# Context-Dependent Attractor Dynamics in Visual Cortex

Satohiro Tajima [a,b], Kowa Koida [c], Chihiro I. Tajima [d], Hideyuki Suzuki [e], Kazuyuki Aihara [f, g], and Hidehiko Komatsu [g]

a. Department of Basic Neuroscience, University of Geneva. CMU, 1 rue Michel Servet, 1211 Genève, Switzerland.
b. JST PRESTO, Japan Science and Technology Agency, 4-1-8 Honcho, Kawaguchi, Saitama, 332-0012, Japan.
c. EIIRIS, Toyohashi University of Technology. 1-1 Hibarigaoka, Tempaku, Toyohashi Aichi, 441-8580, Japan.
d. Graduate School of Information Science and Technology, the University of Tokyo. 7-3-1 Hongo, Bunkyo, Tokyo 113-8656, Japan.
e. Department of Information and Physical Sciences, Graduate School of Information Science and Technology, Osaka University. 1-5 Yamada-oka, Suita, Osaka 565-0871, Japan.
f. Institute of Industrial Science, the University of Tokyo. 4-6-1 Komaba, Meguro, Tokyo 153-8505, Japan.
g. National Institute for Physiological Sciences. 38 Nishigonaka Myodaiji, Okazaki, Aichi, 444-8585, Japan.

### *Correspondence*

Satohiro Tajima: satohiro.tajima@gmail.com, +41-22-37-94701, Département des Neurosciences fondamentales, University of Geneva. CMU, 1 rue Michel Servet, 1211 Genève, Switzerland.

## 1 Abstract

2 The capacity for flexible sensory-action association in animals has been related to context-dependent

3 attractor dynamics outside the sensory cortices. Here we report a line of evidence that flexibly

4 modulated attractor dynamics during task switching are already present in the higher visual cortex in

5 macaque monkeys. With a nonlinear decoding approach, we can extract the particular aspect of the

6 neural population response that reflects the task-induced emergence of bistable attractor dynamics in a

7 neural population, which could be obscured by standard unsupervised dimensionality reductions such

8 as PCA. The dynamical modulation selectively increases the information relevant to task demands,

9 indicating that such modulation is beneficial for perceptual decisions. A computational model that

10 features nonlinear recurrent interaction among neurons with a task-dependent background input

11 replicates the key properties observed in the experimental data. These results suggest that the context-

12 dependent attractor dynamics involving the sensory cortex can underlie flexible perceptual abilities.

## Introduction

14 Animals are able to adapt their behavior flexibly depending on task contexts, even when the physical

15 stimuli presented to them are identical. The physiological mechanisms underlying this flexible

16 translation of sensory information into behaviorally relevant signals are largely unknown. Recent

17 studies indicate that context-dependent behavior is accounted for by adaptive attractor-like dynamics

18 in the prefrontal areas (Mante et al., 2013; Stokes et al., 2013), which associate sensory representation

19 with behavioral responses depending on task contexts (Freedman et al., 2001, 2002, 2003; Wallis et

20 al., 2001; Wallis and Miller, 2003; Meyers et al., 2012). In contrast to the prefrontal cortex, the visual

21 areas have been suggested to show no or only modest task-related modulations (Sasaki and Uka,

22 2009; McKee et al., 2014). This supports the view that sensory information is processed sequentially

23 across the cortical hierarchy; that is, the physical properties of stimuli are encoded by the sensory

24 cortex, and read out by the higher areas such as the prefrontal cortex.

25 An alternative to this sequential processing model is a view that the sensory cortex is dynamically

26 involved in the neural mechanisms for the flexible sensory-action association. Unlike the former

27 model, the latter does not assume a strong differentiation between sensory and higher areas, which is

28 described in the "encoding-vs.-readout" framework, but allows the decision process to arise from

29 mutual interaction among them. In particular, assuming the involvement of sensory areas in the task-

30 dependent behavior predicts that the neural representations in those areas are modulated by task

31 contexts. Indeed, some studies report that neurons in the sensory areas can change their activities

32 depending on task demands (Koida and Komatsu, 2007; Mirabella et al., 2007; Brouwer and Heeger,

33 2013). For example, it is reported that performing a color categorization task modulates the neural

34 responses to color stimuli in the ventral visual pathway, including macaque inferior temporal (IT)

35 cortex (Koida and Komatsu, 2007) and human V4 and VO1 (Brouwer and Heeger, 2013).

36 However, no clear consensus has been reached on the functional interpretations of those sensory

37 modulations. Some researchers suggest that the task-dependent modulation of neural activities could

38 reflect multiple confounding factors (Sasaki and Uka, 2009). For example, although the task demands

39 can modulate the neuronal response amplitudes in the IT cortex (Koida and Komatsu, 2007), the

40 response amplitudes in individual neurons could be affected by the changes in arousal levels

41 (Greenberg et al., 2008), visual awareness (Lamme et al., 1998; Lamme and Zipser, 2002), task

42 difficulty (Chen et al., 2008) and feature-based attention (Treue and Martínez Trujillo, 1999; Kastner

43 and Ungerleider, 2000; Reynolds and Heeger, 2009).

44 To understand the functions and mechanisms of the task-dependent modulations in the sensory

45 neurons, we need to elucidate the structures of collective dynamics in the neural population—in

3

46    particular, the dynamical structures reflecting the perceptual functions to accomplish the tasks. To this

47    end, in the present study we analyze the spatiotemporal structures of collective neural activity

48    recorded from the macaque IT cortex during context-dependent behavior. To focus on functional

49    aspects of the collective dynamics, we first characterize the evolution of neuronal states within a

50    perceptual space that is reconstructed from the neural population activities. The analysis reveals a

51    task-dependent dynamics of sensory representation in the IT neurons, demonstrating the emergence of

52    discrete attractors during categorical perceptions. Moreover, those attractor dynamics are found to

53    reflect adaptive information processing and explain behavioral variabilities. Finally, through a data

54    analysis and a computational modeling, we suggest a potential mechanism in which the task-

55    dependent attractor structures emerge from a bifurcation in recurrent network dynamics among the

56    sensory and downstream areas.


57    **Results**

58    We analyzed the responses of color-selective neurons recorded in the macaque IT cortex, which

59    change their activities depending on the task demands (Koida and Komatsu, 2007). In the

60    experiments, the monkeys made saccadic responses based on either of two different rules

61    (Categorization or Discrimination) that associate the stimulus with different behavior. In both tasks,

62    the monkeys were presented a sample color stimulus for 500 ms. In the Categorization task, the

63    monkeys then classified the sample color into one of two color categories, "Red" or "Green" (**Figure**

64    **1a**). In the Discrimination task, the monkeys discriminated precise color differences by reporting

65    which of two choice stimuli was the same color as the sample stimulus (**Figure 1b**). We then

66    analyzed the neural responses to the sample colors in the two tasks—where the visual stimuli were

67    physically identical between those tasks.

68    A previous study reported that about 64% of recorded IT cells changed their response magnitudes

69    significantly depending on the task demands (**Figure 1c**) (Koida and Komatsu, 2007). Although the

70    earlier reports have demonstrated that the modulations in individual sensory neurons could be

71    correlated to the hypothetical models that encode categorical information (Koida and Komatsu, 2007;

72    Tajima et al., 2016), the mechanisms and functional impacts of the neural population-response

73    modulation remain to be understood. To elucidate the functional impacts of neural activity

74    modulations, the present study directly investigates the dynamical structure of the collective responses

75    of large numbers of neurons from a decoding perspective.

76 *Reconstructing population activity dynamics from a decoding perspective*

77 To reconstruct the stimulus representation by the neural population, we projected the population

78 activity to stimulus space by extending the idea of likelihood-based decoding (Jazayeri and Movshon,

79 2006; Ma et al., 2006; Brouwer and Heeger, 2009; Graf et al., 2011; Fetsch et al., 2012) such that it

80 captures cross-conditional differences and time-varying properties in neural population

81 representations (**Materials and Methods**). To obtain the joint distribution of neural activities, we

82 generated "pseudo-population" activities from the collection of single neuron firing-rate distributions

83 by randomly resampling the trials (Fetsch et al., 2012).We assumed no noise correlation in our main

84 analyses although we also confirmed by additional analyses that adding noise correlation did not

85 affect the conclusion of this study (see **Discussion**). The basic procedures are as follows (**Figure 1d**):

86 to reconstruct the subjects' percepts, we first built a maximum-likelihood decoder of stimulus based

87 on the spike-count statistics of the correct trials in the Discrimination task, in which the subjects

88 reported precise color identity during stimulus presentation; next, the same decoder was used to

89 analyze the data from the Categorization task. Note that the decoded values are matched to both the

90 presented and the perceived stimuli in the Discrimination task because we used only the correct trials

91 from that task and the monkeys' correct rates were overall high (80-90%). Including the incorrect

92 trials did not affect our conclusion based on the subsequent analyses. On the other hand, in the

93 Categorization task the perceived stimuli could differ from the presented stimuli. Although in the

94 Categorization task we had no access to the precise percepts of the stimulus identities but the

95 categorical reports, we could reconstruct the putative percepts in the decoded stimulus space. The

96 relationship between the decoder outputs and subjective percepts was also supported by follow-up

97 analyses on the choice variability.

98 The decoding-based approach has two major advantages for interpreting the neural population state.

99 First, the decoding provides a way to reduce the dimensionality of neural representation effectively by

100 mapping the high-dimensional population state to a low-dimensional space of the perceived stimuli

101 (which is, in the present case, one-dimensional space of color varying from red to green), which

102 enhances visualization and analysis of the dynamical structures. Second, the decoding-based method

103 enables clear functional interpretation of neural representation because the decoded stimuli are

104 directly related to the subject's judgment of stimulus identity (note that it is often difficult to interpret

105 global distance in a reduced space in nonlinear dimensionality-reduction methods; e.g., (Roweis and

106 Saul, 2000; Tenenbaum et al., 2000; van der Maaten and Hinton, 2008)). In particular, the decoded

107 stimulus identity was what the subject had to respond to in the Discrimination task, and thus we can

108 compare the decoder output and the subjects' behavior (see also **Materials and Methods**). If the

109 decoding is successful, it means that the neural population responses to different stimuli are

110 effectively differentiated within the space of the decoder output. Indeed, cross validation of the

111    decoder performance (by dividing the data from the Discrimination task into two non-overlapping sets

112    of trials) showed a high correct rate (> 75% on average across stimuli), which was comparable to the

113    actual subject performance in the Discrimination task. We will also compare the results to those of

114    other dimensionality reduction techniques in a later subsection.

115    ***Task context modulates the attractor dynamics of the sensory neural population***

116    To characterize the dynamical properties of the decoder output changes for the individual stimuli, we

117    reconstructed the time evolution of the neural states within the space of decoder-output vs. mean

118    firing rate (**Figure 2a**). The population state trajectories during the Discrimination task were

119    accurately matched the presented stimuli, confirming the successful mapping from neural

120    representations to stimuli (**Figure 2a**, bottom). Remember that here we used the trials in which the

121    subjects correctly identified the sample stimuli in the subsequent fine discrimination in the

122    Discrimination task (**Figure 1b**), thus the decoded stimulus identity should also correspond to the

123    stimuli perceived by the subjects. In contrast to the Discrimination task, we found that the same

124    analysis for the Categorization task yielded strikingly different state trajectories (**Figure 2a**, top),

125    which suggests that the neural representation was altered between the two tasks. In particular, the

126    population state trajectories in the Categorization task showed attractor-like dynamics in which the

127    state relaxes toward either of two stable points respectively corresponding to the "Red" and "Green"

128    categories along the "line" structure (in the horizontal direction in the figure) connecting those two

129    stable points (**Figure 2a**, top). The relationship between the mean firing rate and the decoded stimulus

130    identity was also kept in the Discrimination task and showed a similar "line" structure with little

131    bistability (**Figure 2a**, bottom). Interestingly, green stimuli tended to evoke larger neural responses

132    than the red stimuli, consistently in both Discrimination and Categorization tasks, although the reason

133    for this is not clear. Finally, these properties of the dynamics were robust to various changes in the

134    decoder construction and neural noise-correlation structures in the data, indicating that the present

135    results do not rely on the specific designs of the decoder (see **Discussion**). We observed that the

136    results in the eye-fixation task were similar to those of the Categorization task (data not shown),

137    replicating the previous report that the neural tunings in the eye-fixation task shared properties with

138    the Categorization task (Koida and Komatsu, 2007).

139    Remarkably, the attraction toward stable points continued throughout the stimulus presentation

140    period, even after the population average firing rate had stabilized (as demonstrated by the horizontal

141    shifts in **Figure 2a**, top). This also confirms that the dynamics in decoded stimuli are not merely

142    reflecting the changes in the overall firing rate in the population (which could be potentially

143    concerned to affect the decoding analysis through the changes in signal-to-noise ratio in the data). The

144    polarity of the modulation depended strongly on the presented stimulus identity (**Figure 2b**). The

145    modulation was not large at the beginning of the stimulus presentation (light plot along diagonal in

146    **Figure 2b**), but was magnified in the late period (dark plots in "S" shape, **Figure 2b**). The evolution

147    of the modulation continued across the entire period of stimulus presentation, and was not directly

148    associated with the dynamics of the mean firing rate, which became stable about 250 ms after the

149    stimulus onset (**Figure 2c**).


150    *The recurrent model explains the stimulus-dependent dynamics*

151    Standard models of a recurrent dynamical system in which the system's energy function relaxes as the

152    state evolves toward either stable point, naturally accounted for the dynamics converging to stable

153    point attractors in the Categorization task. In addition, the dependency on presented stimulus identity

154    indicates that the modulation was dynamically driven by the visual input, rather than by *pre-readout*

155    (i.e., stimulus-invariant) modulation of neural response gains, such as conventional feature-based

156    attention (Treue and Martínez Trujillo, 1999). These facts are more consistent with the recurrent

157    model than conventional gain-modulation models as an explanation of the population dynamics

158    reported here. To verify this, we next examined how gain-modulation and recurrent models could

159    account for the quantitative aspects of modulation dynamics.


160    To analyze the dynamics of neural modulation quantitatively, we considered three gain-modulation

161    models (in which neural response gains could depend on the task and either of time and stimulus;

162    **Figure 2d**) and a recurrent model (response modulation via self-feedback through mutual connections

163    to two hidden units, whose weights depended on the task but neither on time nor on the stimulus

164    identity; **Figure 2d**). Note that we did not assume explicit stimulus-dependency of model parameters

165    in any of the three models. We derived the model parameters based on the recorded neural responses,

166    such that the modulated neural responses in the Discrimination task fit the responses in the

167    Categorization task (full details of the modeling are provided in the **Materials and Methods**). Using

168    these four models, we determined to what extent the gain modulations and recurrent modulation

169    predict the temporal evolution of decoder output changes in the Categorization task. The model-fitting

170    performances were assessed using cross-validation based on two separate sets of trials: the first set

171    was used to train models, and the second was used to test each model's fitting performance. We

172    computed the cross-validation errors, $E_{CV}$, directly based on the difference between the predicted and

173    actual neural population activities, thus the measure is independent of the assumptions about the

174    decoder (**Materials and Methods**).


175    We found that the recurrent model showed the smallest cross-validation error among the four models

176    ($E_{CV} = 2.78, 2.86, 3.77,$ and $2.08$ in the gain-modulation models 1–3 and the recurrent model,

177    respectively; the recurrent model's error was significantly below each of the gain-modulation models,

178    p<0.001, permutation test). Indeed, neither gain-modulation model could account for the large

179    increase in decoder output change in the late period (about >150 ms) after the stimulus onset (**Figure**

180    **2d**, the green and blue curves). The time- or stimulus-dependency of the gain parameters did not make

181    a major difference in the prediction performance among the three gain-modulation models, suggesting

182    that the modulation depends both on stimulus and time. On the other hand, the recurrent model

183    explained the large continuous increase in decoder output change (**Figure 2d**, the magenta curve). It

184    should be noted that the parameters in the recurrent model were constant across time, in contrast to

185    the time-variant gain-modulation model. This means that the time-invariant recurrent model is

186    superior even to the time-variant gain-modulation model at explaining the task-dependent modulation

187    of neural population dynamics. The reason for this is that the effective modulation signals in the

188    recurrent model could vary across different stimuli because the recurrent architecture allowed the

189    modulation to depend on the neurons' past activities evoked by stimulus, leading to an "implicit"

190    dependency on stimulus and time. It is remarkable that the recurrent model is capable of describing

191    the dynamic activity modulations without assuming any explicit parameter change across stimuli and

192    time, even better than the time- and stimulus-dependent gain modulation, which had much more

193    parameters than the recurrent model. The results were similar when we assumed full-connected

194    pairwise interactions instead of the restricted connection via the hidden units. All the results were

195    cross-validated, making it unlikely that the difference in model performance was caused by

196    overfitting. In addition, the superiority of the recurrent model was robustly observed with changes in

197    the decoder construction and neural noise correlations (**Discussion**). These results support the idea

198    that the task-dependency of neural dynamics originates from a recurrent mechanism, although we do

199    not exclude the possibility of more complex gain-modulation mechanisms (that depend on both the

200    stimulus and time) as substrates of the context-dependent dynamics observed here (see also

201    **Discussion**). Note that the analysis here compares the data-fitting performance of gain-modulation

202    and recurrent models in the decoded stimulus domain, but does not aim to model the mechanisms of

203    the emergent bistable attractor structure in the Categorical task. A possible mechanism underlying the

204    organization and task-dependent modulation of attractor dynamics is discussed later.


205    *Reconstructed collective dynamics explains choice variability*

206    We also found that the neural state represented in the space of the decoded stimulus was closely

207    related to the subjects' subsequent behavior. First, the locus of the behavioral classification boundary

208    in the Categorization task, which moderately prefers the "Green" category, was replicated by the

209    stimulus classification based on decoder output (**Figure 3a, b**). This suggests the decoded-stimulus

210    space used here was closely related to the behavioral response dimension. Second, the modulation of

8

211   the dynamics reflected the subjects' trial-to-trial response variability. The subject's choices between

212   the "Red" and "Green" categories were variable across trials, particularly for the stimuli around the

213   classification boundary (stimuli #4–6), even when the task condition and the presented stimulus were

214   the same. To investigate the mechanism underlying this behavioral variability, we reanalyzed the

215   neural responses during the Categorization task using the same decoding protocol used in the previous

216   sections, but now separated the trials into two groups according to the subsequent choice behavior.

217   We found that the behavioral fluctuation was clearly reflected in the preceding population dynamics

218   in the decoded-stimulus space (**Figure 3c**). The neural state was shifted toward the "Red" extreme

219   before the subject classified the stimulus into the "Red" category, whereas the state was shifted

220   toward "Green" before classifying it into the "Green" category. The difference was small in the

221   beginning of the response, but gradually increased as time elapsed (**Figure 3d**, **e**). Gradual

222   amplification of small differences in the initial state is a general property of a recurrent dynamical

223   system having two distinct stable attractors, which further supports the recurrent model. Note that the

224   current decoding analysis shares some concept with conventional choice-probability analysis in single

225   neurons, but the current decoder analysis focuses more on the collective representation by neural

226   population. In addition, the decoding analysis allows us to specify not only choice polarities but also

227   the estimated perceptual contents (color identities) at each moment.

### *Dynamical modulation enhances task-relevant information*

229   The evidence so far indicates that the neural population in the IT cortex flexibly modulates its

230   recurrent dynamics depending on the task context. What is this modulation for? We hypothesized that

231   the modulation is a consequence of stimulus information processing adapted to the changing task

232   demands. To test this possibility, we computed the mutual information between the neural population

233   firing and the stimulus identity (hue) or stimulus category. The mutual information provides an upper

234   limit for the information extracted from the neural state trajectories, which indicates how the

235   dynamical modulations could contribute to the task-relevant information processing. We found that

236   the modulatory effect was accompanied by selective increases in the task-relevant stimulus

237   information conveyed by the neural population (**Figure 4**). Namely, category information increased in

238   the categorization task compared with the discrimination task, whereas hue information increased in

239   the discrimination task. The fact that the modulation of the neural dynamics increases the task-

240   relevant information indicates that the modulation benefits the subjects switching the tasks depending

241   on the context.

9

242    *Comparison to other methods of dimensionality reduction*

243    We have shown that the decoding approach captures the task-dependent attractor-like dynamics in the

244    neural population. To examine how the other dimensionality reduction methods capture the task-

245    dependent natures of the collective neural dynamics, we first applied the principal component analysis

246    (PCA) to the neural responses during the stimulus presentation. **Figure 5a** shows the reconstructed

247    trajectories of the neural population states in the space spanned by PCs 1−3. The trajectories for

248    categorization and discrimination tasks largely overlapped, and the task-dependent attractor-like

249    structure is not obvious in this space despite that these top three PCs together explained more than

250    60% of the total variance (**Figure 5b**). This indicates that the task-dependent components of the

251    dynamics are hidden in the other dimensions. Similarly, it was not straightforward to demonstrate the

252    emergence of two discrete attractors in the Categorization task with other unsupervised

253    dimensionality reduction methods, including PCA based on the differential responses between the

254    tasks (**Figure 5c**) and nonlinear methods such as t-stochastic neighbor embedding (tSNE) (van der

255    Maaten and Hinton, 2008) (**Figure 5d**). These results implicate that the task-dependent components

256    could be obscured when visualized naively with some of those conventional methods.

257    *Bifurcation of attractor dynamics in a recurrent model*

258    The analyses in the previous sections have indicated the flexible recurrent interactions that modulate

259    the structures of attractors depending on the task context. What mechanism could explain such a

260    dynamic changes in neural dynamics? Here we show a simple potential mechanism that accounts for

261    the flexible changes in attractor structures in the collective neural dynamics.

262    We extended a model of prefrontal attractor dynamics that was proposed in the context of two-interval

263    discrimination (Machens et al., 2005) by introducing a recurrent interaction that involves a population

264    of hue-selective neurons. **Figure 6a** illustrates a potential mechanism for the context-dependent

265    change in attractor structure. We assume that the hue-selective neurons (hereafter referred to as "hue-

266    neurons") in the IT cortex have mutual interaction with category-selective neurons (hereafter,

267    "category-neurons") in the frontal or other cortical area. The hue neurons receive sensory input from

268    earlier visual areas. The connectivity weights between hue- and category-neurons are modeled using

269    the functions of the preferred hues in hue-neurons such that a "red" category-neuron exhibits

270    excitatory interactions with hue-neurons preferring reddish hues and inhibitory interactions with

271    neurons preferring greenish hues (similar for "green" category-neuron). We assume that the category

272    neurons also receive a common background input, and respond based on an activation function with

273    response threshold and saturating nonlinearity, which characterizes the categorical response in cortical

274    neurons (Freedman et al., 2001) (see **Materials and Methods**).

275    This system has different numbers of stable attractors depending on the strength of common

276    inhibitory background input (parameter $B$), with the connectivity among neurons unaffected (**Figure**

277    **6b**–**d**). The neural state converges to a single stable equilibrium point under a strong background

278    inhibition (**Figure 6b**) whereas two distinct stable equilibrium points emerges under a weak or no

279    background input, yielding bistability that depends on the initial state (**Figure 6c**).

280    We confirmed that the model replicated multiple aspects of the collective neural dynamics observed

281    in IT cortex.  First, the representation of modeled hue neurons (hypothetical IT neurons) showed the

282    gradually evolving biases toward either of two extreme stimuli ("red" or "green"; (**Figure 6e**) as well

283    as the moderately higher mean activity in the Categorization task (**Figure 6h**). Second, the recurrent

284    dynamics replicated the gradual development of the choice-related neural variability (**Figure 6f, g**).

285    Third, the circuit enhanced the task relevant information (**Figure 6i**). Finally, the task-dependent

286    components of dynamics could be obscured when visualized with PCA (**Figure 6j**), which is also

287    consistent with the results in IT neurons (**Figure 5a**).

## Discussion

289    We demonstrated that the task context modulates the structures of collective neural dynamics in the

290    macaque IT cortex. The neural population in the IT cortex exhibited the dynamics with two discrete

291    attractors that respectively corresponded to the two task-relevant color categories in the

292    Categorization task. The trial-to-trial variability in the dynamics confirmed that those two stable

293    attractors co-existed under a single stimulus, thus the observed bistability reflects an inherent property

294    of neural circuit. Remarkably, we found that the patterns of the neural state evolution was explained

295    by a recurrent mechanism, but not fully accounted for by conventional gain-modulation models such

296    as the ones assumed for top-down attention (Treue and Martínez Trujillo, 1999; Reynolds and Heeger,

297    2009). The present hierarchical recurrent model rather shares some features with other recent models

298    including the recurrent interactions between top-down and bottom-up signals (Wimmer et al., 2015;

299    Haefner et al., 2016). A unique point in the present model is that it explains the context-dependent

300    structure of collective neural dynamics in terms of the bifurcation of attractors caused by a simple

301    change in the background input to the categorical neurons. Lastly, although the present results suggest

302    a profound contribution by a recurrent mechanism to the context-dependent modulation of sensory

303    cortex dynamics, which has not been emphasized in previous studies, we do not exclude the potential

304    contributions of a gain-modulation mechanism; rather, it is quite possible that the brain uses a

305    combination of both the recurrent and feedforward mechanisms.

306 Recent studies emphasize a variety of stimulus-dependent contextual modulations, particularly in the

307 early visual cortex (Toth et al., 1996; Sceniak et al., 1999, 2002; Sadakane et al., 2006; Tajima et al.,

308 2010; Solomon and Kohn, 2014; Coen-Cagli et al., 2015). However, it is yet to be elucidated whether

309 the same mechanisms also apply to the context-dependent categorical processing in IT cortex as

310 studied here, and how such a modulation could be implemented in biological systems without any

311 recurrent mechanisms. Note that, in principle, a stimulus-dependent gain modulation requires a form

312 of self-referencing process (which is naturally implemented by recurrent mechanisms) because it

313 implies the stimulus encoding being modulated by the encoded stimulus itself, whether the source of

314 modulation is the fluctuations in choice-related activity (Nienborg and Cumming, 2009) or attention

315 (Ecker et al., 2016). Nonetheless, the mathematically equivalent effects could be achieved by a

316 feedforward mechanism in physiological circuits that feature an information duplication (e.g., two

317 parallel feedforward pathways converging at IT cortex, in which one has longer latency than others).

318 We do not exclude this possibility. Our current results demonstrate that the task-dependent neural

319 dynamics were at least not fully accounted for by conventional forms of stimulus-invariant gain

320 modulations such as assumed in a previous study.

321 As a key methodology, we took a decoding approach to reconstruct the perceptual space form neural

322 population activity. One may concern a possibility that the results rely on the selection of decoder. To

323 examine this point, we replicated the same analyses with different decoders, and confirmed that the

324 results reported in this paper were robust to various changes in the decoder construction, such as

325 introducing noise correlations in neural responses, removing the half of cells to use, assuming non-

326 Gaussian models, and ignoring the time dependence (as summarized in **Figure 7**). This suggests that

327 the present results do not require fine tunings of the decoder constructions or assuming the

328 independent noises across neurons. On the other hand, the task dependence of attractor structures

329 could be unclear when visualized with-conventional unsupervised dimensionality reduction methods,

330 despite that PCA could extract cluster structures in a previous human neuroimaging with a color

331 naming task (Brouwer and Heeger, 2013). The effectiveness of the decoding approach shares some

332 aspects with other recent labeled dimensionality-reduction approaches applied to neural population

333 data (Brendel and Machens, 2011; Mante et al., 2013; Okazawa et al., 2015; Kobak et al., 2016).

334 Although it is beyond the scope of the current study to compare all the possible dimensionality

335 reduction methods, we suggest that analyzing neural-population state-space from a decoding

336 perspective could be useful to extract the hidden dynamical properties that are relevant to the

337 functions of collective neural responses.

338 Previous studies have proposed that context-dependent decision making is achieved through flexible

339 modulations of recurrent attractor dynamics within the prefrontal cortex (Mante et al., 2013; Stokes et

12

340    al., 2013). The present results imply that the dynamical mechanisms of context-dependent

341    computation can include not only the prefrontal areas but also the sensory cortex, potentially

342    organizing the distinct representational layers such as hypothesized in the present model (**Figure 6**).

343    Although an earlier study reported attractor-like dynamics in the IT cortex during object

344    categorization (Akrami et al., 2009), the flexible modulation of a dynamical structure depending on

345    task context has not been demonstrated. It should be noted that the present task design differs from

346    those of many other task-switching studies: in contrast to the previous studies, in which the subjects

347    switched behavioral rules between two different categorization tasks (e.g., categorizing motion or

348    color/depth) (Sasaki and Uka, 2009; Mante et al., 2013; Siegel et al., 2015), the present study is based

349    on switching between Categorization and Discrimination. This difference in task design may underlie

350    the apparent discrepancy between the present and the previous studies regarding the involvement of

351    sensory cortex in task switching.

352    The way of neural modulation such that the population response becomes more sensitive to color

353    around the categorical boundary in the Categorization task is consistent with previous human

354    psychophysics showing that the stimulus discriminability is higher around category boundaries

355    (Uchikawa and Sugiyama, 1993; Uchikawa and Shinoda, 1996). Moreover, the present results add a

356    dynamical viewpoint in neural population representations, which predicts that the perceptual illusion

357    depends on time as well as task demands. Beyond color perception, this modulation of dynamics in

358    sensory representation implies a potential physiological substrates of task-dependent perceptual

359    illusion. For instance, perceived motion direction is biased away from the classification boundary

360    during a motion categorization task (Jazayeri and Movshon, 2007). Theoretically, this illusion could

361    be explained both by considering direct modulation of sensory representation (Jazayeri and Movshon,

362    2007), and by assuming a readout mechanism without direct modulation of the sensory neural

363    representation itself (Stocker and Simoncelli, 2008). The first model would be preferred if the motion

364    perception is based on a population coding mechanism similar to the one demonstrated in this study

365    which suggests the neural population representation is indeed modulated at the level of the sensory

366    cortex.

367    The involvement of the sensory cortex in decision-related neural dynamics is consistent with the idea

368    that responses within the sensory cortex are not only read out by the higher areas in a feedforward

369    manner but also affected by decision-related signals through feedback connections from areas outside

370    the sensory cortex (Nienborg and Cumming, 2009; Siegel et al., 2015; Wimmer et al., 2015).

371    Unfortunately, we cannot fully conclude from the present data whether the observed choice-related

372    attractor-dynamics are the *cause* or the *effect* of decision-making (Nienborg and Cumming, 2009).

373    Nonetheless, the fact that modulation of the neural dynamics enhances the task-relevant information

13

374  in sensory neurons may hint at the potential contribution of this modulation to the task performances.

375  In addition, our data suggest that the choice-related difference in the dynamics had already begun

376  during the early period (< 250 ms; **Figure 2**), which is thought to affect the decision (Nienborg and

377  Cumming, 2009). Therefore, it is likely that the task-dependent modulation of neural dynamics (at

378  least during the early period after the stimulus onset) contributed to improving the behavioral

379  performance rather than merely reflected the decision signal. More generally, theoretical studies have

380  proposed that a common recurrent neural circuit can serve as the basis for multiple functions, such as

381  sensory information encoding, categorization and decision (Wang, 2002, 2008; Machens et al., 2005;

382  Furman and Wang, 2008), enabling a flexible use of the neural dynamics depending on context. The

383  present findings suggest involvement of sensory cortex in the context-dependent behavior, leading to

384  a new view that the sensory neurons could contribute to context-dependent behavior by flexibly

385  modulating their collective attractor dynamics.

## Materials and Methods

### *Subjects, stimuli and behavioral task*

388  To study the neural basis of context-dependent behavior, we analyzed neural responses from the

389  anterior inferior temporal (IT) cortices in two female monkeys (*Macaca fuscata*) performing visual

390  tasks. Details of the experimental procedures have been previously published (Koida and Komatsu,

391  2007). All procedures for animal care and experimentation were in accordance with the National

392  Institutes of Health Guide for the Care and Use of Laboratory Animals and were approved by

393  Institutional Animal Experimentation Committee. The monkeys were trained in Categorization and

394  Discrimination tasks. In both tasks, the same 11 sample colors were used as visual stimuli. The 11

395  sample colors ranged from red [color 1, (x = 0.631, y = 0.343 in the CIE 1931 xy chromaticity

396  diagram)] to green [color 11, (x = 0.286, y = 0.603)] and were spaced at equal intervals on the CIE

397  1931 xy chromaticity diagram. The colors all had the same luminance (30 cd/m$^2$). Tasks were

398  alternated in blocks in a fixed sequence that included Categorization and Discrimination tasks, as well

399  as an eye-fixation task in which the monkey passively viewed the same color stimuli. There was no

400  explicit cue to indicate the ongoing task. Each block consisted of 88 correct trials—eight repetitions

401  of the 11 sample color stimuli. The 11 sample colors were presented in a pseudorandom order. If a

402  monkey made an incorrect response to a given color, the trial using that color was repeated after some

403  intervening trials. These repeated trials and other incomplete trials, such as those with fixation errors,

404  were excluded from the subsequent data analyses. The stimulus was usually a disk with a diameter

405  spanning 2.0° of visual angle, but for cells with shape selectivity, an optimal shape was chosen from

406  among seven geometrical shapes (Komatsu and Ideura, 1993; Koida and Komatsu, 2007). The

407    background was uniform 10 cd/m$^2$ gray (x = 0.3127, y = 0.3290). Stimuli were calibrated using a

408    spectrophotometer (Photo Research PR-650). Personal computers controlled the task, presented the

409    visual stimuli and recorded neural signals and eye positions. Eye movements were recorded using the

410    scleral search coil method (Judge et al., 1980). The monkeys were required to maintain fixation within

411    a 2.8° window throughout the trial, except for the saccade response. At the beginning of each trial, a

412    small fixation spot was presented at the center of screen. When the monkeys had gazed at the fixation

413    spot for 500 ms, it turned off, and a sample color stimulus was presented at the center of the display

414    for 500 ms in both Categorization and Discrimination tasks.

415    In the Categorization task, after the sample stimulus was turned off, two small spots of light appeared,

416    one at the center of the visual field, the other 5° to the right (**Figure 1a**). If the sample color belonged

417    to the "reddish" category (sample colors 1-4), the monkeys were rewarded for maintaining fixation on

418    the center spot for another 700 ms ("no-go" response). If the sample color belonged to the "greenish"

419    category (sample colors 8-11), the monkeys were rewarded for making a saccade to the spot on the

420    right ("go" response). For the intermediate colors (sample colors 5-7), the monkeys were rewarded

421    randomly regardless of its behavioral response. In an early phase of the recordings from one monkey

422    (15 neurons), there were no intermediate colors; the "no-go" response was assigned to colors 1-5, the

423    "go" response to colors 6-11.

424    In the Discrimination task, after the sample stimulus was turned off, two choice stimuli appeared 3°

425    above and below the fixation position (**Figure 1b**). The choice stimuli were the same shape and size

426    as the sample stimulus; one was the same color as the sample stimulus, the other a slightly different

427    color. The monkeys were required to make a saccade to the choice stimulus that was the same color as

428    the sample. The two choice colors were three steps apart along the 11 sample colors – that is, the eight

429    choice color pairs included colors #1-4 , #2–5, #3-6, #4-7, #5-8, #6-9, #7-10 and #8-11. This color

430    interval was chosen so as to yield a relatively high discriminability (about 80-90% correct).

431    Throughout the present paper, the term "Discrimination" is used for consistency with our previous

432    study (Koida and Komatsu, 2007). Note that the task is also known as "matching to sample."

433    *Electrophysiological recording*

434    Neuronal activity was recorded with single unit recording from the anterior part of the IT cortex in the

435    monkeys. We could record from 125 neurons in total. The recording region was slightly lateral to the

436    posterior end of the anterior middle temporal sulcus (anterior 9-14 mm in the stereotaxic coordinates,

437    area TE), which is a region where color-selective neurons are concentrated (Komatsu et al., 1992;

438    Matsumora et al., 2008). The activities of single neurons were first isolated with online monitoring

439    during recordings, then subject to offline spike sorting using a template matching algorithm, which

440    confirmed that all of the data reported in this paper were single neuron activities.

441    All data analyses were based on neural responses to the sample colors and the fact that the monkeys

442    saw the same visual stimuli in the Categorization and Discrimination tasks. For this purpose, we

443    analyzed neural spikes recorded up to 550 ms after the sample onset, taking into account the neural

444    response delay to the visual stimuli.

445    Our main results are based on a collection of single unit recordings (not a simultaneous recording of

446    multiple neurons). In the population decoding analyses, we generated "pseudo-population" activities

447    from those single neuron data by randomly resampling the trials, following a procedure reported in a

448    previous study (Fetsch et al., 2012). A caveat of the analysis based on "pseudo-population" is that it

449    omits the noise correlation (i.e., the correlation in trial-to-trial fluctuations) across neurons. As widely

450    recognized, the noise correlation can have profound influences on the information coding by neural

451    population, affecting particularly the resolution of sensory representation. From the decoding

452    perspective, in many cases the noise correlation is generally considered to affect the accuracy of

453    decoding (e.g., error bars added when plotting the decoder outputs) although how noise correlation

454    actually limits the stimulus information is a subject of ongoing debate (Moreno-Bote et al., 2014). In

455    this study, we do not primarily focus on the resolution of neural coding (reflected in the lengths of

456    error bars) but on the "biases" induced by the change in the mean activity in each neuron, which is

457    captured by the present single-unit recording. In addition, a control analysis confirmed by that

458    artificially inducing noise correlations in the studied pseudo-population did not affect the overall

459    results (**Figure 7a**).

460    *Likelihood-based decoding*

461    To visualize and characterize high-dimensional representation by neural populations, we mapped the

462    neural population activity in the stimulus space by decoding the neural activity. From Bayes' rule, the

463    posterior probability on stimulus $s$ under a given neural population activity $r(t)$ is $P(s|r(t)) \propto$

464    $P(r(t)|s)P(s)$. In a full-normative framework, the prior distribution over the stimulus could be

465    further modeled by assuming the hierarchical model with categorical prior on stimulus, $P(s|c)$; that

466    is, $P(s) = \int \mathrm{d}c P(s|c)P(c)$, where $c$ denotes the category information (Tajima et al., 2016). In the

467    present experiments, however, the stimulus was sampled from a uniform distribution, thus the

468    problem reduces to maximizing the likelihood $P(r(t)|s)$. In our analysis, a maximum-likelihood

469    decoder (Földiák, 1991; Sanger, 1996; Jazayeri and Movshon, 2006; Ma et al., 2006; Graf et al., 2011;

470    Fetsch et al., 2012) of the stimulus was constructed based on the neural responses in the

16

471 Discrimination task and then applied to the data for Categorization task to reconstruct the neural

472 population states in the perceptual stimulus space (**Figure 1d**; see also the later descriptions for the

473 rationale behind this procedure).

474 The decoder was constructed based on a standard likelihood-based population decoding approach as

475 follows (Graf et al., 2011; Fetsch et al., 2012). Let $r_i(t)$ be the spike counts for the cell $i$ response at

476 time bin $t$ in a trial. The spike count was derived from a 50-ms boxcar window whose starting point

477 moved with 10-ms step from 0 to 500 ms after the onset of a sample-color stimulus. We first

478 estimated a probability distribution, $P_{\text{Dis}}(r_i(t)|s)$, of responses evoked by stimulus $s$ for each cell and

479 each time bin, based on the data obtained during the discrimination task. This is approximated by a

480 Gaussian distribution with a mean $\mu_i(t; s)$ and variance $\sigma_i(t; s)^2$, which were respectively estimated

481 from the mean and variance in the neural spike count data. The mean responses $\mu_i(t; s)$ to 11 sample

482 stimuli were converted to smooth functions of the stimulus (a real number varying from 1 to 11)

483 through cubic interpolation over the stimulus space, to obtain smooth likelihood functions in the later

484 analysis. The variance estimate was denoised by fitting a linear function, $\sigma_i(t; s)^2 = \alpha_i \times \mu_i(t; s)^2 +$

485 $residual$, with a stimulus- and time-invariant scalar variable (Fano factor), $\alpha_i$, for each cell, in order

486 to capture the potential variability in the Fano factor across neurons. To ensure that the decoder output

487 matches the subject's perception about color identity, we used the trials in which the subjects

488 answered correctly in the task. The Gaussian model naively implies the potential for negative neural

489 activity, the biological meaning of which is unclear. However, this does not cause a problem in the

490 practical data analysis because the analyzed neural responses are always positive, and we can safely

491 equate the analysis with the one based on a rectified Gaussian model that satisfies the non-negativity

492 of the neural responses. In addition, we also tested a Poisson distribution as a generative model of

493 spike count, and confirmed that the results were not qualitatively affected (**Results**).

494 Combining these models of spike-count distributions derived from individual neurons and time bins

495 yielded the likelihood of a population response.

$$L\big(s; \mathbf{r}(t)\big) = -\big(\boldsymbol{r}(t) - \boldsymbol{\mu}(t; s)\big)^{\top} \boldsymbol{\Sigma}(t; s)^{-1} \big(\boldsymbol{r}(t) - \boldsymbol{\mu}(t; s)\big) - \frac{1}{2} \log|\boldsymbol{\Sigma}(t; s)| - \frac{N}{2} \log 2\pi, \qquad (1)$$

496 where $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$ are the mean and covariance of neural population response, respectively. In the main

497 analysis, for simplicity, we assumed independent trial-to-trial variability in the neural firing (Sanger,

498 1996; Dayan and Abbott, 2001; Jazayeri and Movshon, 2006; Ma et al., 2006; Brouwer and Heeger,

499 2009; Fetsch et al., 2012)—we also observed that our main results were not affected by the decoder

500 that takes into account the correlated variability among neurons (**Figure 7**). The joint log-likelihood $L$

501 of a population response of $N$ neurons, $\boldsymbol{r}(t) := (r_1(t), \dots, r_N(t))^{\top}$, given stimulus $s$ is

$$L(s; \boldsymbol{r}(t)) := \log P_{\text{Dis}}(\boldsymbol{r}(t)|s)$$

$$= -\sum_{i=1}^{N} (r_i(t) - \mu_i(t;s))^2 / 2\sigma_i(t;s)^2 - \sum_{i=1}^{N} \log \sigma_i(t;s) - \frac{N}{2} \log 2\pi \tag{2}$$

502   Here, column vector $\boldsymbol{r}(t) := (r_1(t), \dots, r_N(t))^\top$ represents the population activity of all $N$ neurons at

503   time $t$ ($N = 125$ in the present data). . Based on this population activity, the decoder output (stimulus

504   estimate), $s^*$, is given by maximizing the aforementioned likelihood function, $L(s; \boldsymbol{r}(t))$:

$$s^*(\boldsymbol{r}(t)) := \operatorname{argmax}_s L(s; \boldsymbol{r}(t)). \tag{3}$$

505   This equation represents a mapping from the $N$-dimensional population state $\boldsymbol{r}(t)$ to a one-

506   dimensional value, $s^*$, in the stimulus space. We iterated this decoding procedure for each time $t$. We

507   used different decoders for individual time bins, by constructing a generative model of the spike

508   counts for each time bin. The main results were unaffected when we used a single time-invariant

509   decoder (constructed based on the average spike count statistics across 200–550 ms after stimulus

510   onset) for all time bins (**Results**).


511   To analyze the neural responses in the Categorization task, we used the same function, $s^*(\mathbf{r}(t))$ [i.e.,

512   the same mean and variance parameters, $(\mu_i(t;s), \sigma_i(t;s)^2)$, for each neuron] as the decoder that was

513   constructed based on neural activity in the Discrimination task. The decoder constructed based on the

514   Discrimination task data does not necessarily provide an unbiased estimate of the stimulus for the

515   Categorization task. We made use of this potential decoding bias to characterize the difference in

516   neural population responses between the two tasks. If there were any systematic bias, it would suggest

517   that the neural population changes the stimulus representation depending on task context. It is

518   reasonable to construct the stimulus decoder based on neural responses in the Discrimination task

519   because the perceived stimulus identity to be decoded could be confirmed with what the subjects

520   reported in the Discrimination task. By comparing the decoder output to the subjects' behavior, we

521   were able to map the neural population response to the subjects' perception of the stimulus identity.


522   The rationale behind those procedures is as follows: in the Discrimination task, the subject was

523   presented a sample color (e.g., light green), then later identified it by selecting from a pair of similar

524   colors (e.g., the same light green vs. a slightly deeper green). When the subjects correctly identified

525   the presented sample color, by construction, the presented color matched the chosen color, which

526   suggests that they correctly perceived the sample color such that it could be discriminated from other

527   similar colors in the perceptual space. Although such a correspondence between choice and

528   perception is not always guaranteed if the subject's choice is nearly random, it was not the case in this

529   study because the subject showed high correct rate (about 80-90%) in the Discrimination task.

18

530    Nonetheless, there were a few error trials in which the presented colors differed from the chosen

531    colors. In those error trials, it is not straightforward to tell what color was perceived by the subjects; it

532    could be the chosen color, but alternatively, they might have actually perceived the presented color

533    but made a mistake in the response, or they might have simply unattended to the task. Thus, we

534    excluded those error trials from the present analysis, and focused on the correct trials in which the

535    presented and chosen colors were identical. We also confirmed that the overall results were

536    qualitatively maintained when we replicated the same analysis using half the recorded neural

537    population without extremely high or low activity (by eliminating the neurons showing average firing

538    rates outside the 25th-75th percentile of the whole population; **Results**), which excluded the

539    possibility that a small subset of strongly-responding neurons determined the results of decoding.


540    *Fitting the task-dependent components in neural dynamics*

541    To investigate what form of neural response modulation explains the difference between the decoded

542    stimulus dynamics in Discrimination and Categorization, we fitted the population responses in the

543    Categorization task ($r(t)_{|s,\text{Cat}}$) by modulating those in the Discrimination task ($r(t)_{|s,\text{Dis}}$), based on

544    four different models: three feedforward gain-modulation models and a recurrent model.


545    *Gain modulation model 1 (time-invariant, stimulus-independent gains)*: In the time-invariant gain-

546    modulation model, the neural response data in the Discrimination task were modulated so that they

547    simulate the Categorization-task responses. The simulated Categorization-task response of neuron $i$,

548    $\hat{r}_i(t)_{|s,\text{Cat}}$, was provided as

$$\hat{r}_i(t)_{|s,\text{Cat}} := \bar{g}_i \, r_i(t)_{|s,\text{Dis}}, \tag{4}$$

549    where $\bar{g}_i$ denotes a constant gain-modulation for each cell $i$. The gain $\bar{g}_i$ was estimated as follows:

$$\bar{g}_i := \frac{1}{|S|} \sum_{s \in S} \left( \bar{r}_{i|s,\text{Cat}} / \bar{r}_{i|s,\text{Dis}} \right), \tag{5}$$

550    where $S = \{\#1, \#2, \ldots, \#11\}$, and $|S| = 11$. $\bar{r}_{i|s,\text{Cat}} = \langle r_i(t)_{|s,\text{Cat}} \rangle_t$ and $\bar{r}_{i|s,\text{Dis}} = \langle r_i(t)_{|s,\text{Dis}} \rangle_t$ are the

551    time-averaged responses of neuron $i$ to stimulus $s$ in the Categorization and Discrimination tasks,

552    respectively, where $r_i(t)_{|s,\text{Cat}}$ and $r_i(t)_{|s,\text{Dis}}$ are the respective responses of neuron $i$ to stimulus $s$ at

553    time $t$ during the Categorization and Discrimination tasks. The numbers of parameters were 125

554    (corresponding to the number of recorded neurons $N$) in the time-invariant model. We analyzed the

555    simulated population activity in the same procedure used for the actual response during the

556    Categorization task.

19

557     *Gain modulation model 2 (time-variant, stimulus-independent gains)*: Similarly, in the time-variant

558     gain modulation model, the predicted Categorization-task response, $\hat{r}_i(t)_{|s,\text{Cat}}$, was given by

$$\hat{r}_i(t)_{|s,\text{Cat}} := g_i(t)\, r_i(t)_{|s,\text{Dis}}. \tag{6}$$

559     The gain term $g_i(t)$ for each neuron $i$ was estimated as follows:

$$g_i(t) := \frac{1}{|S|} \sum_{s \in S} \left( r_i(t)_{|s,\text{Cat}} / r_i(t)_{|s,\text{Dis}} \right). \tag{7}$$

560     In our main analysis, the number of neurons was $N$=125, thus the numbers of parameters were $125 \times$

561     $51$=6375 (corresponding to the number of recorded neurons $\times$ the number of time bins) in this model.

562     *Gain modulation model 3 (time-invariant, stimulus-dependent gains)*: we also considered a gain

563     modulation depending on the presented stimulus as a control (see Discussion for its biological

564     interpretation). Note that the modulation component for each neuron can be trivially fitted by the gain

565     modulation depending on both stimulus and time, since they are the only variables (except for the task

566     demands) in the present experiment. Thus, here we tried to fit the data with a gain-modulation model

567     in which the neuronal gains depend on the stimulus but not on time. In this model, the predicted

568     Categorization-task response, $\hat{r}_i(t)_{|s,\text{Cat}}$, was given by

$$\hat{r}_i(t)_{|s,\text{Cat}} := g_i(s)\, r_i(t)_{|s,\text{Dis}}. \tag{8}$$

569     The gain term $g_i(t)$ for each neuron $i$ was estimated as follows:

$$g_i(s) := \bar{r}_{i_{|s,\text{Cat}}} / \bar{r}_{i_{|s,\text{Dis}}}, \tag{9}$$

570     The numbers of model parameters were $125 \times 11$=1375 (corresponding to the number of recorded

571     neurons $\times$ the number of sample colors).

572

573     *Recurrent model*: Lastly, we also fitted the neural dynamics with a model that features a recurrent

574     feedback. In the recurrent model, a self-feedback term was added to the responses in the

575     Discrimination task so that the resulting modulated activities fit those recorded in the Categorization

576     task. We assumed a restricted recurrent circuit with a single hidden layer consisting of two nonlinear

577     hidden units. In this model, we assumed mutual connections between the recoded IT neurons and the

578     two hidden units (which could be interpreted as the neural activity outside IT cortex, e.g., the frontal

579     cortex, as modeled in further details later). There was no direct connection between the hidden units,

20

580    resembling two-layer restricted Boltzmann machines (Smolensky, 1986; Hinton, 2002). The model is

581    a simplified version of the circuit model (**Fig. 6a**) that we used to demonstrate the task-dependent

582    change in attractor structures (see the later section); here we use this simplified version in purpose of

583    the quantitative fitting.

584    Based on this model, the hypothetical neural activity in the Categorization task, $\hat{r}(t)_{|s,\text{Cat}} :=$

585    $\left(\hat{r}_1(t)_{|s,\text{Cat}}, \dots, \hat{r}_N(t)_{|s,\text{Cat}}\right)^\top$, was provided as

$$\hat{r}(t)_{|s,\text{Cat}} := r(t)_{|s,\text{Dis}} + Wh(t),$$

$$h(t) := f\left(W^\top r(t-1)_{|s,\text{Cat}} + b\right), \tag{10}$$

586    where the $N \times 2$ matrix $W$ denotes the connectivity weights between the neurons to the two hidden

587    units; the weights are symmetric between the bottom-up and top-down connections (from the neurons

588    to the hidden units, and from the hidden units to the neurons, respectively). $h(t) = (h_1(t), h_2(t))^\top$ is

589    the activities of hidden units at time $t$. The function $f(\cdot) := \tanh(\cdot)$ is the activation function for the

590    hidden units. $b = (b_1, b_2)^\top$ is the bias inputs to the hidden units. $W$ and $b$ were learned from the data,

591    but kept constant across time and different stimuli. To optimize those parameters, we minimized the

592    sum of squared error between the actual and predicted neural activities in the Categorization task,

593    $\left\|r(t)_{|s,\text{Cat}} - \hat{r}(t)_{|s,\text{Cat}}\right\|^2$, with a standard gradient descent method on $W$ and $b$. The number of

594    parameters was $2N + 2 = 252$, corresponding to the total number of connections and the bias inputs.

595    Note that it is not necessarily straightforward to relate those two hidden units directly to the "red" and

596    "green" category neurons modeled because such categorical information is represented in a mixed

597    way in the circuit learned from the real data. Nonetheless, the goodness of fitting with this model

598    demonstrates that the recurrent network with the restricted architecture is capable of describing the

599    neural data quantitatively. It should be also noted that we do not consider that the task switching

600    requires changes in all the connectivity weights among the neurons. Instead, we could assume a more

601    parsimonious mechanism that features the attractor structure in the circuit is modulated through the

602    change in a background input to the circuit (see the later subsection).

603    *Assessment of model-fitting performances*: We assessed the model-fitting performances based on the

604    cross-validation procedure as follows: we randomly divided the data into two non-overlapping sets of

605    trials ("trial set 1" and "trial set 2"), the first of which was used to train models, and the second of

606    which was used to test each model's fitting performances. This procedure ensured that a difference in

607    fitting performance did not reflect overfitting or a difference in the number of parameters. The model-

608    fitting errors, $E_{\text{CV}}$, were quantified by the root mean square errors between the predicted and actual

609    neural population activities, normalized by the "baseline" variability across trials:

21

$$E_{CV} = \frac{\left\langle \left( r_i(t)_{|s,\text{Cat,,trial set 2}} - \hat{r}_i(t)_{|s,\text{Cat,,trial set 2}} \right)^2 \right\rangle_{i,t,s}^{\frac{1}{2}}}{\left\langle \left( r_i(t)_{|s,\text{Cat,,trial set 2}} - r_i(t)_{|s,\text{Cat,,trial set 1}} \right)^2 \right\rangle_{i,t,s}^{\frac{1}{2}}}, \tag{11}$$

610  where $\langle \cdot \rangle_{i,t,s}$ is the average over the cells, time bins, and stimuli. The numerator corresponds to the

611  error in the model prediction, whereas the denominator represents the "baseline" variability within the

612  condition due to the trial-to-trial fluctuations in neural firing. Note that this measure itself is

613  independent of the assumptions about decoders because it is computed directly from the neural

614  population activities.

### *Mutual information analysis*

616  The amount of information about a stimulus carried by the neural population response was also

617  evaluated using mutual information, which does not require any specific assumptions about the

618  decoder or the models of dynamical modulations. The mutual information between the stimulus hue

619  and the neural responses within each time bin $t$ during the Categorization task was given by

$$I_{\text{Cat}}(\text{hue}; t) = I\big(s, \mathbf{r}(t)\big) = \sum_{s,i} P_{\text{Cat}}(r_i(t)|s)P(s)\{\log P_{\text{Cat}}(r_i(t)|s) - \log P_{\text{Cat}}\big(r_i(t)\big)\}. \tag{12}$$

620  where $P_{\text{Cat}}(r_i(t)|s)$ is the probability distribution of the $i$th neuron's response (spike counts) evoked

621  by stimulus $s$ during the Categorization task. The "hue" in the parenthesis indicates that this is the

622  mutual information about the stimulus hue. Similarly, the mutual information between the stimulus

623  category $c \in \{\text{Red}, \text{Green}\}$ and the neural responses within each time bin $t$ was given by

$$I_{\text{Cat}}(\text{cat}; t) = I\big(c, \mathbf{r}(t)\big) = \sum_{c,i} P_{\text{Cat}}(r_i(t)|c)P(c)\{\log P_{\text{Cat}}(r_i(t)|c) - \log P_{\text{Cat}}\big(r_i(t)\big)\}, \tag{13}$$

624  where $P_{\text{Cat}}(r_i(t)|c) = \sum_{s \in S_c} P_{\text{Cat}}(r_i(t)|s)$ , and $c \in \{\text{Red}, \text{Green}\}$ denotes the stimulus category;

625  $S_{\text{Red}} = \{\#1, \#2, \#3, \#4\}$ and $S_{\text{Green}} = \{\#8, \#9, \#10, \#11\}$ are the sets of stimuli belonging to the

626  "Red" and "Green" categories, respectively. The "cat" in the parenthesis indicates that this is the

627  mutual information about the stimulus category. The mutual information values for the Discrimination

628  task, $I_{\text{Dis}}(\text{hue})$ and $I_{\text{Dis}}(\text{cat})$, were provided by substituting $P_{\text{Cat}}(r_i(t)|s)$ in the above equations with

629  the corresponding spike count distributions, $P_{\text{Dis}}(r_i(t)|s)$, obtained during the Discrimination task.

630  The differential mutual information for hue and category were defined by $\Delta I(\text{hue}; t) = I_{\text{Cat}}(\text{hue}; t) - $

631  $I_{\text{Dis}}(\text{hue}; t)$ and $\Delta I(\text{cat}; t) = I_{\text{Cat}}(\text{cat}; t) - I_{\text{Dis}}(\text{cat}; t)$, respectively.

632     We evaluated the cumulative values of mutual information over time (e.g., $\sum_{t'=0}^{t} I_{\mathrm{Dis}}(\mathrm{hue}; t')$ for the

633     cumulative hue information in the Discrimination task). This cumulative mutual information reflects

634     the amount of information obtained by observing the sequence of neural population responses. For

635     this purpose, we used non-overlapping consecutive time bins (each with a duration of 20 ms). Note

636     that the variability in neural responses can be temporally correlated even if we use the non-

637     overlapping time bins, although the magnitude of the autocorrelation generally decreases

638     exponentially over time (Murray et al., 2014). Therefore, this cumulative mutual information should

639     be interpreted as an upper bound of the total information obtained by observing the sequence of the

640     neural population response.

641     *Unsupervised dimensionality reduction analyses*

642     We also conducted several unsupervised dimensionality reduction analyses to compare their results

643     with that of the likelihood-based decoding. First, the standard principal component analysis (PCA)

644     was applied to the set of trial-averaged data points (i.e., population response vectors

645     $\{r(t)_{|s,\mathrm{Cat}}, r(t)_{|s,\mathrm{Dis}}\}|_{s\in\{\#1,\dots,\#11\},\ 0\ \mathrm{ms} \leq t \leq 550\ \mathrm{ms}})$ that varied over time $t$. Second, we performed

646     PCA based on the differential responses between Categorization and Discrimination tasks (i.e.,

647     $\{r(t)_{|s,\mathrm{Cat}} - r(t)_{|s,\mathrm{Dis}}\}|_{s\in\{\#1,\dots,\#11\},\ 0\ \mathrm{ms} \leq t \leq 550\ \mathrm{ms}})$. Lastly, we conducted t-stochastic neighbor

648     embedding (t-SNE) (van der Maaten and Hinton, 2008) on the trial averaged data

649     $(\{r(t)_{|s,\mathrm{Cat}}, r(t)_{|s,\mathrm{Dis}}\}|_{s\in\{\#1,\dots,\#11\},\ 0\ \mathrm{ms} \leq t \leq 550\ \mathrm{ms}})$ to examine the effects of nonlinearity in the

650     unsupervised dimensionality reduction.

651     *A model of context-dependent attractor dynamics*

652     We introduced a simple recurrent model that provides a parsimonious explanation for the observed

653     context-dependent change in attractor dynamics (see **Figure 6a**, **Results**). The model assumed

654     bidirectional interactions between *n* hue-selective neurons (hereafter, *hue neurons*) in IT cortex and

655     two groups of category-selective neurons (*category neurons*) outside IT cortex; for example, such

656     neurons that encode category have been found in the prefrontal cortex (Freedman et al., 2001; McKee

657     et al., 2014). This circuit share the basic architecture with our previous model that was proposed for

658     general categorical inference (Tajima et al., 2016); here we extend this model to explain the context

659     dependent bifurcation of attractor dynamics. Note that the category- and hue-neurons in this model

660     should not be confused with the terms 'Categorization-' and 'Discrimination-task preferred cells' used

661     in the previous study (Koida and Komatsu, 2007), which were the labels on the IT neurons introduced

662     to describe the polarity of task-dependent modulation for each cell, and not relevant to the current

663     model.

23

664    The dynamics of category neurons were described by differential equations as follows:

$$T_C \dot{C}_1 = -C_1 + f\big(\boldsymbol{W}_1^{\mathrm{BU}} \cdot \boldsymbol{H} + B\big), \tag{14}$$

$$T_C \dot{C}_2 = -C_2 + f\big(\boldsymbol{W}_2^{\mathrm{BU}} \cdot \boldsymbol{H} + B\big), \tag{15}$$

$$\boldsymbol{H} = \boldsymbol{W}_1^{\mathrm{TD}} C_1 + \boldsymbol{W}_2^{\mathrm{TD}} C_2 + \boldsymbol{I}(s, t), \tag{16}$$

665    where the dots between variables denote inner products of vectors. $T_C$ is the time constant for the

666    dynamics of category neurons, which was set as $T_C = 100$ ms in the simulation, roughly matched to

667    the order of time constants in cortical neurons (Murray et al., 2014). $C_1$ and $C_2$ are scalar values

668    representing mean activity of red- and green-preferring category neurons, respectively. The time

669    constant for hue-neurons was neglected for the sake of the tractability in nullclines analysis. The

670    faster dynamics in sensory neurons compared to those in higher-area is consistent with a previous

671    report (Murray et al., 2014). We also confirmed that assuming non-zero time constant in hue neurons

672    did not change the qualitative behavior of the model. The activation function in the simulation was

673    given by a sigmoid function, $f(x) = \exp(kx)/(1 + \exp(kx))$, where $k = 0.2$, though the precise

674    form of the activation function was not critical for the emergence of bistability as long as the neural

675    activity was described by a monotonic saturating function. $\boldsymbol{H} \coloneqq (H_1, \dots, H_n)^{\top}$ is a vector

676    representing the population activity of hue-neurons with different preferred stimuli (varying from red

677    to green), which receive sensory input, $\boldsymbol{I}(s, t) \coloneqq (I_1, \dots, I_n)^{\top}$, from the earlier visual cortex. The hue

678    neurons interact with category-neuron groups $C_1$ and $C_2$ through bottom-up and top-down connections

679    with weights $(\boldsymbol{W}_1^{\mathrm{BU}}, \boldsymbol{W}_2^{\mathrm{BU}})$ and $(\boldsymbol{W}_1^{\mathrm{TD}}, \boldsymbol{W}_2^{\mathrm{TD}})$, respectively, where the connectivity weights were

680    expressed as vectors (e.g., $\boldsymbol{W}_1^{\mathrm{BU}} \coloneqq \big(W_{11}^{\mathrm{BU}}, \dots, W_{1n}^{\mathrm{BU}}\big)^{\top}$). The category neurons also receive a common

681    background input, $B$. We assume that this background input is the only component that depend on

682    task demand in this circuit.

683    In the simulation, the numbers of hue-neurons were set to $n = 300$, although the size of neural

684    population did not have major effect on the results of simulation. Sensory input to hue-neuron $i$ was

685    modeled using a von Mises function, $I_i(s, t) = g(t) \exp(\kappa \cos(s - s_i^{\mathrm{pref}}))$, where the sharpness

686    parameter $\kappa = 2$; $s \in [-\pi/2, \pi/2]$ is the stimulus hue, which varied from red to green, and $s_i^{\mathrm{pref}}$ is

687    the preferred hue of neuron $i$; $g(t) = 0.5e^{(t-50)/100} + 0.5$ for $t > 50$, $g(t) = 0$ for $t \le 0$. The

688    preferred hues were distributed uniformly across the entire hue circle, $[-\pi, \pi]$. Each category-neuron

689    group contained 150 cells, which were uniform within each group. The connectivity weight between

690    hue neuron $i$ and category-neuron group $j$ was modeled by $W_{ji}^{\mathrm{BU}} = W_{ji}^{\mathrm{TD}} = a \cos(s_j^{\mathrm{Cat}} - s_i^{\mathrm{pref}})$,

691    where $a = 10/n$, $s_j^{\mathrm{Cat}} = (-1)^j$ is the preferred hue of $C_j$. For simplicity, the bottom-up and top-

24

692    down weights were assumed to be symmetric. We assumed that all the model parameters except for

693    the background input $B$ were the same between different task conditions. The differential equations

694    were solved with the Euler method with a unit step size of 0.25 ms.

## Acknowledgements

## References

Akrami A, Liu Y, Treves A, Jagadeesh B (2009) Converging neuronal activity in inferior temporal cortex during the classification of morphed stimuli. Cereb Cortex 19:760–776.

Brendel W, Machens CK (2011) Demixed Principal Component Analysis. Adv Neural Inf Process Syst.

Brouwer GJ, Heeger DJ (2009) Decoding and reconstructing color from responses in human visual cortex. J Neurosci 29:13992–14003.

Brouwer GJ, Heeger DJ (2013) Categorical clustering of the neural representation of color. J Neurosci 33:15454–15465.

Chen Y, Martinez-Conde S, Macknik SL, Bereshpolova Y, Swadlow HA, Alonso J-M (2008) Task difficulty modulates the activity of specific neuronal populations in primary visual cortex. Nat Neurosci 11:974–982.

Coen-Cagli R, Kohn A, Schwartz O (2015) Flexible gating of contextual influences in natural vision. Nat Neurosci:1–11.

Dayan P, Abbott LF (2001) Theoretical Neuroscience: Computational and Mathematical Modeling of Neural Systems. Cambridge: MIT Press.

Ecker AS, Denfield GH, Bethge M, Tolias AS (2016) On the structure of population activity under fluctuations in attentional state. J Neurosci 36:1775–1789.

Fetsch CR, Pouget A, DeAngelis GC, Angelaki DE (2012) Neural correlates of reliability-based cue weighting during multisensory integration. Nat Neurosci 15:146–154.

Földiák P (1991) Learning invariance from transformation sequences. Neural Comput 3:194–200.

Freedman D, Riesenhuber M, Poggio T, Miller EK (2002) Visual categorization and the primate prefrontal cortex: neurophysiology and behavior. J Neurophysiol 88:929–941.

Freedman D, Riesenhuber M, Poggio T, Miller EK (2003) A comparison of primate prefrontal and inferior temporal cortices during visual categorization. J Neurosci 23:5235–5246.

Freedman DJ, Riesenhuber M, Poggio T, Miller EK (2001) Categorical representation of visual stimuli in the primate prefrontal cortex. Science (80- ) 312:312–316.

Furman M, Wang XJ (2008) Similarity Effect and Optimal Control of Multiple-Choice Decision

Making. Neuron 60:1153–1168.

Graf ABA, Kohn A, Jazayeri M, Movshon JA (2011) Decoding the activity of neuronal populations in macaque primary visual cortex. Nat Neurosci 14:239–245.

Greenberg DS, Houweling AR, Kerr JND (2008) Population imaging of ongoing neuronal activity in the visual cortex of awake rats. Nat Neurosci 11:749–751.

Haefner RM, Berkes P, Fiser J (2016) Perceptual decision-making as probabilistic inference by neural sampling. Neuron 90:649–660.

Hinton GE (2002) Training products of experts by minimizing contrastive divergence. Neural Comput 14:1771–1800.

Jazayeri M, Movshon JA (2006) Optimal representation of sensory information by neural populations. Nat Neurosci 9:690–696.

Jazayeri M, Movshon JA (2007) A new perceptual illusion reveals mechanisms of sensory decoding. Nature 446:912–915.

Judge SJ, Richmond BJ, Chu FC (1980) Implantation of magnetic search coils for measurement of eye position: an improved method. Vision Res 20:535–538.

Kastner S, Ungerleider LG (2000) Mechanisms of visual attention in human visual cortex. Annu Rev Neurosci 23:315–341.

Kobak D, Brendel W, Constantinidis C, Feierstein CE, Kepecs A, Mainen ZF, Romo R, Qi X, Uchida N, Machens CK (2016) Demixed principal component analysis of neural population data. Elife 5:1–37.

Koida K, Komatsu H (2007) Effects of task demands on the responses of color-selective neurons in the inferior temporal cortex. Nat Neurosci 10:108–116.

Komatsu H, Ideura Y (1993) Relationships Between Color, Shape, and Pattern Selectivities of Neurons in the Inferior Temporal Cortex of the Monkey. J Neurophysiol 70:677–694.

Komatsu H, Ideura Y, Kaji S, Yamane S (1992) Color selectivity of neurons in the inferior temporal cortex of the awake macaque monkey. J Neurosci 12:408–424.

Lamme VA, Zipser K, Spekreijse H (1998) Figure-ground activity in primary visual cortex is suppressed by anesthesia. Proc Natl Acad Sci U S A 95:3263–3268.

Lamme VAF, Zipser K (2002) Masking interrupts figure-ground signals in V1. J Cogn Neurosci 14:1044–1053.

Ma WJ, Beck JM, Latham PE, Pouget A (2006) Bayesian inference with probabilistic population codes. Nat Neurosci 9:1432–1438.

Machens CK, Romo R, Brody CD (2005) Flexible control of mutual inhibition: a neural model of two-interval discrimination. Science 307:1121–1124.

Mante V, Sussillo D, Shenoy K V., Newsome WT (2013) Context-dependent computation by recurrent dynamics in prefrontal cortex. Nature 503:78–84.

Matsumora T, Koida K, Komatsu H (2008) Relationship Between Color Discrimination and Neural Responses in the Inferior Temporal Cortex of the Monkey. J Neurophysiol 100:3361–3374.

McKee JL, Riesenhuber M, Miller EK, Freedman DJ (2014) Task Dependence of Visual and Category Representations in Prefrontal and Inferior Temporal Cortices. J Neurosci 34:16065–16075.

Meyers EM, Qi X-L, Constantinidis C (2012) Incorporation of new information into prefrontal cortical activity after learning working memory tasks. Proc Natl Acad Sci U S A 109:4651–4656.

Mirabella G, Bertini G, Samengo I, Kilavik BE, Frilli D, Della Libera C, Chelazzi L (2007) Neurons in area V4 of the macaque translate attended visual features into behaviorally relevant categories. Neuron 54:303–318.

Moreno-Bote R, Beck J, Kanitscheider I, Pitkow X, Latham P, Pouget A (2014) Information-limiting correlations. Nat Neurosci 17:1410–1417.

Murray JD, Bernacchia A, Freedman DJ, Romo R, Wallis JD, Cai X, Padoa-Schioppa C, Pasternak T, Seo H, Lee D, Wang X-J (2014) A hierarchy of intrinsic timescales across primate cortex. Nat Neurosci 17:1661–1663.

Nienborg H, Cumming BG (2009) Decision-related activity in sensory neurons reflects more than a neuron's causal effect. Nature 459:89–92.

Okazawa G, Tajima S, Komatsu H (2015) Image statistics underlying natural texture selectivity of neurons in macaque V4. Proc Natl Acad Sci 112:E351–E360.

Reynolds JH, Heeger DJ (2009) The Normalization Model of Attention. Neuron 61:168–185.

Roweis ST, Saul LK (2000) Nonlinear dimensionality reduction by locally linear embedding. Science (80- ) 290:2323–2326.

Sadakane O, Ozeki H, Naito T, Akasaki T, Kasamatsu T, Sato H (2006) Contrast-dependent, contextual response modulation in primary visual cortex and lateral geniculate nucleus of the cat. Eur J Neurosci 23:1633–1642.

Sanger T (1996) Probability density estimation for the interpretation of neural population codes. J Neurophysiol 76:2799–2793.

Sasaki R, Uka T (2009) Dynamic readout of behaviorally relevant signals from area MT during task switching. Neuron 62:147–157.

Sceniak MP, Hawken MJ, Shapley R (2002) Contrast-dependent changes in spatial frequency tuning of macaque V1 neurons: effects of a changing receptive field size. J Neurophysiol 88:1363–1373.

Sceniak MP, Ringach DL, Hawken MJ, Shapley R (1999) Contrast's effect on spatial summation by macaque V1 neurons. Nat Neurosci 2:733–739.

Siegel M, Buschman TJ, Miller EK (2015) Cortical information flow during flexible sensorimotor decisions. Science 348:1352–1355.

Smolensky P (1986) Information processing in dynamical systems: foundations of harmony theory. In: Parallel Distributed Processing: Explorations in the Microstructure of Cognition, Volume 1: Foundations (Rumelhart DE, McLelland JL, eds), pp 194–281. MIT Press.

Solomon SG, Kohn A (2014) Moving sensory adaptation beyond suppressive effects in single neurons. Curr Biol 24:R1012–R1022.

Stocker AA, Simoncelli EP (2008) A Bayesian Model of Conditioned Perception. Adv Neural Infromation Process Syst.

Stokes M, Kusunoki M, Sigala N, Nili H (2013) Dynamic Coding for Cognitive Control in Prefrontal Cortex. Neuron 78:364–375.

Tajima CI, Tajima S, Koida K, Komatsu H, Aihara K, Suzuki H (2016) Population code dynamics in categorical perception. Sci Rep 6:1–13.

Tajima S, Watanabe M, Imai C, Ueno K, Asamizuya T, Sun P, Tanaka K, Cheng K (2010) Opposing effects of contextual surround in human early visual cortex revealed by functional magnetic resonance imaging with continuously modulated visual stimuli. J Neurosci 30:3264–3270.

Tenenbaum JB, de Silva V, Langford JC (2000) A global geometric framework for nonlinear dimensionality reduction. Science (80- ) 290:2319–2323.

Toth LJ, Rao SC, Kim DS, Somers D, Sur M (1996) Subthreshold facilitation and suppression in primary visual cortex revealed by intrinsic signal imaging. Proc Natl Acad Sci U S A 93:9869–9874.

Treue S, Martínez Trujillo JC (1999) Feature-based attention influences motion processing gain in macaque visual cortex. Nature 399:575–579.

Uchikawa K, Shinoda H (1996) Influence of Basic Color Categories on Color Memory Discriminatibn. Color Res Appl 21:430–439.

Uchikawa K, Sugiyama T (1993) Effects of eleven basic color categories on color memory. Invest Ophthalmol Vis Sci 34:745.

van der Maaten L, Hinton G (2008) Visualizing Data using t-SNE. J Mach Learn Res 9:2579–2605.

Wallis JD, Anderson KC, Miller EK (2001) Single neurons in prefrontal cortex encode abstract rules. Nature 411:953–956.

Wallis JD, Miller EK (2003) Neuronal activity in primate dorsolateral and orbital prefrontal cortex during performance of a reward preference task. Eur J Neurosci 18:2069–2081.

Wang X (2002) Probabilistic Decision Making by Slow Reverberation in Cortical Circuits. Neuron 36:955–968.

Wang XJ (2008) Decision Making in Recurrent Neuronal Circuits. Neuron 60:215–234.

Wimmer K, Compte A, Roxin A, Peixoto D, Renart A, de la Rocha J (2015) The dynamics of sensory integration in a hierarchical network explains choice probabilities in MT. Nat Commun 6:1–13.
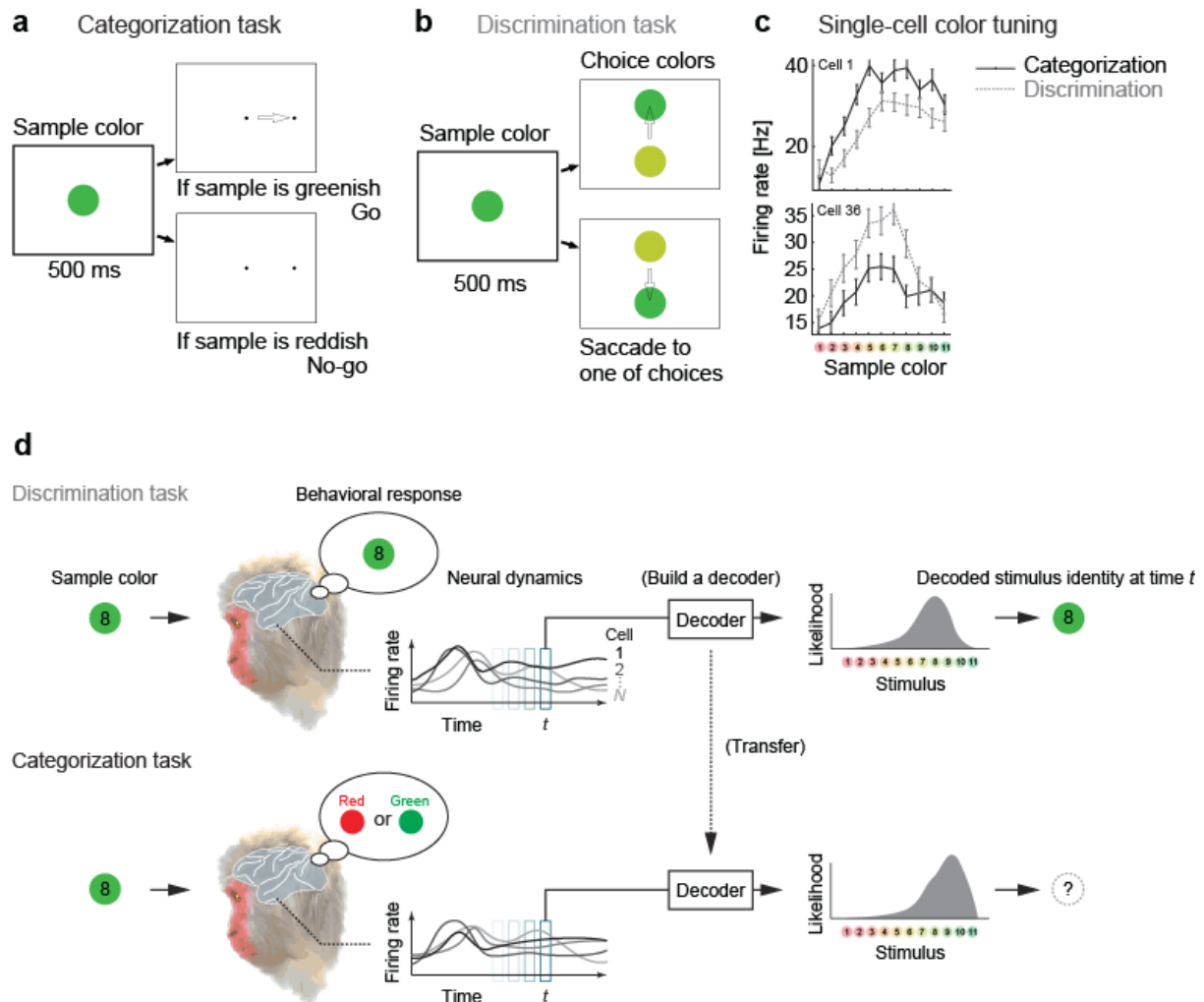
29

# Figures



**Figure 1. Color selectivity of IT neurons and the decoding-based stimulus reconstruction.**

(a) In the Categorization task, subjects classified sample colors into either a "reddish" or "greenish" group.

(b) In the Discrimination task, they selected the physically identical colors.

(c) Color tuning curves of four representative neurons in the Categorization and Discrimination tasks. The color selectivity and task effect are variable across neurons. The average firing rates during the period spanning 100–500 ms after stimulus onset are shown. The error bars indicate the s.e.m across trials.

(d) The likelihood-based decoding for reconstructing the stimulus representation by the neural population.

**Figure 2. Population dynamics in the perceptual domain.**

(a) State-space trajectories during the Categorization and Discrimination tasks. Small markers show the population states 100-550 ms after stimulus onset in 10-ms steps. Large markers indicate the endpoint (550 ms). The colors of the trajectories and numbers around them refer to the presented stimulus.

(b) During the Categorization task, the decoded stimulus was shifted toward either the "reddish" or "greenish" extreme during the late responses but not during the early responses. The thickness of the curve represents the 25th–75th percentile on resampling. The yellow arrow on the horizontal axis indicates the sample color corresponding to the categorical boundary estimated from the behavior (subject's 50% response threshold) in the Categorization task.

(c) Evolution of the task-dependent difference in the decoded stimulus (the curve with a shade), as compared to the population average firing rate (the black solid and dashed curves). The difference in the decoded stimulus was larger in the late period (450–550 ms after the stimulus onset) than the early period (100–200 ms) (P=0.002, bootstrap test). The figure shows data averaged across all stimuli. The black curve and shaded area represent the median ± 25th percentile on resampling.

(d) The time-evolution of the gain modulation models applied to the Discrimination-task data (the recurrent model and the three different gain-modulation models) compared to the actual evolution in the Categorization task (black curve, the same as in Fig. 3c). The curve and shaded areas represent the median ± 25th percentile on resampling.
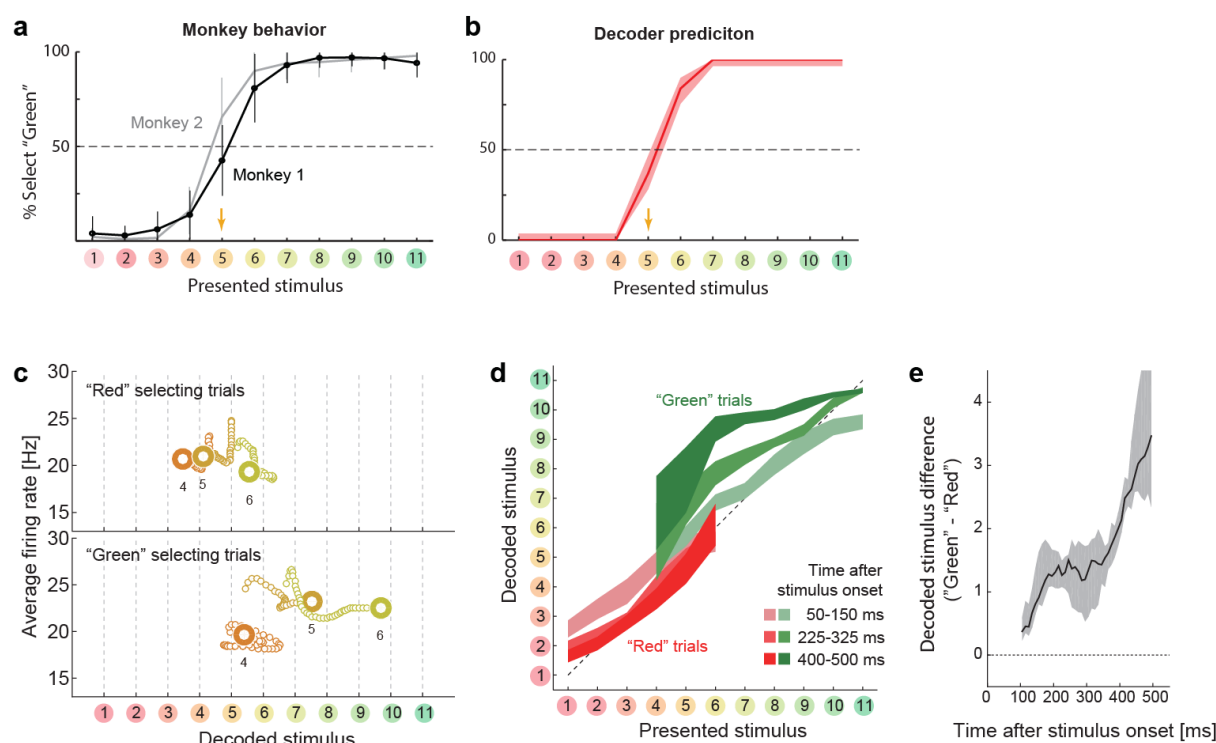
**Figure 3. Choice-related dynamics.**

(a) Actual monkey behavior. Note that the monkeys' subjective category borders were consistent with the decoder output. The error bar is standard error of mean. The yellow arrow on the horizontal axis indicates the sample color corresponding to the putative categorical boundary based on the behavior.

(b) Fraction of selecting Green category predicted by the likelihood-based decoding. The shaded area indicates the 25th–75th percentile on resampling.

(c) The same analysis as **Figure 2a** (top) but with trial sets segregated based on whether the monkeys selected the "Red" or "Green" category. The results for stimuli #4–6 are shown.

(d) The same analysis as **Figure 2b**, except that the trials were segregated based on the behavioral outcome. For stimuli #1–3 (#7–11), only the "Red" ("Green") selecting trials were analyzed because the subjects rarely selected the other option for those stimuli.

(e) Evolution of difference in the decoded color. Data were averaged across stimuli 4-6. The difference in the decoded stimulus was larger during the late period (450-550 ms) than the early period (50–150 ms) (P=0.002, permutation test).
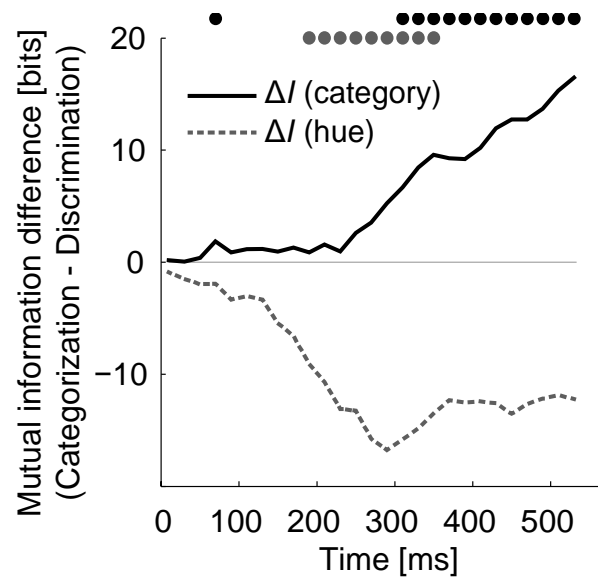
**Figure 4. Modulation increases task-relevant information.**

The figure shows the evolution of the cumulative mutual information difference after the stimulus onset. The dots indicate the statistical significance (P<0.05, permutation test; top black dots: the category information; bottom gray dots: the hue information).
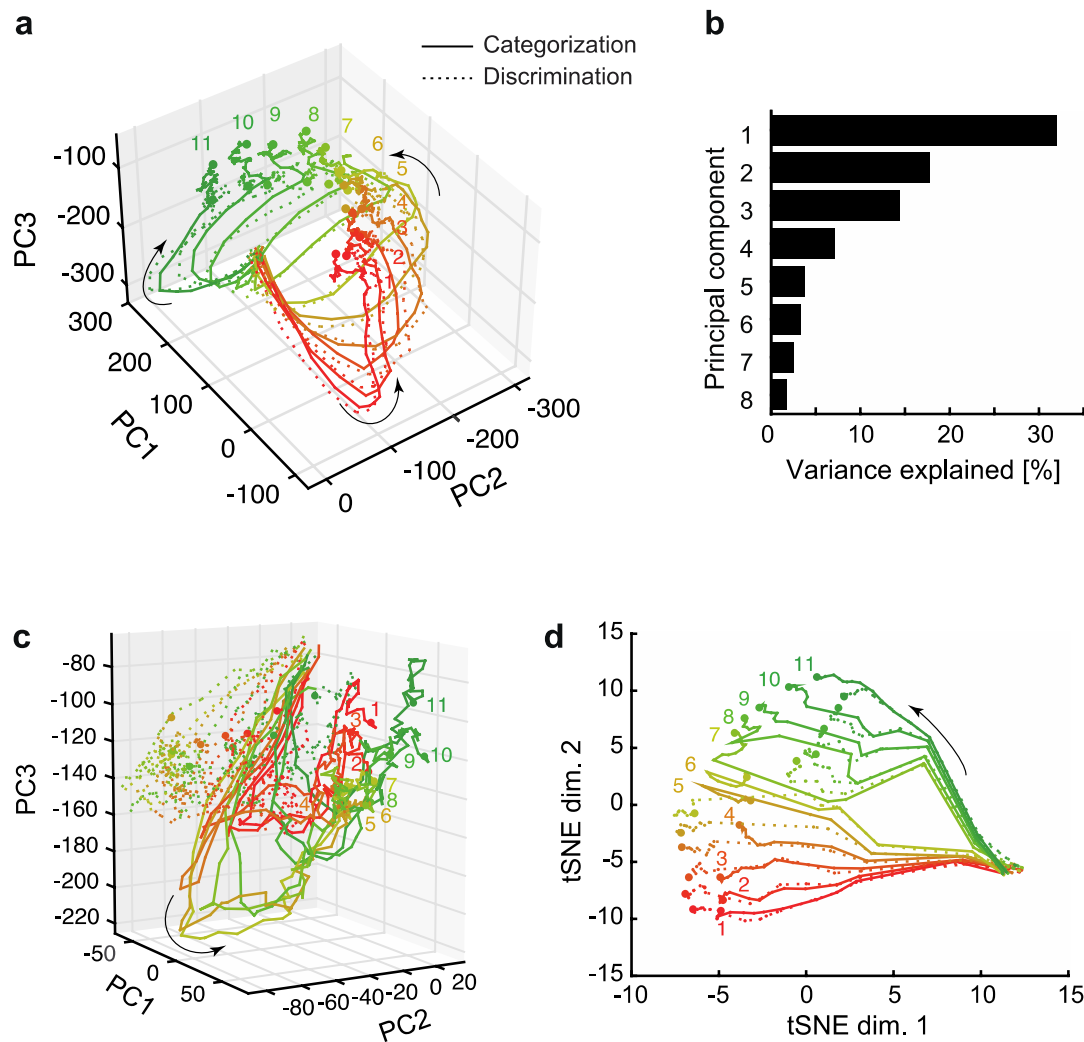
**Figure 5. Comparison to other dimensionality reduction methods.**

(a) Reconstructed neural population dynamics in 0–550 ms after stimulus onset are shown as trajectories in the space spanned by the three principal components. The numbers correspond to the stimulus index.

(b) The fraction of data variance explained by each eigenvector.

(c) PCA based on the differential activities between Categorization and Discrimination tasks.

(d) The result of dimensionality reduction with t-stochastic neighbor embedding (tSNE).

**Figure 6. Bifurcation of attractor dynamics in a simple neural circuit model.**

(a) Schematic of the model circuit architecture. IT hue-selective neurons (hereafter, hue-neurons), $H$, with different preferred stimuli (varying from red to green) receive sensory input, $I(s)$, from the earlier visual cortex. The hue neurons interact with category neuron groups $C_1$ and $C_2$ through bottom-up and top-down connections with weights ($W_1^{\mathrm{BU}}$, $W_2^{\mathrm{BU}}$) and ($W_1^{\mathrm{TD}}$, $W_2^{\mathrm{TD}}$), respectively. The category neurons also receive a common background input, $B$, whose magnitude depends on task context. Note that the modeled hue-neurons covered entire hue circle, $[-\pi, \pi]$, although the figure shows only the half of them, corresponding to the stimulus range from red to green.

(b) Activity evolution represented in the space of category-neurons in the Discrimination task (where the background input $B = -8$). The red (dashed) and green (solid) curves represent nullclines for category-neurons 1 and 2, respectively. The black line shows a dynamical trajectory, starting from (0, 0) and ending at a filled circle. The gray arrows schematically illustrate the vector field.

35

(c) The same analysis as in panel c but in the Categorization task (where $B = -1$). The black and blue lines show two different dynamical trajectories, starting from (-0.01, 0.01) and (-0.01, 0.01), respectively (indicated as numbers "1" and "2" in the figure), and separately ending at filled circles.

(d) The number of stable fixed points is controlled by the parameter $B$. Here, the parameter $B$ was continuously varied as the bifurcation parameter while the other parameters were kept constant. The vertical axis shows the difference of category neuron activities, $C_2 - C_1$, corresponding to the fixed points. The solid black and blue curves show the stable fixed points; the dashed line indicates the unstable fixed point. The stimulus value was $s = 0$.

(e–k) The model replicates recorded neural population dynamics.

(e) Presented and decoded stimuli. The same analysis as in **Figure 2b** was applied to the dynamics of the modeled hue-neurons.

(f) The same as panel e, except that the trials were segregated based on the choices (i.e., to which fixed point the neural states were attracted). The plot corresponds to **Figure 3d**.

(g) Evolution of difference in the decoded color, corresponding to **Figure 3e**.

(h) Mean activity of the entire neural population, corresponding to **Figure 2c**, inset.

(i) Differences in mutual information about category and hue between the Categorization and Discrimination tasks, corresponding to **Figure 4**.

(j) The activity trajectories of the modeled hue-neurons population in PCA space, corresponding to **Figure 5**a.

Note that the scaling of the stimulus coordinate (ranging from $-\pi/2$ to $\pi/2$) used in the model is not necessarily identical to that of experimental stimuli (index by colors #1 – #11), and point of this modeling is to replicate the qualitative aspects of the data.

**Figure 7. Robustness of the results to changes in the decoder.**

We replicated the main results of the paper using four different decoders. Both the stimulus-dependent clustering effect and the temporal evolution were replicated with those decoders. (Left) State-space trajectories during the Categorization task (corresponding to **Figure 2a**, top). (Right) Time-evolution of the gain modulation models applied to the Discrimination-task data (corresponding to **Figure 2d**).

(a) Results obtained by simulating noise correlation among neurons. Here we assumed that the covariance $\sigma_{ij}^2$ between two different neurons, $i$ and $j$, is proportional to the correlation between their mean spike counts: $\sigma_{ij}^2 = k\sqrt{\alpha_i \alpha_j}\mu_i\mu_j$, where $k$ is a constant shared across all neuron pairs (here, $k = 1$), and $\alpha_i$ is the Fano factor for neuron $i$.

(b) Results based on a subset of the recorded cell population; excluded are cells showing extremely high or low activity, as compared to the typical firing rate of the population. We only used cells whose average firing rates (the average across all stimuli and time bins) were within the 25th-75th percentile of the whole population.

(c) Results with a decoder based on Poisson spike variability. The generative model of neuron $i$'s spike count in response to stimulus $s$ at time $t$ was given by $P(r_i(t)|s) = \mu_i(t;s)^{r_i(t)} \exp(-\mu_i(t;s)) / r_i(t)!$ (i.e., the log likelihood was provided by $L(s; \mathbf{r}(t)) := \log P(\mathbf{r}(t)|s) = \sum_{i=1}^{N} r_i(t) \log \mu_i(t;s) - \sum_{i=1}^{N} \mu_i(t;s) + const.$)

(d) Results with a time-invariant decoder. The mean and variance of each neuron's spike count were computed by pooling all the time bins during the period spanning 200–550 ms after stimulus onset.