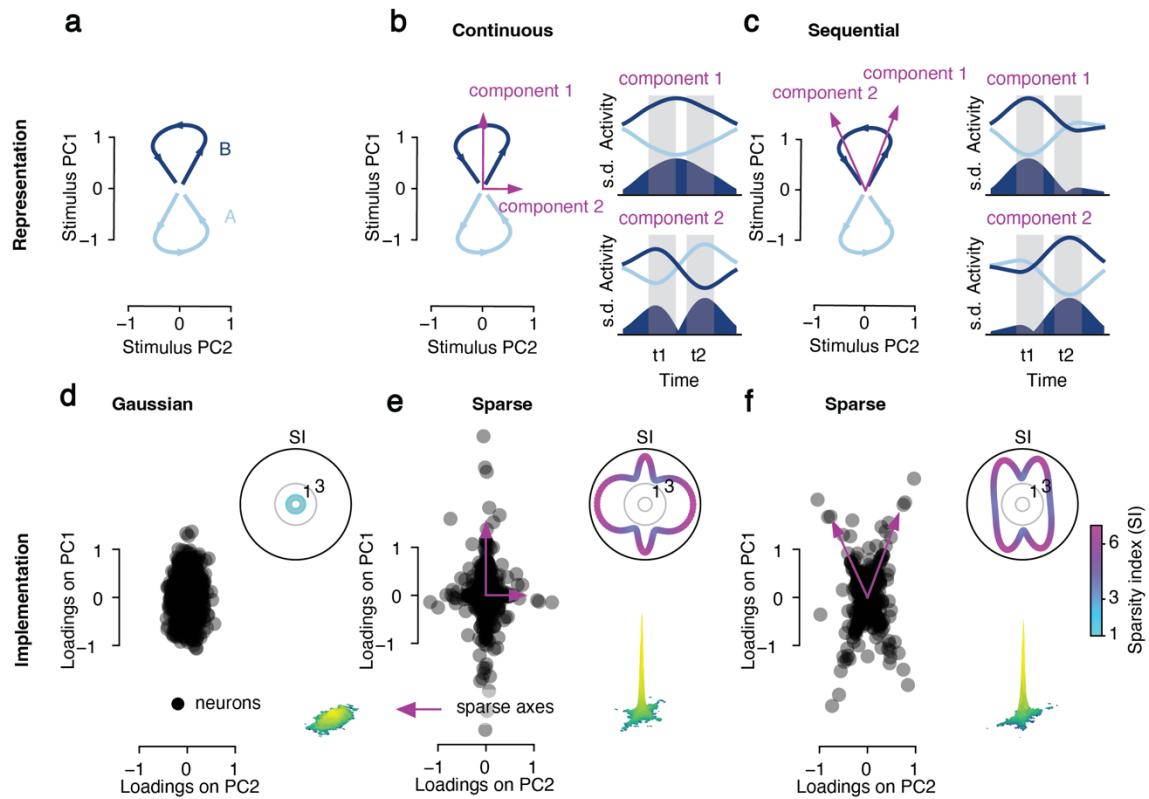
941 In order not to distort the strong connection between sample and distractor numerosity coding  
 942 (e.g., Fig. 3b, Fig. S1), the loadings of these two parts of the data and their interaction were  
 943 shuffled together to create three types of substitute datasets. The RNN model was then trained  
 944 on the substitutes.

945 *Gaussian distribution of loadings.* The Gaussian substitutes were created as described for  
 946 SCA, except for that singular value decomposition was performed on  $X_{sample} + X_{distractor} +$   
 947  $X_{sd\_interaction} = X_{all} - X_t = U\Sigma V^T$ .

948 *Sparse distribution with random alignment.* For  $k$  dimensions of the numerosity coding part of  
 949 the data (determined by cross validation), a  $k \times k$  unitary matrix  $R$  was randomly drawn from  
 950 a Haar distribution and combined with an identity matrix  $I$  to create  $R' = \begin{pmatrix} R & 0 \\ 0 & I \end{pmatrix}$ . Then,  $V' =$   
 951  $VR'$  was substituted for  $V$ . This leaves the sparse structure in the original  $k$  dimensional  
 952 numerosity representing subspace intact, but rotates the sparse structure in  $V_{:,1:k}$  to random  
 953 orientations.

954 *Sparse distribution with original alignment.* The rows of  $V_{:,1:k}$ , i.e., the neuronal identities, were  
 955 permuted by substituting  $V' = (V_{permute,1:k}, V_{:,k+1:p})$  for  $V$ .

**Figure 1**

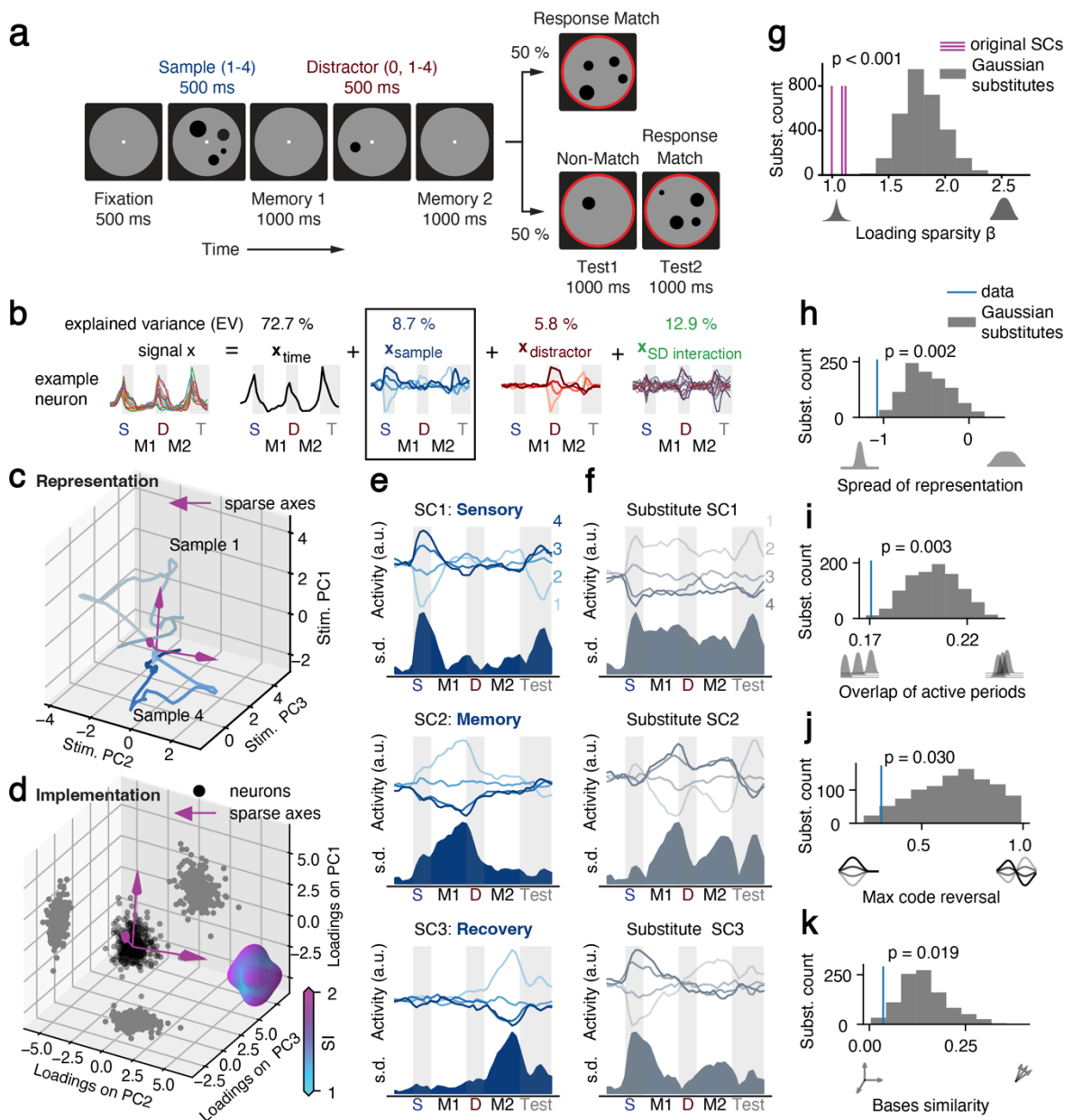


956

957 **Fig. 1 | Different neuronal implementations of the same representational geometry**

958 (a) Representational geometry for two trials with stimuli A and B on the plane specified by  
 959 stimulus PC1 and PC2. Time runs along the individual trajectories. (b) Left: example pair of  
 960 components that express the representational geometry (magenta arrows). Right: activities  
 961 on the corresponding components and standard deviation (s.d.) across components as a  
 962 measure of amount of information carried by them. Components are aligned with the PCs.  
 963 (c) Same layout as in (b) for a non-aligned pair of components. (d-f) Neuronal implementation  
 964 underlying the representational geometry in (a-c), specified by the distribution of neuronal  
 965 loadings on the stimulus PCs. Insets: sparsity index (SI) of all axis orientations in the space  
 966 spanned by PC1 and PC2. Axes with high SI (sparse axes, magenta arrows) in (e) and (f)  
 967 correspond to the components 1 and 2 in (b) and (c), respectively.

**Figure 2**



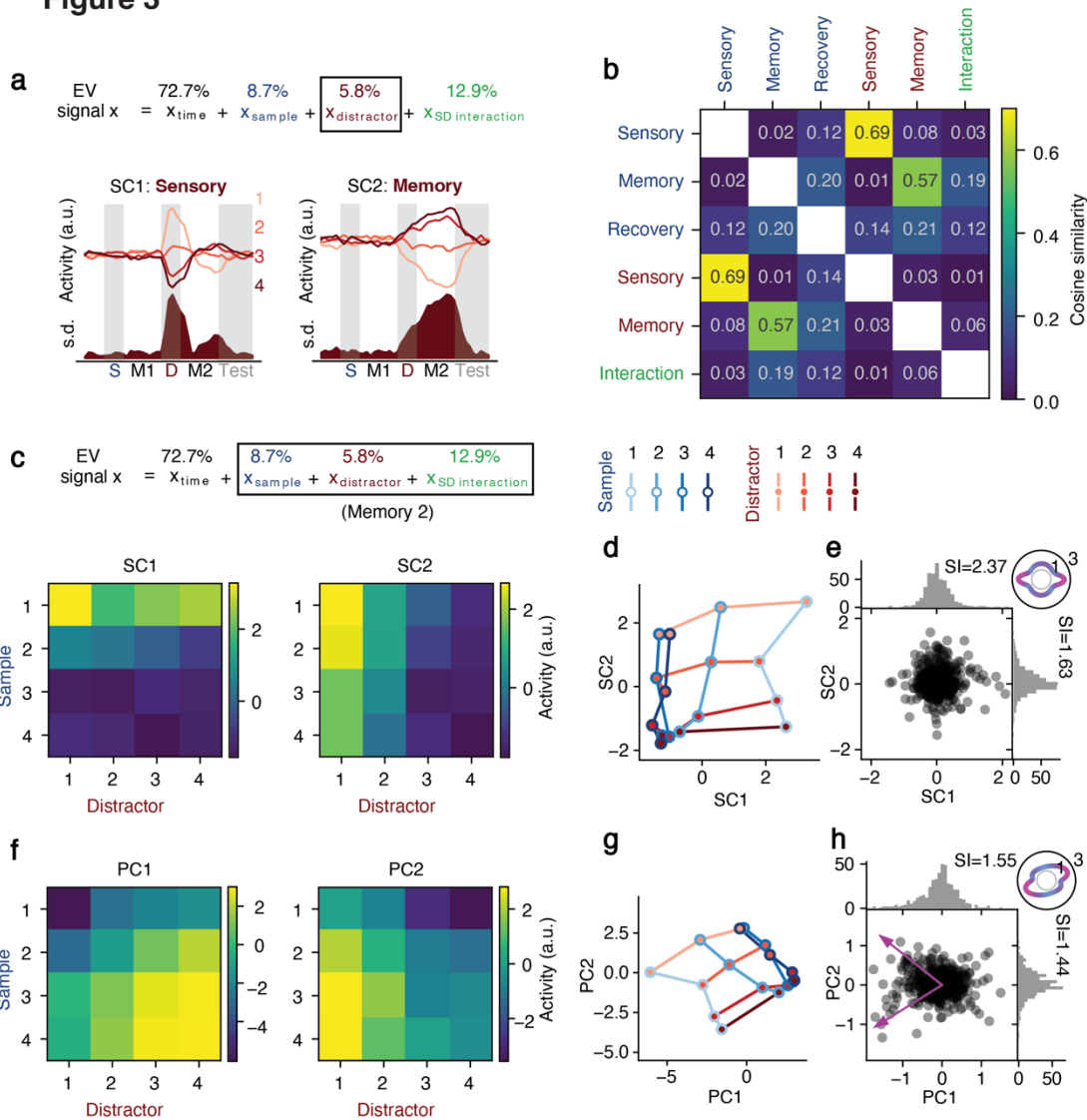
968

969 **Fig. 2 | The neuronal implementation of working memory**

970 (a) Delayed-match-to-numerosity task with distractors. (b) Demixing procedure separating the  
 971 activity of each neuron into the parts coding time, sample numerosity, distractor numerosity  
 972 and sample-distractor interaction. The sample coding part is used for the following analyses.  
 973 Top: percentage of explained variance for each part. (c) Representational geometry for  
 974 sample numerosities 1 and 4 in PC space, averaged across trials of the same condition.  
 975 (d) Loadings of all recorded neurons on the top three PCs (black dots) including distributions  
 976 projected onto the planes formed by PC pairs (gray dots). Sparse axes (magenta arrows;  
 977 determined by SCA) have high SI. Inset: surface plot of SI for all axes in the space. (e) Activity  
 978 of the three identified sparse components (SCs), averaged across trials for each sample  
 979 numerosity condition (top; numbers indicate sample numerosity) and relative information  
 980 across conditions measured as standard deviation (s.d.). (f) SCs of an example substitute  
 981 dataset with non-structured Gaussian implementation. (g) Sparsity  $\beta$  of the neuronal loadings

982 on the SCs (fit to generalized normal distribution) for the original data and the substitute  
983 datasets (permutation test with  $n = 3 \times 1000$  permutations). (**h-k**) Activity measures for the SCs  
984 of the original data and the substitute datasets (permutation test with  $n = 1000$  permutations).

**Figure 3**

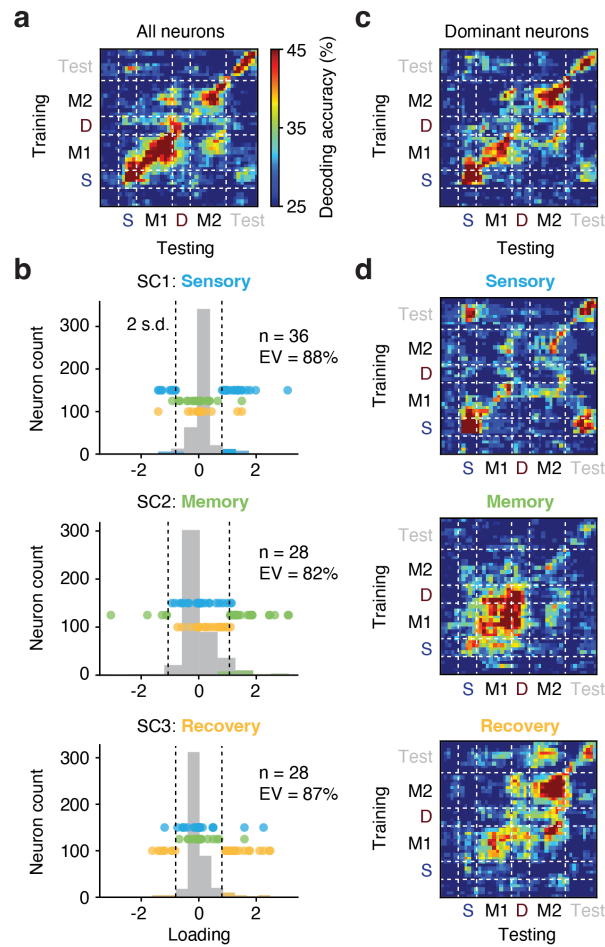


985

986 **Fig. 3 | The effect of distraction on sample representations**

987 (a) Top: the demixed distractor representing part used in the analysis. Bottom: distractor  
 988 numerosity sparse components (SCs). Numbers indicate distractor numerosity. (b) Cosine  
 989 similarity between loadings of sample numerosity SCs (blue), distractor numerosity SCs (red)  
 990 and the sample-distractor interaction SC (green). (c) Activity of the two SCs identified using  
 991 firing rates averaged across the second memory delay for all sample-distractor combinations  
 992 without demixing the stimulus presentations. (d) Representational geometry in SC space. Blue  
 993 and red colors indicate sample and distractor numerosity, respectively. (e) Neuronal loadings  
 994 on the 2 SCs. Dots: joint distribution in SC space. Histograms: marginal distribution of neuronal  
 995 loadings on SC1 and SC2. Inset: SI for all axes. (f-h) Same layout as in (c-e) but for PCs.  
 996 Magenta arrows in (H) indicate sparse axes.

**Figure 4**

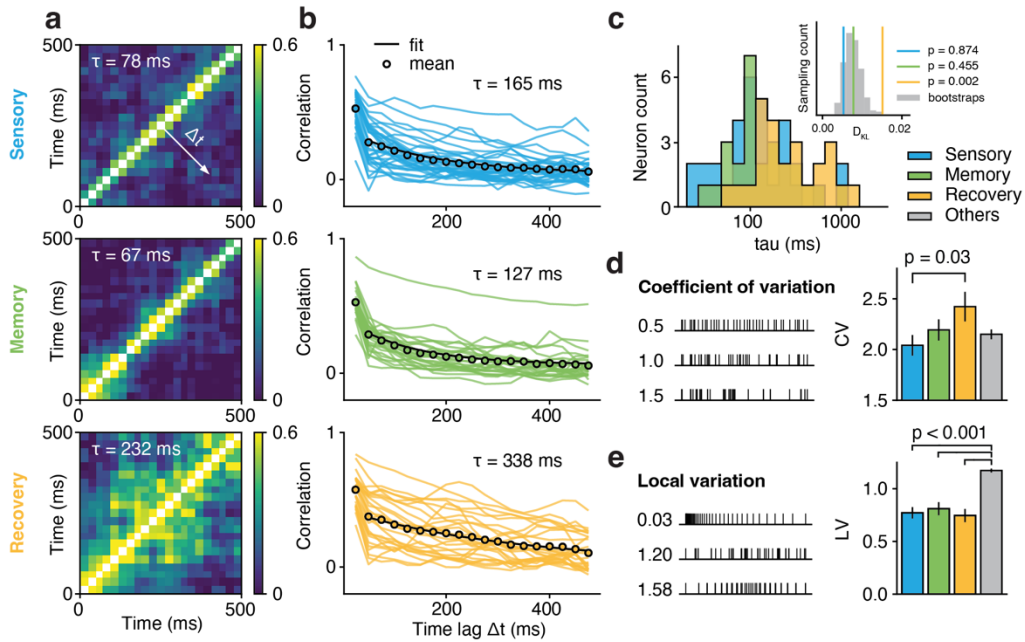


997

998 **Fig. 4 | Subpopulations of neurons dominating working memory coding**

999 (a) Accuracy of cross-temporal linear discriminant analysis (LDA) decoding of sample  
1000 numerosity using all recorded neurons (y axis: training, x axis: testing). (b) Neuronal loadings  
1001 on the three identified sample numerosity SCs. Colored dots indicate the 'dominant' neurons  
1002 selected in each SC (cut-off: two s.d.). The percentage of variance explained within each SC  
1003 is given for each subpopulation. (c) Accuracy of cross-temporal LDA decoding of sample  
1004 numerosity using only the dominant neurons. Compare to (a). (d) Sample numerosity  
1005 decoding accuracy using the dominant subpopulations of each SC. Same color scale in (a),  
1006 (c) and (d).

**Figure 5**

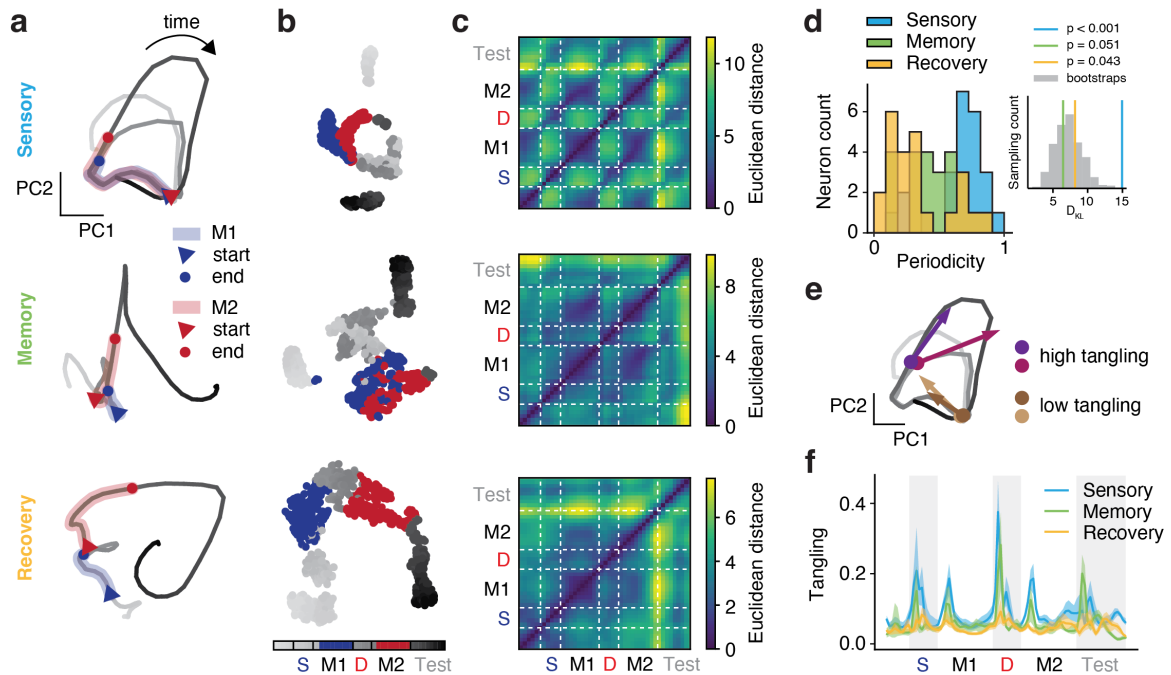


1007

1008 **Fig. 5 | Subpopulation-specific electrophysiological properties**

1009 (a) Between-timepoint Pearson correlations of the trial-to-trial fluctuation of firing rates in the  
 1010 fixation epoch for the three dominant subpopulations. (b) Auto-correlograms obtained by  
 1011 averaging across diagonal offsets in (a). Auto-correlograms of individual neurons are given  
 1012 (single lines) together with the subpopulation average and the fitted exponential decay (black  
 1013 dots and line, respectively). (c) Distribution of fitted decay constants of individual neurons in  
 1014 each dominant subpopulation. Inset: Kullback-Leibler divergence ( $D_{KL}$ ) between the  
 1015 distribution of each subpopulation and the whole population (null distribution for significance  
 1016 testing created with  $n = 1000$  bootstraps from the whole population). (d) Coefficient of  
 1017 variation (CV) of inter-spike intervals (ISI) of the dominant subpopulations and the non-  
 1018 dominant other neurons (two-tailed  $t$ -Test). Left: example spike trains for different CVs.  
 1019 (e) Same layout as in (d) for the local variation (LV) of ISI.

**Figure 6**



1020

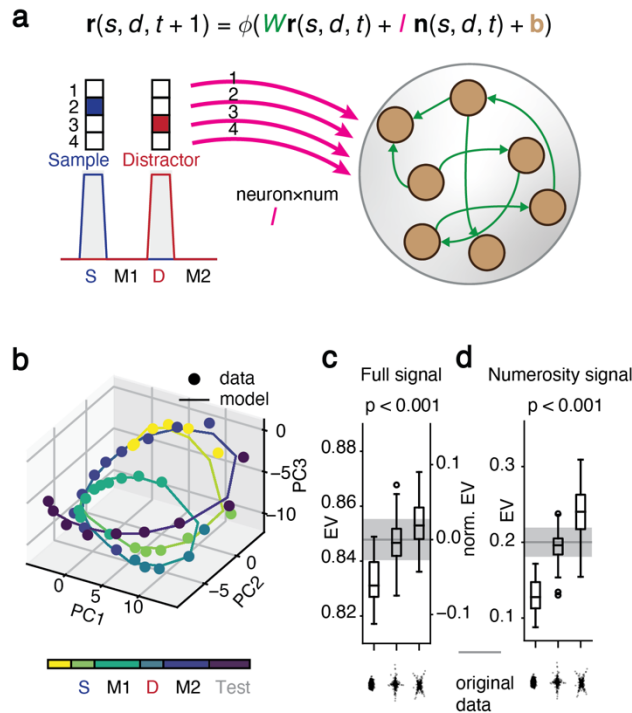
1021

**Fig. 6 | Subpopulation-specific temporal dynamics**

1022 (a) Temporal part of the demixed neuronal activity, averaged across conditions, of each  
 1023 dominant subpopulation projected onto their respective top two PCs. Time runs along the  
 1024 individual trajectories (bin width 50 ms). First and second memory delay are marked in blue  
 1025 and red, respectively. (b) Full signal averaged within each condition and embedded in 2D t-  
 1026 SNE space. Bins as in (a). (c) Euclidean distances between timepoints on the trajectory in (a)  
 1027 of each subpopulation. (d) Distribution of periodicity (relative power of 1/1.5 Hz and harmonics)  
 1028 of individual neurons in each subpopulation. Inset: Kullback-Leibler divergence ( $D_{KL}$ ) between  
 1029 the distribution of each subpopulation and the whole population (null distribution for  
 1030 significance testing created with  $n = 1000$  bootstraps from the whole population). (e) Example  
 1031 timepoints on the trajectory of the sensory subpopulation with high and low tangling. (f) Time  
 1032 resolved tangling of the trajectory of each subpopulation.



**Figure 7**

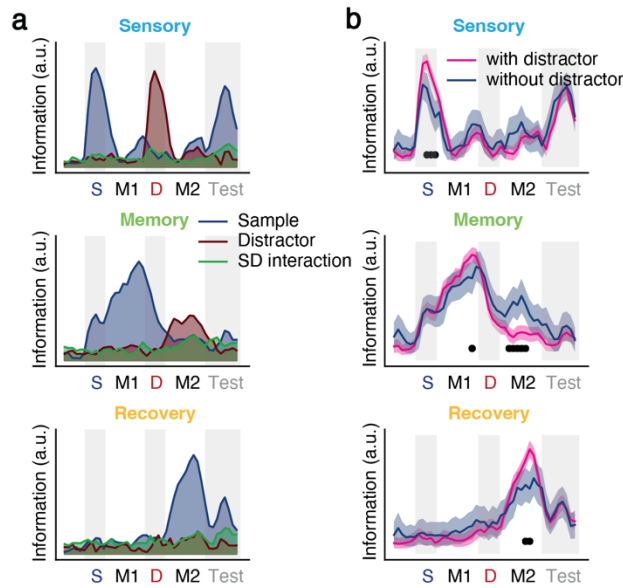


1033

1034 **Fig. 7 | Recurrent neural network modeling**

1035 (a) RNN model governing equation and structure. Magenta and green arrows indicate  
 1036 numerosity-specific inputs and connectivity weights to be trained, respectively. (b) Model fit  
 1037 (solid trajectory) to original data (dots) averaged across all conditions. (c) Percentage of  
 1038 variance of the full signal explained by the model for non-structured Gaussian implementations  
 1039 of numerosity representations (left bar), sparse implementations with random orientations of  
 1040 sparse axes (middle bar) and sparse implementations with the same orientation of sparse  
 1041 axes as in the original data (right bar). Left and right axis show explained variance relative to  
 1042 the full signal and to the manipulated signal, respectively (one-way ANOVA across substitutes).  
 1043 (d) Same layout as in (c) for the percentage of variance of the numerosity signal explained by  
 1044 the model.

Figure S1



1045

1046 **Fig. S1 | The effect of distraction on sample numerosity sparse components**

1047 (a) Information (standard deviation across conditions) about sample numerosity, distractor  
1048 numerosity and their interaction in each of the three sample numerosity sparse components  
1049 (SCs) in trials with a distractor. (b) Sample numerosity information as in (a) for the three SCs  
1050 in trials with and without a distractor. Shaded area indicates [2.5 %, 97.5 %] confidence  
1051 interval. Black dots indicate timepoints with significant differences ( $p < 0.00125$ , bootstrap).

**Figure S2**



1052

1053 **Fig. S2 | Sample-distractor interaction sparse component**

1054 SCA performed on the demixed sample-distractor interaction part of the data identified one  
1055 component that optimally reconstructed the data using cross-validation. The activity of this SC  
1056 is shown for all sample-distractor combinations.