



# Prefrontal cortex exhibits multidimensional dynamic encoding during decision-making

Mikio C. Aoi<sup>1,2</sup>✉, Valerio Mante<sup>3</sup> and Jonathan W. Pillow<sup>1</sup>

**Recent work has suggested that the prefrontal cortex (PFC) plays a key role in context-dependent perceptual decision-making. In this study, we addressed that role using a new method for identifying task-relevant dimensions of neural population activity. Specifically, we show that the PFC has a multidimensional code for context, decisions and both relevant and irrelevant sensory information. Moreover, these representations evolve in time, with an early linear accumulation phase followed by a phase with rotational dynamics. We identify the dimensions of neural activity associated with these phases and show that they do not arise from distinct populations but from a single population with broad tuning characteristics. Finally, we use model-based decoding to show that the transition from linear to rotational dynamics coincides with a plateau in decoding accuracy, revealing that rotational dynamics in the PFC preserve sensory choice information for the duration of the stimulus integration period.**

A large body of work has aimed to identify the precise computational roles of various brain regions during perceptual decision-making<sup>1–8</sup>. Recent interest has centered on the PFC, which has been shown to carry a wide range of sensory, cognitive and motor signals relevant for integrating sensory information and making decisions<sup>1,4,6,7,9–12</sup>. A barrier to understanding the PFC's functional role, however, is that PFC neurons exhibit mixed selectivity, characterized by heterogeneous tuning to multiple task variables<sup>13</sup>. These idiosyncratic single-neuron responses make it difficult to gain insight into the population-level representation of different sensory and cognitive variables.<sup>14–16</sup>

In this study, we analyzed the population-level representation of information in the PFC using model-based targeted dimensionality reduction (mTDR), a general method for identifying the dimensions of population activity that encode information about different task variables over time. We applied this method to electrophysiology data recorded during a context-dependent perceptual decision-making task<sup>1</sup>, in which a context cue determined what kind of sensory information (color or motion) should be used for making a binary decision on each trial (Fig. 1a,b). In contrast to previous findings, our analysis revealed that the encoding of decisions, context and relevant as well as irrelevant stimulus variables exhibited rotational dynamics in a multidimensional subspace, involving modulation of two or more orthogonal neural activity patterns over time.

We also introduce a new unsupervised method, sequential principal component analysis (seqPCA), for decomposing multidimensional representations into an ordered set of axes that capture the temporal order in which information about each variable becomes available. This method reveals that multidimensional trajectories can be divided into an early linear phase, followed by a later rotational phase. We used model-based decoding under the mTDR framework to show that the transition between these phases corresponded to a saturation in decoding accuracy for both sensory and decision information, suggesting that the population did not continue to accumulate sensory information during the rotational phase.

Taken together, these results substantially extend the prevailing picture of decision encoding in the PFC: rather than integrating evidence along a single dimension of population activity, with

amplitude that reflects accumulated evidence<sup>17</sup>, neural population activity enters a phase of rotational dynamics that maintains information about the choice as well as relevant and irrelevant sensory information over the entire course of a single trial.<sup>18–20</sup>

## Results

**mTDR.** To characterize population-level representations in the PFC, we introduce mTDR, a dimensionality-reduction method that seeks to identify the dimensions of population activity that carry information about distinct task variables. We illustrate the basic intuition for mTDR with a hypothetical three-neuron population in a perceptual decision-making task (Fig. 2). For this example, there are two task variables of interest: a sensory stimulus  $x_s$  and a binary decision variable  $x_c$ . These variables modulate the firing rates in different ways, producing a diverse pattern of population responses across conditions (Fig. 2a).

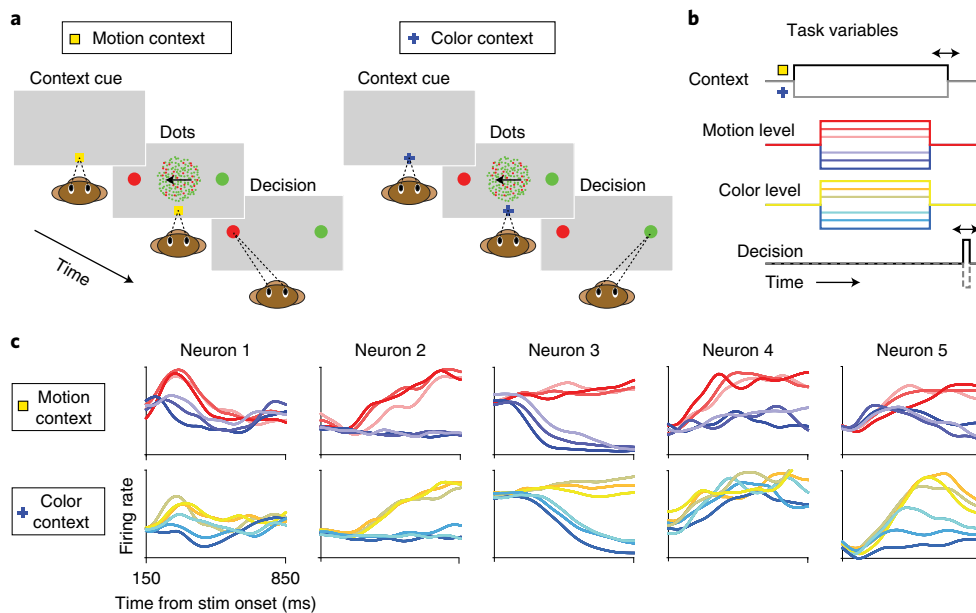
The population-level response can be described as trajectories in a three-dimensional (3D) state space, where the coordinates along each axis correspond to the firing rates of the three neurons (Fig. 2b). Although the full space is 3D, the trajectories exhibit low-dimensional structure that is not apparent from the firing rates alone. Specifically, the population activity is confined to a two-dimensional (2D) plane defined by two axes: a one-dimensional (1D) 'stimulus axis' (blue arrow) captures information about the stimulus strength, whereas a 1D 'decision axis' (red arrow) captures information about the choice. Projecting the population response onto each of these axes reveals the time course of information about stimulus level and choice, respectively (Fig. 2c).

The goal of mTDR is to identify these encoding subspaces from neural population data. For our three-neuron example, the mTDR model describes the time evolution of the population response  $\mathbf{y}(t)$ , a vector of three neural firing rates, as:

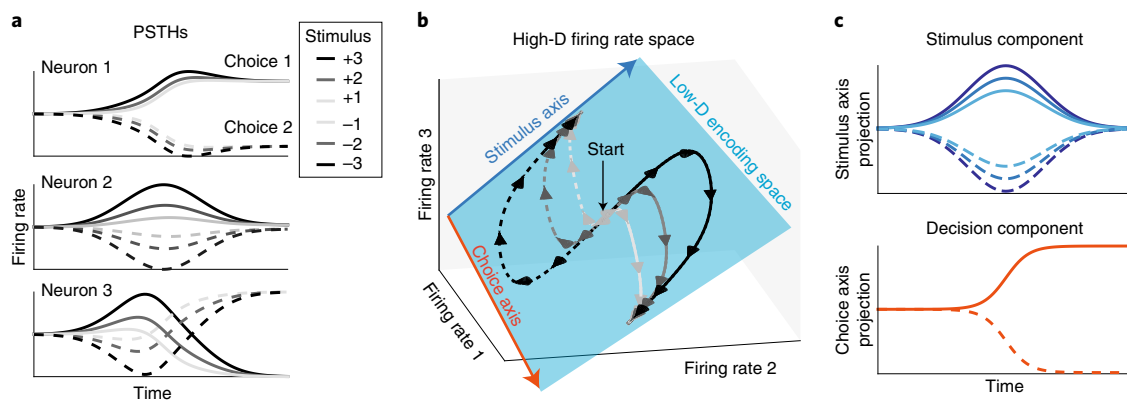
$$\mathbf{y}(t) = x_{\text{stim}}(\mathbf{w}_{\text{stim}} \cdot \mathbf{s}_{\text{stim}}(t)) + x_{\text{choice}}(\mathbf{w}_{\text{choice}} \cdot \mathbf{s}_{\text{choice}}(t)) + \mathbf{1} \quad (1)$$

where  $x_{\text{stim}}$  is the stimulus variable, which takes one of six values from  $[-3, -2, -1, +1, +2, +3]$ , indicating the level of sensory

<sup>1</sup>Department of Psychology and Princeton Neuroscience Institute, Princeton University, Princeton, NJ, USA. <sup>2</sup>Division of Biological Sciences & Halıcıoğlu Data Science Institute, University of California, San Diego, La Jolla, CA, USA. <sup>3</sup>Institute of Neuroinformatics and Neuroscience Center Zurich, University of Zurich and ETH Zurich, Zurich, Switzerland. ✉e-mail: [maoi@ucsd.edu](mailto:maoi@ucsd.edu)



**Fig. 1 | Context-dependent decision-making task and neural responses.** **a**, On each trial, the animal was presented with a context cue (yellow dot or blue cross) indicating which dimension of the stimulus the animal was to attend to, followed by a stimulus of colored, moving dots. On motion-context trials, the animal was cued to respond to the dominant dot motion direction. In color-context trials, the animal was cued to respond to the dominant color of the dots. **b**, The strength of both the color (red or green) and motion (left or right) stimulus took on one of six possible coherence levels, making for many possible task conditions (2 choices  $\times$  2 contexts  $\times$  6 motion strengths  $\times$  6 color strengths = 144 possible combinations). **c**, PSTHs of representative neurons for monkey A. Motion-context PSTHs were sorted by motion coherence and averaged over color coherence. Color-context PSTHs were sorted by color coherence and averaged over motion coherence. Red-indigo color scale indicates motion coherence, where red indicates the preferred motion direction. Gold-blue color scale indicates color coherence, where gold indicates the preferred color direction.



**Fig. 2 | Schematic illustrating low-dimensional population-level encoding in a binary sensory decision-making task.** **a**, Conditional PSTHs for three neurons that exhibit mixed selectivity to a stimulus variable (taking on six different values) and a choice variable (taking on two values). **b**, Modulations of the PSTHs by the task variables span a 2D ‘encoding subspace’, which is low-dimensional relative to the 3D space of firing rates. In this case, a 1D stimulus-encoding subspace (blue arrow) captures all information about the stimulus value, whereas a 1D choice-encoding subspace (red arrow) captures all information about the decision. Note, for example, that the neuron 2 firing rate axis is nearly orthogonal to the choice axis, meaning that neuron 2 carries almost no information about choice. **c**, Projections onto the stimulus and choice subspaces reveal the time course of information about stimulus and choice, respectively. These time courses can be seen as temporal basis functions for the single-neuron PSTHs shown in **a**. mTDR aims to recover these encoding subspaces even in the presence of additional components that take neural activity outside the plane spanned by these two axes and is not restricted to 1D subspaces.

evidence, and  $x_{\text{choice}}$  denotes the decision variable, which takes on values of  $\pm 1$ , indicating a positive or negative choice. The activity vectors  $w_{\text{stim}}$  and  $w_{\text{choice}}$  are patterns of activity across the three neurons, specifying the stimulus and choice axes (blue and red arrows in Fig. 2b), and the time-varying functions  $s_{\text{stim}}(t)$  and  $s_{\text{choice}}(t)$  are temporal profiles for the activity along stimulus and choice axes,

respectively (Fig. 2c). Noise is added to each firing rate to account for variability not due to task variables.

Although the choice and decision subspaces in this example are both 1D, the mTDR model easily generalizes to higher dimensionality and for an arbitrary number of task variables. Let  $\mathbf{Y}$  denote a neurons  $\times$  time matrix of firing rates for a single condition defined

by task variables  $\{x_1, \dots, x_p\}$ . The mTDR model decomposes population activity as:

$$\mathbf{Y} = x_1 \mathbf{W}_1 \mathbf{S}_1^\top + \dots + x_p \mathbf{W}_p \mathbf{S}_p^\top + \text{noise} \quad (2)$$

where  $\mathbf{W}_p$  is a neurons  $\times$   $r_p$  matrix whose columns span an  $r_p$ -dimensional encoding subspace for task variable  $x_p$ , and  $\mathbf{S}_p$  is a time  $\times$   $r_p$  matrix of temporal profiles that describe the time course of population activity within this subspace (Supplementary Fig. 1). This model-based formalism represents a generalization of TDR<sup>1</sup>, which allows us to identify both the number of activity patterns used to encode different variables and the time courses with which these patterns are recruited (for details, see Methods and Supplementary Note 2).

**Population coding of task variables in the PFC.** To investigate population-level coding in the PFC, we applied mTDR to neural data recorded from an area in and around the frontal eye fields (FEFs) of two monkeys performing a context-dependent decision-making task<sup>1</sup> (see Methods and ‘Experimental details’). In this task, monkeys were presented with a visual stimulus that contained colored, moving dots on each trial (Fig. 1a). A context cue (yellow square or blue cross) appeared before each trial and instructed the monkeys to attend to either the color (red versus green) or the motion (left versus right) of the dots. In the color context, the animal had to attend to color and ignore motion, making a saccade to the red (green) target if most of the dots were red (green). In the motion context, the animal had to attend to motion and ignore color (eg, making a left (right) saccade if the dot motion was left (right)). Task difficulty was controlled by varying the fraction of red versus green (and coherently moving) dots, across six levels of coherence for each stimulus dimension (Fig. 1b). After a randomized delay, the monkey was cued to indicate its decision by making a saccade to one of the two targets.

Classical approaches to analyzing data of this type involve analyzing average firing rates, or peristimulus time histograms (PSTHs), for different task conditions (conditional PSTHs), such as ‘all trials with the strong rightward motion and a rightward choice’. For this data set, the conditional PSTHs of individual neurons exhibited heterogeneous tuning to the different task variables<sup>1</sup> (Fig. 1c). This heterogeneity, and the fact that each neuron encodes several task variables, makes it difficult to obtain a clear picture of the population-level representation of task variables.

To overcome these limitations, we used mTDR to determine the dimensionality of population-level representations of the task variables. We included a regressor for each of six task variables: color strength, motion strength, context and choice, as well as two additional terms for the absolute values of color and motion strength. Absolute value terms were included because of the observation that some neurons displayed nonlinear encoding of stimuli, consistent with observations of nonlinear mixed selectivity<sup>13</sup>. The model also included a term for the condition-independent firing rate, which reflects temporal modulation not due to the task variables (for details, see Methods). To determine the dimensionality of the encoding of each task variable, we used a greedy step-wise selection method based on the Akaike information criterion<sup>21</sup> (AIC) that added dimensions based on their contribution to the model prediction performance<sup>22</sup>. We validated this approach with simulation experiments and with cross-validation on the real data, which we found to slightly underestimate dimensionality owing to the need to divide data into training and test sets (Fig. 3c).

We found that population-level representations of all task variables were at least 2D, and at least 3D in monkey A (Fig. 3 and Supplementary Table 1). Fig. 3a shows the task variable-specific components revealed by mTDR for an example neuron. The first three columns show the time course of this neuron’s activity within

the first three dimensions of the corresponding variable’s encoding space. The time courses represent the columns of the temporal component matrices  $\mathbf{S}_p$ , scaled by the levels of each of the task variables  $x_p$  (Equation 2). Thus, each trace represents the inferred contribution of each dimension to the neuron’s PSTH from the different settings of the associated task variable. The rightmost column of Fig. 3a shows the model-based estimate of the neuron’s net time-varying response to each task variable. Summing these responses gives the model-based reconstruction of the neuron’s PSTH for each task condition; this matches the neuron’s true PSTH to high accuracy (bottom). Because each neuron weights these components differently, the fitted model can account for a wide variety of conditional PSTHs (Fig. 3b). Note that the data were not temporally smoothed, and no smoothness constraints were included in the model, indicating that the smoothness of the time courses is a property of the data.

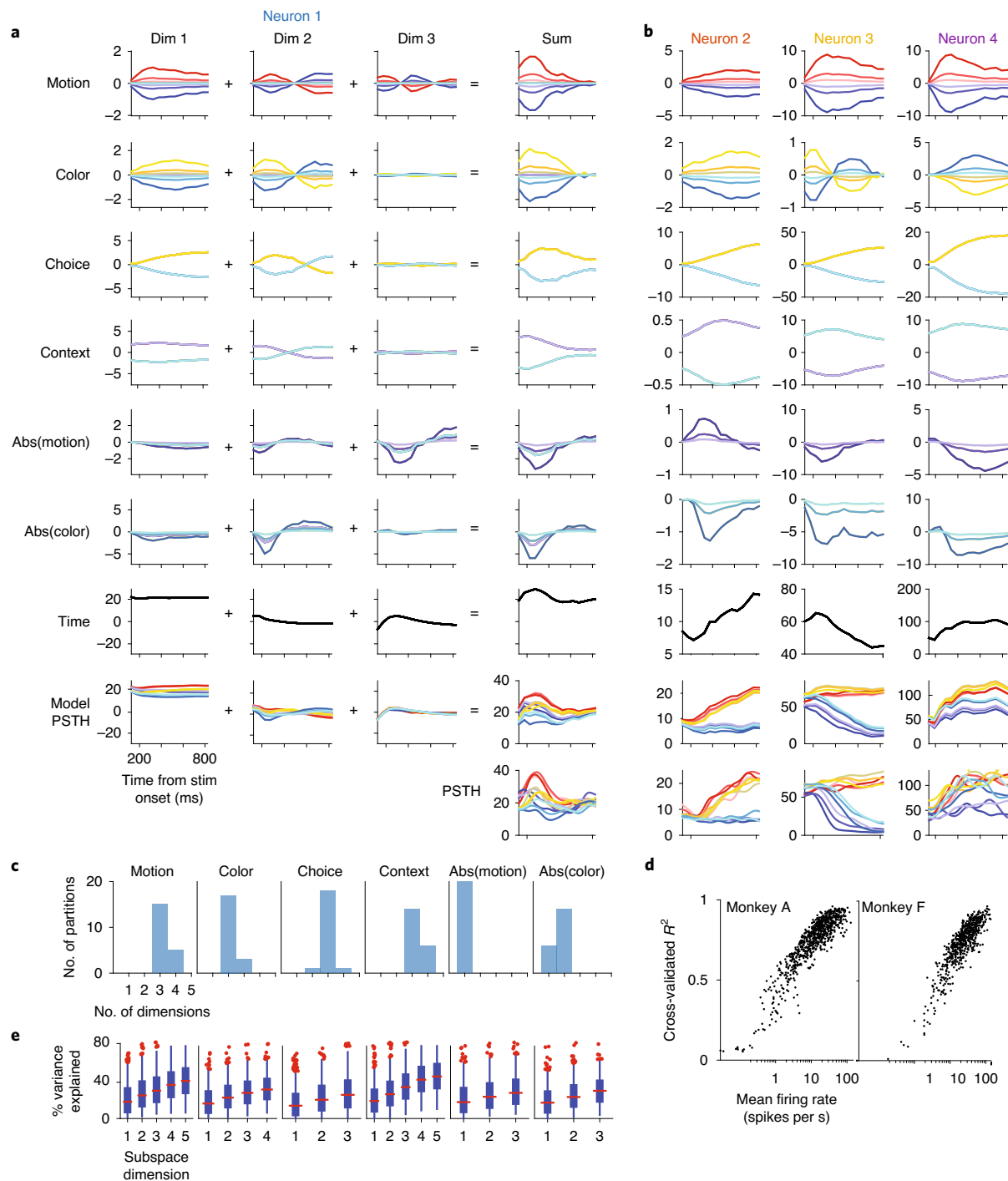
To examine whether the model was flexible enough to capture the diverse response profiles observed across the population, we calculated the  $R^2$  of the conditional PSTH for each condition using held-out data. We found that the  $R^2$  of PSTH reconstructions increased with firing rate, with some neurons achieving  $R^2$  greater than 0.9 (Fig. 3d). The dependence of  $R^2$  on firing rate likely reflects higher signal-to-noise ratio in higher firing rate neurons.

We also measured how much of the variance from held-out trials could be explained by each of the learned subspaces alone (Fig. 3e). We defined each subspace by a set of orthonormal vectors ordered by the amount of variance explained (for details, see Methods and Supplementary Note 4.1). We found that all dimensions contributed to the variance of at least some neurons but that different neurons had their variance distributed differently across components. For example, for the decomposition in Fig. 3a, dimension three of the abs(motion) axis contributes more than dimension one despite the first dimension describing most of the variance across the population. These findings verify that the mTDR model captures high-variance dimensions and that the model describes a large fraction of the variance of the PSTHs for most neurons, despite the model being relatively low dimensional.

**State-space trajectories reveal dynamic encoding.** To explore the dynamics of population-level encoding during decision formation, we examined projections of neural activity from held-out trials onto the estimated subspaces (Fig. 4a–d, Extended Data Fig. 1, Extended Data Fig. 3 and Supplementary Videos). In contrast to previous findings<sup>1,8</sup>, we found that the encoding of the stimulus variables (motion and color) was not transient but persisted throughout the recording epoch. Projections of population activity onto a single motion or color axis identified with classic TDR suggested that sensory axis projections decay rapidly after stimulus onset<sup>1</sup>. However, mTDR revealed that stimulus information persists by rotating within multidimensional motion and color subspaces (Fig. 4a,b).

More generally, we observed that, for nearly all subspaces, the neural trajectories on nearly all task conditions initially moved outward along a single axis and then began rotating in a consistent direction (Fig. 4a–d). This observation prompted us to identify the precise orientation of this initial axis and when, or if, the trajectories curved into a second dimension. We, therefore, sought a procedure that would identify an orthogonal set of axes ordered by the times at which population activity first projects onto them. The resulting method, which we call ‘sequential principal component analysis’ (seqPCA), identifies the direction the trajectories are moving and the time at which a change in direction occurs (see Methods and Supplementary Note 9). We used seqPCA to obtain an interpretable set of axes for the subspaces identified by mTDR.

Using seqPCA, we identified an orthonormal basis for each subspace, with axes that we labeled as ‘early’, ‘middle’ and ‘late’, based on the times at which they became active during the task period (Fig. 4a–d). By definition, the early axis accounted for most of the



**Fig. 3 | Model fit for monkey A. a**, Example of a neuron's fitted responses decomposed into a set of weighted basis functions (same as neuron 1 from Fig. 1c). These basis functions are shared by the whole population but are weighted differently for each neuron. Weighted basis functions are summed to form the neuron's response to each task variable. The responses for each task variable are then added together to give the model reconstructed 'model PSTHs'. The true conditional PSTHs of this neuron are shown below for comparison. **b**, Summed responses for three additional example neurons (same as neurons 2–4, from Fig. 1c), which display a diversity of dynamics. **c**, Dimension estimation based on 5x four-fold (20 estimates) cross-validation. Dimensionality is slightly smaller than estimated using all data but is tightly distributed around a single estimate. **d**,  $R^2$  of the model reconstructions for the PSTHs as a function of mean firing rate for each neuron. **e**, Percent variance explained for PSTHs of each neuron ( $n = 762$ ) by projection onto each subspace dimension. Red horizontal bars indicate the median. Box edges indicate the 25th and 75th percentiles. Whiskers indicate positions of furthest points from median not considered outliers. Red dots indicate outliers with respect to a normal distribution. Dots have been horizontally jittered to aid with visualization. Results have been averaged for each neuron over four cross-validation folds. Colors in the title text for **a** and **b** correspond to colors of markers in Fig. 5.

variance in neural trajectories during the time period immediately after stimulus onset. Variance that was not described by the early axis but emerges sometime after stimulus onset is captured by the

middle axis. The late axis accounts for activity that is not accounted for by the early and middle axes but is present as the epoch transitions from the stimulus presentation to the delay period. Projections



onto the seqPCA axes show clear times at which task variable information becomes available onto each axis (right-side panels in Fig. 4a–d). For all subspaces, we found that the early epoch is characterized by loading of the projections almost exclusively onto a single axis. In contrast, the middle and late epochs were 2D or higher.

We found that the transience of the early stimulus axes resembles that of the stimulus encodings using the TDR method<sup>1</sup>. Indeed, we found that our early axis was well correlated with the TDR axes (see Supplementary Note 10). It is, therefore, apparent that the middle and late seqPCA axes permit the stimulus information to persist. We compared projections onto the subspaces learned by mTDR with the 1D axes of TDR (see Supplementary Note 10) and found that, whereas stimulus information appeared transient for TDR, the mTDR projections were both larger and more persistent (Fig. 4e and Extended Data Fig. 4). Finally, the population-level representation of choice, context, abs(motion) and abs(color) also exhibited multidimensional structure (Fig. 4c,d and Extended Data Fig. 1). We describe this structure and discuss its consequences in subsequent sections.

**Trajectories exhibit rotational dynamics.** The projections of neural population activity onto motion, color, choice and context exhibited rotations after early-axis activity reached a peak and middle-axis activity began to increase (Fig. 4a–d). This observation is supported by the fact that the trajectories are  $\geq 2$  dimensional during this period (Fig. 4a–d, right panels). Although rotations are inherently  $\geq 2$  dimensional, the fact that we found trajectories to be  $\geq 2$  dimensional need not imply rotations. We, therefore, identified the plane of greatest rotation of the trajectories using jPCA<sup>18</sup> (Extended Data Figs. 2 and 6) and observed clear rotational structure. The two dimensions of the jPCA plane accounted for a relatively large amount of the variance for all task variables (Supplementary Fig. 1). Condition-shuffled projections yielded no apparent sequential or rotational structure (Supplementary Note 5 and Supplementary Figs. 2 and 3).

To rigorously examine the presence of rotational dynamics, we examined the angle of rotation that the trajectories traversed from the beginning of the middle epoch to the end of stimulus viewing (Fig. 4f). We reasoned that, for trajectories to be consistent with rotational dynamics, they would have to have monotonically changing angles of rotation. We compared the angle of rotation to samples from the null distribution corresponding to the maximum entropy distribution with the same second-order moments as the data<sup>23</sup> (Fig. 4f; for details, see Supplementary Note 5). We found evidence

for rotational dynamics in motion, color, choice and context subspaces, although rotations were less consistent with the trajectories of the color encoding for monkey F (Extended Data Figs. 3 and 6). These results indicate that rotational dynamics are not trivially present in these data and that we observed them in most of the linear subspaces examined.

Projections onto the subspaces for the absolute values of motion and color (abs(motion) and abs(color)) were qualitatively different from those of the linear terms (Extended Data Fig. 1). Although they clearly encoded the absolute values of the stimuli, evidence for rotational dynamics was not significant (Fig. 4f and Extended Data Figs. 1 and 2).

**Characterizing neural selectivity across subspaces.** We used the mTDR model and seqPCA to examine the how tuning properties of these cells changed over time and the relationships in tuning between cells. Individual neurons exhibited complex mixtures of early, middle and late responses (Fig. 5a). Although the population tuning of some task variables (abs(motion) and abs(color)) were dominated by the early response, none of the task variables was found to display clustering but a continuous distribution of tuning across all three seqPCA axes. Late axes tended to explain less of the population variance, especially for color, choice and abs(motion), but were responsible for explaining most of the variance for at least some neurons. This is evident from the dearth of units in the ‘late’ vertices of Fig. 5a.

Also notable was the low density of cells near the early/late axis (that is, the left arm of the ternary plots in Fig. 5a). A low density of cells near the early/late edge of the plot indicates that there are few cells that encode a task variable at the beginning and end of stimulus viewing but lose sensitivity to a task variable in the middle of stimulus viewing. The lack of cells with a gap in the timing of encoding implies that individual cells tend to encode task variables in continuous epochs, even if only transiently.

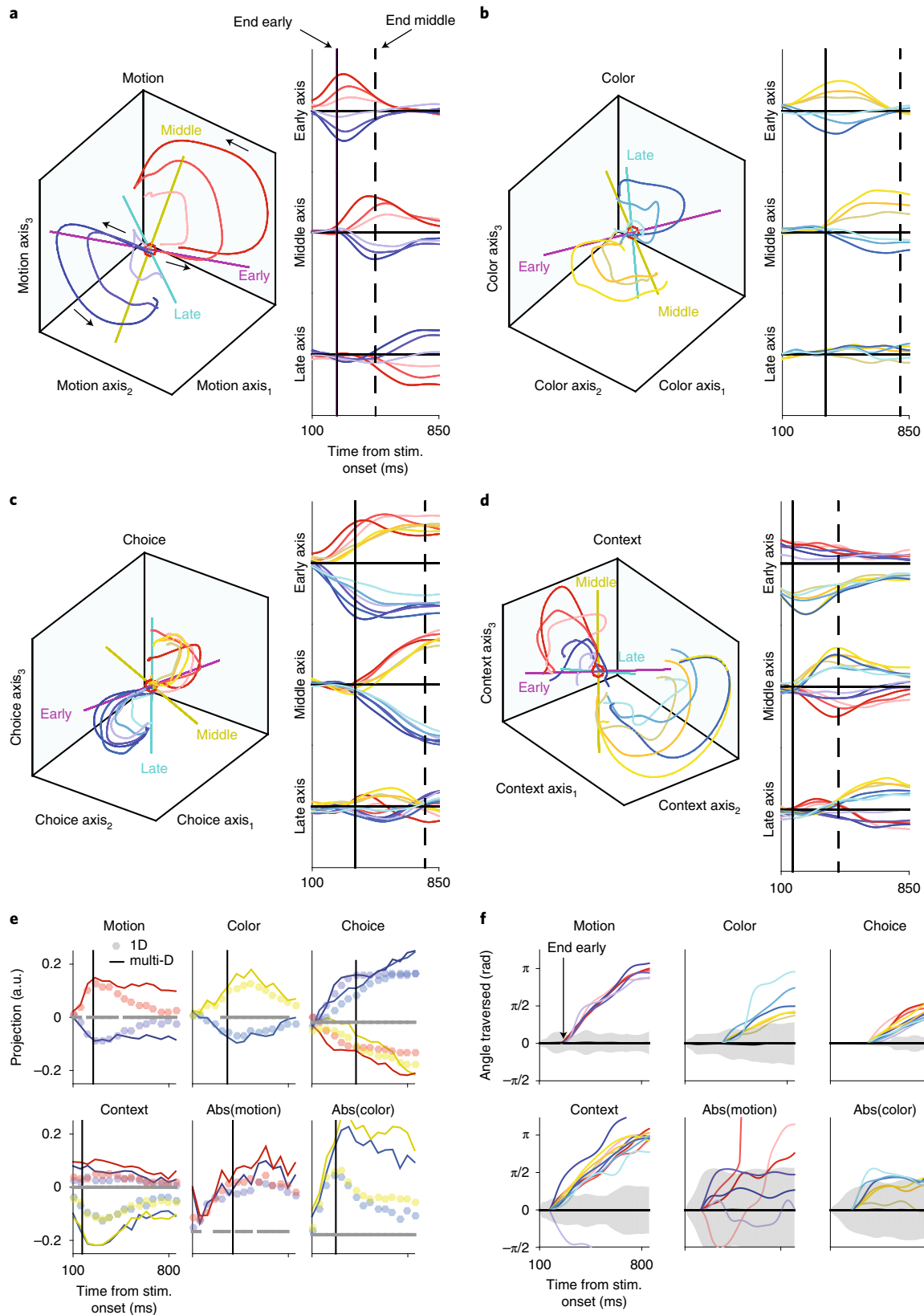
The subspaces identified by mTDR for motion, color and choice were positively correlated (Fig. 5b,c). More specifically, the weights defining the motion and color bases were correlated with the choice weights but not with one another, indicating that motion and color representations both contributed to the choice encoding but that there was little cross-stimulus interference between representations.

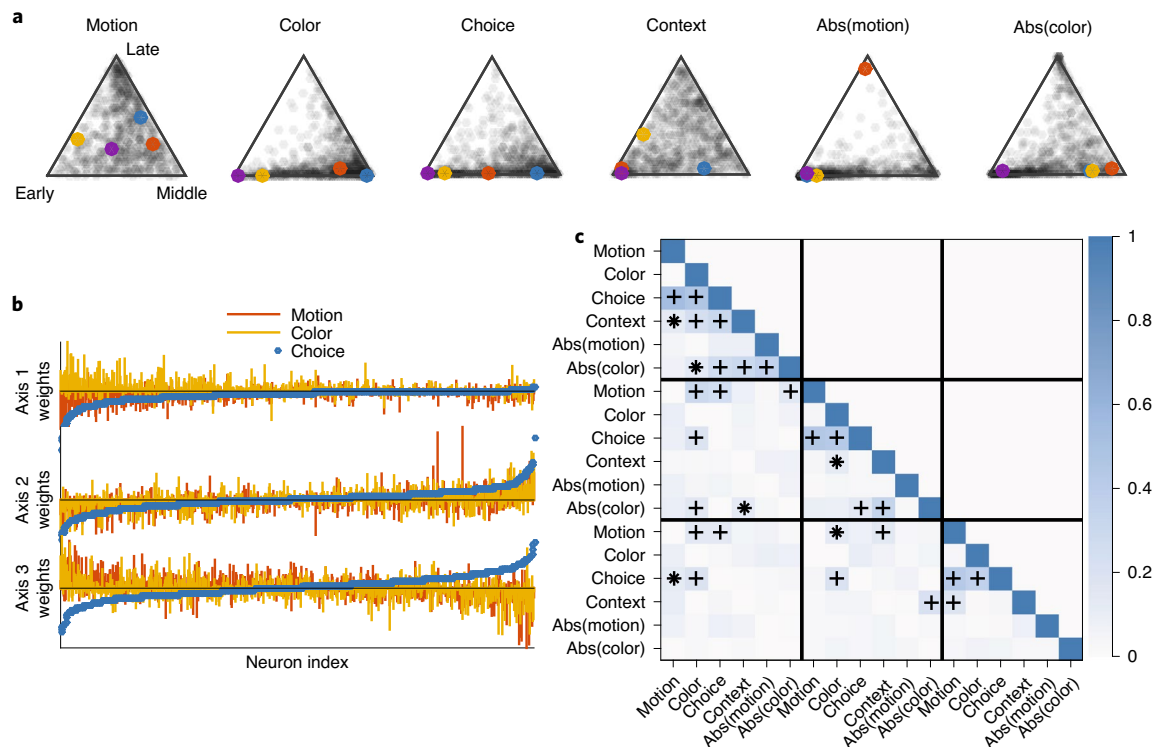
**Accurate stimulus decoding coincided with the onset of rotational dynamics.** The mTDR model provides a framework for

**Fig. 4 | Projections of population PSTHs onto latent encoding subspaces.** Projections onto the first, second and third principal axes of the (a) motion, (b) color, (c) choice and (d) context subspaces. Projections of population PSTHs onto latent encoding subspaces. Motion, color and context subspaces have been orthogonalized with respect to the first dimension of the choice subspace. The choice subspace has been orthogonalized with respect to the context subspace. The context subspace has also been orthogonalized with respect to the motion and color subspaces. Details of orthogonalization are presented in Supplementary Math Note 4.2. Color conventions are the same as those described in Fig. 1. Red dots indicate the origin. Projected PSTHs made from held-out data not used during parameter estimation. **a**, Projections of PSTHs onto the motion subspace, sorted by motion coherence and averaged over color coherence for trials where the motion stimulus was the active context. **b**, Projections onto the color subspace sorted by color coherence and averaged over motion coherence for trials where the color stimulus was the active context. **c**, Projections onto the choice subspace. Motion-context trials are displayed with the same sorting and color conventions as displayed in **a**. Color-context trials are displayed with the same sorting and color conventions as displayed in **b**. Only correct trials are displayed. **d**, Projections onto the context subspace using the same conventions as displayed in **c**. Only correct trials are displayed. Colored axes in 3D plots indicate seqPCA axes. Solid vertical lines accompanying time traces indicate the time points where middle-axis variance starts to increase. Dashed vertical lines indicate the time points where late-axis variance starts to increase. Units of the ordinate are arbitrary, but all time-trace axes are on the same scale. PSTHs were generated with  $\sim 13$ -ms time bins and smoothed with a Gaussian window with s.d. of  $\sim 50$  ms. **e**, Median encoding strength of pseudo-trials onto the first three encoding axes of mTDR compared with the 1D subspace estimated by the max-norm method used by Mante et al.<sup>1</sup> (for details, see Supplementary Math Note 10). For clarity, only trials with the strongest stimulus strengths are shown. Gray bars at  $y = 0$  indicate time points when the mTDR projections had significantly stronger encoding across all stimulus levels than the 1D projections (left-tailed Wilcoxon signed-rank test; positive false discovery rate<sup>45</sup> controlled at 0.01). Multidimensional mTDR projections are larger than 1D projections at nearly all times for all task variables. **f**, Rotation angle traversed through rotational projection using jPCA. The angle was calculated starting from the time when the projection transitions between the early and middle epochs. Coherent traversal across stimulus strengths that is consistent and monotonically increasing is an indication of rotation. Shaded areas are 95% confidence regions calculated using a maximum entropy method<sup>23</sup> ( $n = 100$  samples) under the null hypothesis of no population structure other than the empirical means and covariances across time, neurons and task conditions.

population decoding by maximum likelihood (see Supplementary Note 6). This framework allows our decoding analysis to be consistent with the results of dimensionality reduction. We can, therefore, investigate how and when the features of the low-dimensional trajectories translate into putatively perceived stimuli and behavior

and whether these features might be read out by downstream populations. Although decoding of task variables does not imply a causal role for the encoded variables in PFC function, decoding analysis can provide a clearer picture of the dynamics and fidelity of task variable encoding.





**Fig. 5 | Distribution of variance within and between subspaces.** **a**, Proportion of variance among seqPCA axes. Each marker corresponds to one neuron. The position of each neuron indicates that the distribution of variance from PSTHs across corresponding early, middle, and late axes (e.g., a point that lies closer to the 'early' vertex of the motion plot) has more of its motion-specific variance explained by the early axis, whereas a point in the middle of the simplex has variance equally distributed across all axes. Darker regions indicate higher density of points. Colored dots correspond to cells displayed in Fig. 3. **b**, Weights of the top (in terms of variance explained) three axes for all cells for motion, color and choice subspaces. Cell indexes are sorted according to the choice weights from most positive to most negative. **c**, Magnitude of the Pearson correlation between the top three subspace axes. The magnitude is used because the axes are only identifiable up to a sign. Markers indicate significant correlations controlled by the positive false discovery rate<sup>45</sup> (\* $Q < 0.01$ , \*\* $Q < 0.01$ ). Null distribution is based on the positive half-Gaussian with zero mean and s.d.  $\sigma_0 = 1/n$ , where  $n = 762$  is the number of neurons. Significant correlations are most consistent between color-choice and motion-choice pairs. All tests were one sided.

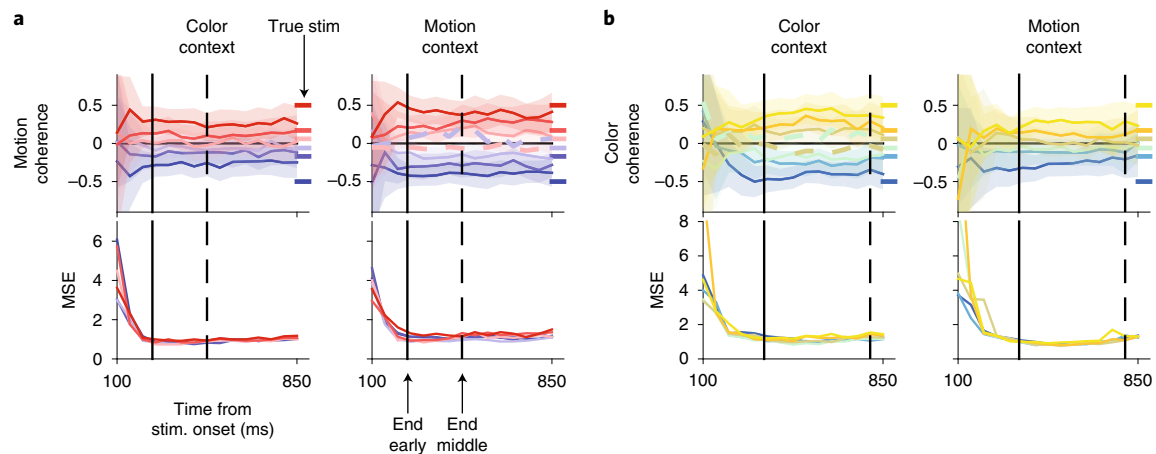
For decoding analyses with monkey A, we used four-fold cross-validation in which the held-out trials were used to produce 100 pseudo-samples (with replacement) for decoding (for monkey F, we used two-fold cross-validation with similar results). The resulting decoded values were averaged over pseudo-samples and cross-validation folds.

Stimuli could be accurately decoded within  $\approx 150$  ms of stimulus onset for the motion stimulus and within  $\approx 200$  ms for the color stimulus, roughly corresponding to the time of transition between the early and middle seqPCA axes (Fig. 6a). The values of the decoded stimuli were constant by the start of the middle epoch for both contexts, and the variance of the decoding decreased dramatically up to this time (Fig. 6b). Thus, the change in population dynamics (early-to-middle transition) within the stimulus subspace was consistent with decoding accuracy and stability. The decoded values are slightly biased toward zero in the irrelevant context, suggesting some gating of information across contexts.

To examine the content of stimulus encoding, it is often informative to examine error trials. We, therefore, examined the decoded stimulus for error trials using the weakest stimulus strengths (dashed lines, Fig. 6c and Extended Data Fig. 7c). For these data, only the weakest stimulus strengths had enough error trials to provide reliable statistical analysis<sup>1</sup>. For monkey A, we found that the decoded stimulus values on error trials were similar to correct-trial decoding but were opposite in sign, suggesting that the origin of errors was (on average) an incorrect percept.

**Choice decoding.** To examine how the decision might have evolved over the course of stimulus viewing, we next studied how and when decision information became available in the PFC and the dynamics of choice encoding. Although we encoded stimuli as continuous-valued variables in our model, choice was encoded as binary. Therefore, we examined the log likelihood ratio (LLR) (for details, see Supplementary Math Note 6.3) over time of pseudo-trials sampled from held-out data between the likelihood of a preferred versus an anti-preferred choice (Fig. 7a).

The magnitude of the LLRs increased monotonically over time, indicating an increasing strength of the decision signal. However, the magnitude of the LLR did not differ strongly with respect to context, direction of decision, stimulus strength or whether the trials were correct or error trials (dashed lines, Fig. 7a). By transforming the LLRs into decision probabilities (Fig. 7b; see Supplementary Note 6.3), we could examine a moment-by-moment probability of the animal's choice and estimate when the decisions were unequivocal. We found that the choices could be discriminated with better than 95% accuracy as early as 300–350 ms after stimulus onset (Fig. 7b). This timing corresponded to the time of transition between the early and middle seqPCA axes for choice. Similar results were observed for monkey F (Extended Data Fig. 8). These results suggest that, on average, the animals had made their decisions well before stimulus offset regardless of the stimulus coherence and that decisions were coincident with a change in dynamics from linear to rotational within the choice subspace.



**Fig. 6 | Instantaneous decoding of stimulus for monkey A.** **a**, Top: decoded motion coherence by mTDR model in both contexts. Bottom: mean squared error (MSE) over time of motion coherence decoding across stimulus levels and context. MSE decreases precipitously and then stabilizes around the time of the first transition. **b**, Same as **a** for color coherence decoding. Color conventions are the same as in Fig. 4. Shaded regions indicate 50% confidence intervals. Dashed lines indicate error trials from the corresponding context for the lowest stimulus strengths. 100 pseudo-trials for each of four-fold cross-validation ( $n = 400$ ) used for all analyses. Solid vertical lines indicate the time of early-middle axis transition for the corresponding stimulus subspace projections. Dashed vertical lines indicate the time of middle-late transition.

Restricting the choice subspace to only the early, middle or late axes, the LLRs displayed the same invariance to choice, stimulus strength, context and correct/error trial identity as the full model. For both monkeys, the early axis provided most of the available information about the decision, and early axis decoding alone is nearly as accurate as the full model (Fig. 7d and Extended Data Fig. 8d). However, the middle and late axes also displayed information about the choice later during stimulus viewing.

Because we can decode the animals' decisions with the early axis alone, it would seem as though the middle and late axis information is redundant, and it is unclear what the purpose of these axes is. Similar multidimensional encoding of decision has been observed previously in the premotor cortex<sup>24</sup>.

**Context decoding.** We examined the context signal using the same LLR method as our analysis of choice (Extended Data Figs. 9a and 10a). The context evidence did not differ strongly across decision, stimulus strength or whether the animal provided a correct or incorrect response. Transforming the LLRs into a probability of the perceived context (Extended Data Figs. 9b and 10b) showed that the correct context could be identified for both monkeys on most pseudo-trials from the first time point, which is consistent with the fact that the context cue was presented 650 ms before stimulus onset<sup>1</sup>. These patterns hold for LLRs of error pseudo-trials as well as for decoding restricted to only the early, middle or late subspaces (Extended Data Figs. 9c,d and 10c,d). These findings demonstrate that accurate context information was available in the PFC for the vast majority of both correct and error trials, suggesting that confusion about context was not a significant source of errors.

## Discussion

Our analyses have shown that the PFC encodes individual task variables in distinct multidimensional subspaces that capture dynamic changes in representation over time. The population activity patterns representing each task variable tended to exhibit a stereotyped 1D linear phase, followed by a rotational phase. Our ability to make these observations relied on a new method for dimensionality reduction based on a probabilistic low-rank model of the data.

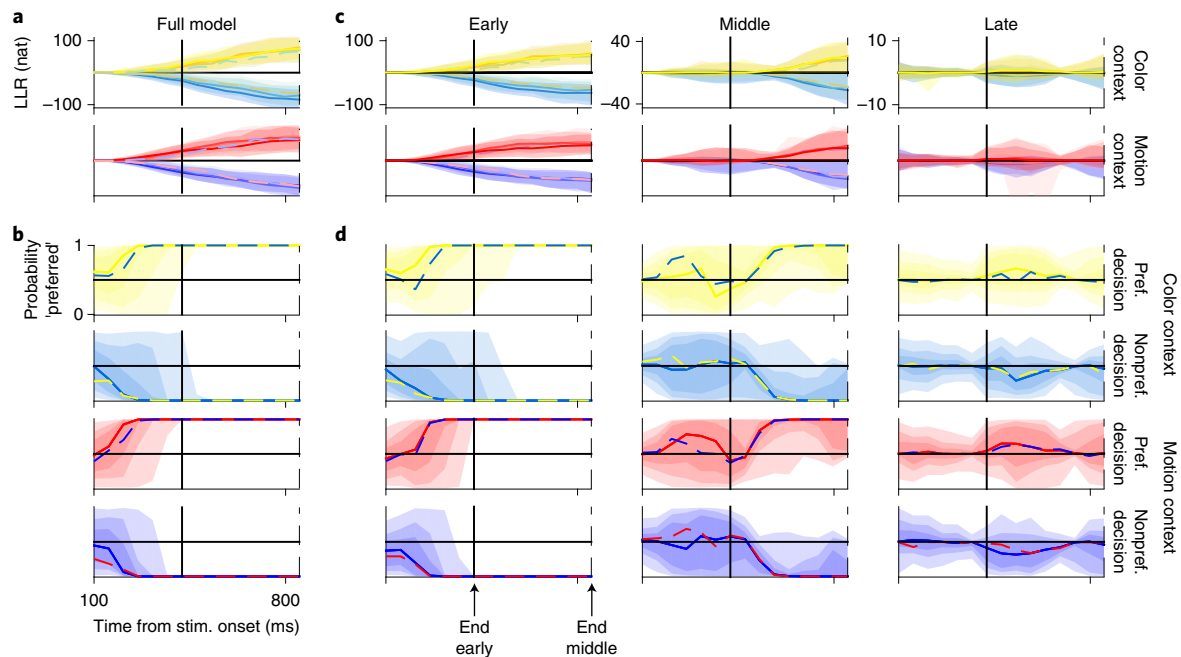
We found that the dynamic nature of encodings in the PFC requires multiple dimensions of neural population activity for accurate characterization. In particular, only multidimensional

encoding, as opposed to 1D encoding, captures the persistence of stimulus information in the PFC throughout the stimulus-viewing epoch (Fig. 4a,b). This finding complements the original report of these data<sup>1</sup>, suggesting that previously reported transient stimulus encoding in the PFC is consistent only with the early encoding axis (Fig. 4e). Although the mechanisms of selection and integration proposed in Mante et al.<sup>1</sup> are consistent with the early evolution of trajectories revealed here, by themselves they cannot readily explain the following rotational dynamics. Our observations resemble multidimensional stimulus coding that mixes transient and persistent components<sup>24</sup> as well as population code 'morphing'<sup>16</sup>, where the optimal weights for decoding from population activity change over time, although the results shown here are on a time scale that is nearly an order of magnitude faster than previously reported.

Although we validated our method for identifying the 'true' dimensionality of the data using simulation experiments, it is unclear whether the dimensionality would differ under different experimental conditions. Specifically, the dimensionalities we learned are likely to be influenced by a variety of factors<sup>25</sup>, including the sample size, the fraction of neurons observed, the intrinsic model dynamics and the task complexity. Some of these factors might explain the differences in dimensionality between the two animals in the present study, where the dimensionalities of monkey F were lower than monkey A in correspondence with smaller sample sizes and fewer recorded cells. However, we emphasize that, during the early encoding, nearly all trajectories are 1D and only afterward are  $\geq 2$ D. The fact that trajectories are multidimensional after the first transition might be a reflection of their rotational nature during this epoch. Rotations are inherently  $\geq 2$ D because they require both sine and cosine parts for each axis of rotation.

The mTDR method is distinct from unsupervised dimensionality reduction methods such as PCA or factor analysis in that it uses information about the experimental variables of interest on each trial. The method is also distinct from previously proposed supervised methods<sup>1,26–29</sup> in its use of an explicit generative model to describe the transformation from task variables to neural activity patterns. This distinction not only allows us to make predictions of population responses to experimental contingencies not observed in the data (something not possible for methods based on the conditional PSTHs like dPCA without model-based interpolations<sup>27</sup>), but it also allows us to apply the tools of probabilistic modeling and





**Fig. 7 | Instantaneous decoding of choice.** **a**, LLRs for monkey A in favor of a preferred choice using single pseudo-trials from color-context (gold and blue, sorted by color coherence) and motion-context (red and violet, sorted by motion coherence) trials. Shaded regions indicate 95% quantile intervals for each stimulus strength. Solid lines indicate the median of correct trials. Dashed lines indicate median of error trials. **b**, Probability of a preferred choice based on corresponding LLRs combined over all stimulus strengths (for details, see Supplementary Note 6.3). Solid lines indicate the median of correct trials. Dashed lines indicate the median of error trials. Shaded regions indicate quantile coverage intervals of correct trials (light-to-dark: 95%, 75% and 50%). Color conventions are the same as in Fig. 4. 100 pseudo-trials for each of four-fold cross-validation folds used for all analyses. **c**, LLRs in favor of a preferred choice where the choice subspace has been restricted to only the early, middle or late axes. **d**, Probability of a preferred choice based on LLRs from **c**.

inference to estimate both the model parameters and the dimensionality of the encoding.

Our approach (Equation 5) is similar to that used by Mante et al.<sup>1</sup> and to other methods based on linear regression models (see examples in refs. 30–33). However, our model is distinguished by its explicit specification of low-rank regression parameters and neuron-specific noise variances. Future improvements to our model may incorporate nonlinear mapping of stimuli to neuronal responses<sup>30</sup>, noise correlations between simultaneously recorded neurons and support for variable trial lengths.

Much theoretical development has rested on the notion that single-neuron spike rates map onto an evidence accumulator, but recent evidence in the FOF, a rodent analogue of the FEF, has challenged this view<sup>6</sup>, suggesting that this region can be better described as maintaining a running motor plan (saccade for FEF and orienting for FOF) based on the evidence accumulated<sup>6,7</sup>. Although our analysis does not aim to suggest a causal role of FEF, the results of the present study could be interpreted as supporting this view, where the early dynamics represent an evolving decision and the rotational dynamics indicate an evolving motor plan, but more work is needed to determine the precise role of FEF, and PFC more generally.

**Functional significance of sequential subspaces.** Our analysis revealed temporally segregated dynamics with early-axis activity transitioning to middle and late axes, with rotations dominating at around 200–400 ms after stimulus onset (Fig. 4 and Extended Data Fig. 2). The temporal separation of the early/linear and rotational subspaces suggests that these are subspaces within which distinct computations are evolving<sup>18,20,34</sup> or have independent sets of downstream targets<sup>19</sup>.

With the present data, we can only speculate about what the nature of these different computations must be, but the present analysis indicates the possibility that the early epochs are concomitant

with the temporal window that decision-making is performed. For example, the timing of transition between early and middle epochs is consistent with the timing of accurate decoding of the animals' decisions from single pseudo-trials (Fig. 7 and Extended Data Fig. 8). This timeframe is consistent with the timing of saturation of the chronometric curve for the traditional kinematogram task<sup>35–37</sup>, with the distribution of step times in the stepping model of evidence accumulation<sup>38</sup> and with early weighting of evidence in visual discrimination tasks<sup>39</sup>. This evidence suggests that the transition from linear to rotational dynamics is a correlate of decision commitment.

A similar sequence of dynamics has been observed in population activity from the premotor cortex that corresponds to distinct 'preparatory' and 'movement' epochs<sup>18–20,34</sup>. However, in these studies, the transitions in dynamics could be linked directly to an overt action (arm movement), whereas our animals would not have made an overt action (saccade to target) until 300–1,500 ms after the end of our analysis window<sup>1</sup>. Therefore, if the animal has made its decision, then it would have done so only covertly.

These distinctions, however, might be superficial. The qualitative features of our results reflect those in the motor cortex strikingly well<sup>18–20,34</sup>, suggesting that common mechanisms might be at work in both motor execution and decision-making. Indeed, FEF is defined as a region that elicits eye movement under stimulation<sup>40,41</sup> and has been implicated as a region important for visual decision-making<sup>4,4,6,9–12,14,15</sup>, oculomotor planning<sup>42</sup> and covert visuospatial attention<sup>43,44</sup>. Thus, we may think of FEF as itself a premotor area responsible for visuospatial attention and motor planning associated with decision-making<sup>4,6,7,9</sup>. The dynamic transitions in our analysis could be interpreted as signaling decision commitment<sup>6</sup>, or as signaling a covert action (saccade preparation), in analogy with the transitions observed between preparatory and movement periods seen in the premotor cortex<sup>18,19,34</sup>. Single-trial population analysis and analysis of delay and saccade epochs of these



experiments might shed light on how the dynamics we observe reflect the animals' decisions.

Some subspaces lacked a distinct late component (eg, color and choice subspaces for monkey A; Fig. 4a,c). However, it is possible that the middle seqPC for some task variables served a similar role as the late seqPC for others, preparing the network for a new set of targets or storing the memory of the stimuli as persistent activity over the course of the delay period. The number of seqPCs needed to describe the population activity might reflect the rate that trajectories rotate into new encoding directions and therefore correspond to a quantitative rather than a qualitative difference in encoding. Future work should be aimed at identifying the significance of the dimensionality of the encoding relative to the sequential dynamics.

The nature of dynamic encoding for the context variable remains mysterious. Context encoding for both animals displayed clear and consistent dynamics (Fig. 4d and Extended Data Fig. 3d), including rotations (Fig. 4f). Furthermore, although most of the predictive capacity of the context encoding lies in the early subspace (Extended Data Figs. 9 and 10), where context is encoded throughout the stimulus viewing period, context encoding at the single-neuron level is broadly distributed across the early, middle and late axes (Fig. 5b and Extended Data Fig. 5), indicating that some neurons do not encode context until well after stimulus onset. Further work is needed to determine what, if any, function these dynamics serve in decision-making and memory. The uniqueness of these phenomena to the present setting is an active area of research.

**Differences in encoding between animals.** The two monkeys in this study displayed similar, but not identical, encoding properties. For example, the encoding trajectories for motion were similar (Fig. 4a and Extended Data 3a), but we found obvious differences between the encoding trajectories for color (Fig. 4b versus Extended Data Fig. 3b). For monkey A, the color trajectories closely resembled the trajectories for motion (Fig. 4a,b), whereas, for monkey F, the color trajectories did not display obvious rotations (Extended Data Fig. 3b). Choice and context trajectories in monkey F appear to be similar to those of monkey A (Fig. 4e,g and Extended Data Fig. 3e,g) but displayed less pronounced rotations (Extended Data Figs. 2 and 6). These across-animal differences verify that rotational dynamics are not trivially present in these data, and, although it is unclear precisely what function they serve, they are a potentially important feature of encoding in the PFC.

Although the reason for differing dynamics between the color encoding for monkey F and the other stimulus encodings is unclear, we do have some behavioral clues as to its effect. For example, the color-context psychometric curve for monkey F was somewhat more shallow than for motion as well as for both motion and color for monkey A (Extended Data Fig. 2d in ref. <sup>1</sup>), and motion served as more of a distraction during the color task for monkey F than for monkey A, suggesting that the color discrimination task was more difficult for monkey F. Furthermore, we found that the decoding accuracy for color in monkey F was worse than for monkey A (Fig. 6 and Extended Data Fig. 7), suggesting that color information was more poorly represented in the PFC for monkey F. Although not definitive, together these results suggest that monkey F might have had more difficulty with color coherence perception and that indistinct encoding features are a correlate of perceptual uncertainty. Future experiments could be aimed at examining this hypothesis.

**Decoding of error trials suggests sources of errors.** There are three ways that the animals might commit an error: the animal perceived the wrong stimulus (eg, perceived left motion on a right-motion trial); the animal was confused about the context (eg, made its decision using the color information in the motion context); or the animal made a random choice (that is, a 'lapse' trial). The results of this analysis for monkey A at the weakest stimulus strengths indicate

that the animal perceived the wrong stimulus. The decoded context, on average, was the correct context (Extended Data Fig. 9), ruling out whether the animal was confused about which stimulus it was supposed to attend. Lapse errors are also unlikely to contribute significantly to the animal's behavior. The psycho-physical curves of monkey A suggest a small lapse rate, if any<sup>1</sup>, and stimulus decoding indicates that the perceived relevant stimulus on error trials was of the opposite sign as the stimulus that was presented (Fig. 6c). Together, these observations indicate that most error trials are based on an incorrect perception of the relevant stimulus. A more direct trial-by-trial analysis of simultaneously recorded neurons would be useful in confirming this hypothesis.

The results for monkey F are more difficult to interpret. The decoded stimuli for error trials appear to be close to zero, indicating an ambiguous stimulus (Extended Data Fig. 7b). Furthermore, the choice signal on error trials appears to be present earlier, on average, than on correct trials and is present on some trials as early as the first time point (Extended Data Fig. 8b), suggesting that the animal might have made its decision before even viewing the stimulus and that lapses were a significant source of errors for monkey F.

Given the present data, it might be impossible to distinguish the neural correlates of decision-making from those of planning for the eventual saccade. Recent work has shown that there might be independent cortical signals for evidence accumulation and decision commitment in other cortical areas<sup>39</sup>. It might be unlikely, using these data, to distinguish between a deliberate effort to make a stimulus discrimination and the formation of a motor plan<sup>31</sup>.

Nevertheless, the results presented here demonstrate the utility of mTDR for the analysis of neuronal population data and provide a description of PFC dynamics that should serve as important constraints on future models of the mechanisms of PFC function.

## Online content

Any methods, additional references, Nature Research reporting summaries, source data, extended data, supplementary information, acknowledgements, peer review information; details of author contributions and competing interests; and statements of data and code availability are available at <https://doi.org/10.1038/s41593-020-0696-5>.

Received: 6 August 2018; Accepted: 21 July 2020;

Published online: 05 October 2020

## References

- Mante, V., Sussillo, D., Shenoy, K. V. & Newsome, W. T. Context-dependent computation by recurrent dynamics in prefrontal cortex. *Nature* **503**, 78–84 (2013).
- Raposo, D., Kaufman, M. T. & Churchland, A. K. A category-free neural population supports evolving demands during decision-making. *Nat. Neurosci.* **17**, 1784–1792 (2014).
- Roy, J. E., Buschman, T. J. & Miller, E. K. PFC neurons reflect categorical decisions about ambiguous stimuli. *J. Cogn. Neurosci.* **26**, 1283–1291 (2014).
- Siegel, M., Buschman, T. J. & Miller, E. K. Cortical information flow during flexible sensorimotor decisions. *Science* **348**, 1352–1355 (2015).
- Katz, L. N., Yates, J. L., Pillow, J. W. & Huk, A. C. Dissociated functional significance of decision-related activity in the primate dorsal stream. *Nature* **535**, 285–288 (2016).
- Hanks, T. D. et al. Distinct relationships of parietal and prefrontal cortices to evidence accumulation. *Nature* **520**, 220–223 (2015).
- Erllich, J. C., Brunton, B. W., Duan, C. A., Hanks, T. D. & Brody, C. D. Distinct effects of prefrontal and parietal cortex inactivations on an accumulation of evidence task in the rat. *eLife* **4**, e05457 (2015).
- Goard, M. J., Pho, G. N., Woodson, J. & Sur, M. Distinct roles of visual, parietal, and frontal motor cortices in memory-guided sensorimotor decisions. *eLife* **5**, e13764 (2016).
- Bruce, C. J. & Goldberg, M. E. Primate frontal eye fields. I. Single neurons discharging before saccades. *J. Neurophysiol.* **53**, 603–635 (1985).
- Kim, J. N. & Shadlen, M. N. Neural correlates of a decision in the dorsolateral prefrontal cortex of the macaque. *Nat. Neurosci.* **2**, 176–185 (1999).

11. Ding, L. & Gold, J. I. Neural correlates of perceptual decision making before, during, and after decision commitment in monkey frontal eye field. *Cereb. Cortex* **22**, 1052–1067 (2011).
12. Stokes, M. G. et al. Dynamic coding for cognitive control in prefrontal cortex. *Neuron* **78**, 364–375 (2013).
13. Rigotti, M. et al. The importance of mixed selectivity in complex cognitive tasks. *Nature* **497**, 585–590 (2013).
14. Purcell, B. A. et al. Neurally constrained modeling of perceptual decision making. *Psychol. Rev.* **117**, 1113–1143 (2010).
15. Heitz, R. P. & Schall, J. D. Neural mechanisms of speed-accuracy tradeoff. *Neuron* **76**, 616–628 (2012).
16. Parthasarathy, A. et al. Mixed selectivity morphs population codes in prefrontal cortex. *Nat. Neurosci.* **20**, 1770–1779 (2017).
17. Beck, J. M. et al. Probabilistic population codes for Bayesian decision making. *Neuron* **60**, 1142–1152 (2008).
18. Churchland, M. M. et al. Neural population dynamics during reaching. *Nature* **487**, 51–56 (2012).
19. Kaufman, M. T., Churchland, M. M., Ryu, S. I. & Shenoy, K. V. Cortical activity in the null space: permitting preparation without movement. *Nat. Neurosci.* **17**, 440–448 (2014).
20. Lara, A. H., Cunningham, J. P. & Churchland, M. M. Different population dynamics in the supplementary motor area and motor cortex during reaching. *Nat. Commun.* **9**, 2754 (2018).
21. Akaike, H. A new look at the statistical model identification. *IEEE Transactions on Automatic Control* **19**, 716–723 (1974).
22. Aoi, M. & Pillow, J. W. Model-based targeted dimensionality reduction for neuronal population data. *Adv. Neural Inform. Process. Syst.* **31**, 6690–6699 (2018).
23. Elsayed, G. F. & Cunningham, J. P. Structure in neural population recordings: an expected byproduct of simpler phenomena? *Nat. Neurosci.* **20**, 1310–1318 (2017).
24. Rossi-Pool, R. et al. Decoding a decision process in the neuronal population of dorsal premotor cortex. *Neuron* **96**, 1432–1446 (2017).
25. Williamson, R. C. et al. Scaling properties of dimensionality reduction for neural populations and network models. *PLoS Comput. Biol.* **12**, e1005141 (2016).
26. Cunningham, J. P. & Byron, M. Y. Dimensionality reduction for large-scale neural recordings. *Nature Neurosci.* **17**, 1500–1509 (2014).
27. Kobak, D. et al. Demixed principal component analysis of neural population data. *eLife* **5**, e10989 (2016).
28. Machens, C. K. Demixing population activity in higher cortical areas. *Front. Comput. Neurosci.* **4**, 126 (2010).
29. Machens, C. K., Romo, R. & Brody, C. D. Functional, but not anatomical, separation of ‘what’ and ‘when’ in prefrontal cortex. *J. Neurosci.* **30**, 350–360 (2010).
30. Romo, R., Brody, C. D., Hernández, A. & Lemus, L. Neuronal correlates of parametric working memory in the prefrontal cortex. *Nature* **399**, 470–473 (1999).
31. Hernández, A., Zainos, A. & Romo, R. Temporal evolution of a decision-making process in medial premotor cortex. *Neuron* **33**, 959–972 (2002).
32. Romo, R., Hernández, A., Zainos, A., Lemus, L. & Brody, C. D. Neuronal correlates of decision-making in secondary somatosensory cortex. *Nat. Neurosci.* **5**, 1217–1225 (2002).
33. Romo, R., Hernández, A. & Zainos, A. Neuronal correlates of a perceptual decision in ventral premotor cortex. *Neuron* **41**, 165–173 (2004).
34. Elsayed, G. F., Lara, A. H., Kaufman, M. T., Churchland, M. M. & Cunningham, J. P. Reorganization between preparatory and movement population responses in motor cortex. *Nat. Commun.* **7**, 13239 (2016).
35. Gold, J. I. & Shadlen, M. N. The influence of behavioral context on the representation of a perceptual decision in developing oculomotor commands. *J. Neurosci.* **23**, 632–651 (2003).
36. Kiani, R., Hanks, T. D. & Shadlen, M. N. Bounded integration in parietal cortex underlies decisions even when viewing duration is dictated by the environment. *J. Neurosci.* **28**, 3017–3029 (2008).
37. Meister, M. L. R., Hennig, J. A. & Huk, A. C. Signal multiplexing and single-neuron computations in lateral intraparietal area during decision-making. *J. Neurosci.* **33**, 2254–2267 (2013).
38. Latimer, K. W., Yates, J. L., Meister, M. L. R., Huk, A. C. & Pillow, J. W. Single-trial spike trains in parietal cortex reveal discrete steps during decision-making. *Science* **349**, 184–187 (2015).
39. Yates, J. L., Park, I. M., Katz, L. N., Pillow, J. W. & Huk, A. C. Functional dissection of signal and noise in mt and lip during decision-making. *Nat. Neurosci.* **20**, 1285–1292 (2017).
40. Vernet, M., Quentin, R., Chanes, L., Mitsumasa, A. & Valero-Cabré, A. Frontal eye field, where art thou? anatomy, function, and non-invasive manipulation of frontal regions involved in eye movements and associated cognitive operations. *Front. Integr. Neurosci.* **8**, 1–24 (2014).
41. Bruce, C. J., Goldberg, M. E., Bushnell, M. C. & Stanton, G. B. Primate frontal eye fields. II. Physiological and anatomical correlates of electrically evoked eye movements. *J. Neurophysiol.* **54**, 714–734 (1985).
42. Rizzolatti, G., Riggio, L., Dascola, I. & Umiltà, C. Reorienting attention across the horizontal and vertical meridians: evidence in favor of a premotor theory of attention. *Neuropsychologia* **25**, 31–40 (1987).
43. Thompson, K. G., Biscoe, K. L. & Sato, T. R. Neuronal basis of covert spatial attention in the frontal eye field. *J. Neurosci.* **25**, 9479–9487 (2005).
44. Schall, J. D. On the role of frontal eye field in guiding attention and saccades. *Vision Res.* **44**, 1453–1467 (2004).
45. Storey, J. D. A direct approach to false discovery rates. *J. R. Stat. Soc. Series B Stat. Methodol.* **64**, 479–498 (2002).

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

© The Author(s), under exclusive licence to Springer Nature America, Inc. 2020

## Methods

### Detailed description of model. High-dimensional description of observations.

Our model describes trial-by-trial neuronal activity with a linear regression with respect to the task variables. We assume that the activity of the  $i^{\text{th}}$  neuron  $y_{i,k}(t)$  at time  $t$  on trial  $k$  can be described by a linear combination of  $P$  task variables  $x_k^{(p)}$ ,  $p = 1, \dots, P$  (eg, stimulus variables, behavioral outcomes and nonlinear combinations thereof), such that

$$y_{i,k}(t) = x_k^{(1)}\beta_{i,1}(t) + x_k^{(2)}\beta_{i,2}(t) + \dots + x_k^{(P)}\beta_{i,P}(t) + \epsilon_{i,k}(t). \quad (3)$$

where the  $P$  values of the task variables  $x_k^{(p)}$  are known, the  $\beta_{i,p}(t)$  are unknown coefficients and  $\epsilon_{i,k}(t)$  is noise. This basic model structure is identical to that of the regression model used in ref. <sup>1</sup> and has been successfully employed in characterizing neuronal activity of single neurons in other studies of perceptual decision-making<sup>46,47</sup>. In cases where we include a time-varying mean rate that is independent of the task variables, we define  $x_k^{(P)} \equiv 1$  for all  $k$ , and the  $P^{\text{th}}$  component becomes the time-varying mean.

To represent all neurons simultaneously, we concatenate the responses into a vector  $\mathbf{y}_k(t)$  and write

$$\mathbf{y}_k(t) = x_k^{(1)}\boldsymbol{\beta}_1(t) + x_k^{(2)}\boldsymbol{\beta}_2(t) + \dots + x_k^{(P)}\boldsymbol{\beta}_P(t) + \boldsymbol{\epsilon}_k(t), \quad (4)$$

where  $\mathbf{y}_k(t) = (y_{1,k}(t), \dots, y_{n,k}(t))^{\top}$ ,  $\boldsymbol{\beta}_p(t) = (\beta_{1,p}(t), \dots, \beta_{n,p}(t))^{\top}$  and  $\boldsymbol{\epsilon}_k(t) = (\epsilon_{1,k}(t), \dots, \epsilon_{n,k}(t))^{\top}$ . For trial epochs of duration  $T$ , we can regard all observations on a given trial to be a matrix,  $\mathbf{Y}_k = (\mathbf{y}_k(1), \dots, \mathbf{y}_k(T))$ , giving the observation model

$$\mathbf{Y}_k = x_k^{(1)}\mathbf{B}_1 + x_k^{(2)}\mathbf{B}_2 + \dots + x_k^{(P)}\mathbf{B}_P + \mathbf{E}_k, \quad (5)$$

where  $\mathbf{E}_k = (\boldsymbol{\epsilon}_k(1), \dots, \boldsymbol{\epsilon}_k(T))$  and  $\mathbf{B}_p = (\boldsymbol{\beta}_p(1), \dots, \boldsymbol{\beta}_p(T))$ . For the present study, we assume the noise is normally distributed  $\boldsymbol{\epsilon}_k(t) \sim \mathcal{N}(0, \mathbf{D}^{-1})$  for all trials  $k$  and times  $t$ , where  $\mathbf{D} = \text{diag}(\lambda_1, \dots, \lambda_n)$  is a  $n \times n$  diagonal matrix of noise precisions.

**Low-dimensional description of observations.** With no additional constraints, our observation model (Fig. 5) is extremely high dimensional and is effectively a separate linear regression for each neuron at every time point. This would only be a sensible model if we thought that neurons were not, in fact, coordinating activity between each other or across time. To define our low-dimensional model, we can describe each  $\mathbf{B}_p$  by a low-rank factorization, that is,  $\mathbf{B}_p = \mathbf{W}_p \mathbf{S}_p$ , where  $\mathbf{W}_p$  and  $\mathbf{S}_p$  are  $n \times r_p$  and  $r_p \times T$ , respectively, where  $r_p = \text{rank}(\mathbf{B}_p)$ . Equivalently, we can say that  $r_p$  is the dimensionality of the encoding of task variable  $p$ . This is equivalent to saying that the characteristic response of each neuron to the  $p^{\text{th}}$  task variable can be expressed as a linear combination of  $r_p$  weighted basis functions  $\beta_{i,p}^p(t) = \sum_{j=1}^{r_p} w_{ij}^{(p)} s_j^{(p)}(t)$ , where  $r_p$  is the dimensionality of the encoding,  $\{s_j^{(p)}(t)\}_{j=1}^{r_p}$  are a common set of time-varying basis functions and  $\{w_{ij}^{(p)}\}_{j=1}^{r_p}$  are neuron-dependent mixing weights.

The mTDR model does not impose any orthogonality between task variables or task variable subspaces. This permits accurate recovery of subspaces even when the encoding dimensions are correlated, which can result in correlations between task variable representations, as we demonstrate in Supplementary Note 8. It is desirable, therefore, to be able to visualize the part of the encoding of each task variable that is unmixed. We, therefore, orthogonalize the subspaces with respect to correlated subspaces for visualization in Fig. 4.

**Marginal estimation of model parameters.** The goal of inference is to estimate the factors of  $\mathbf{B}_p$  and the ranks  $r_p$ . Our proposed estimation strategy, for computational and statistical efficiency, is to estimate only one set of factors ( $\{\mathbf{W}_p\}$  or  $\{\mathbf{S}_p\}$ ). This is possible when we integrate out one set of factors. For example, if we define a prior probability density over the mixing weights  $p(\mathbf{W})$ , then, for data likelihood  $p(\mathbf{Y}|\mathbf{W}, \mathbf{S})$ , the marginal likelihood of the matrix of time-varying basis functions  $\mathbf{S}$  can be obtained by

$$p(\mathbf{Y}|\mathbf{S}, \lambda) = \int_{-\infty}^{\infty} p(\mathbf{Y}|\mathbf{W}, \mathbf{S}, \lambda) p(\mathbf{W}) d\mathbf{W}. \quad (6)$$

In principle, either set of factors may be selected for marginalization. In practice, however, the set of factors with lowest dimension should be selected to keep computational costs low. In this paper, we focus on the case where  $T \ll n$ , and we, therefore, will estimate the set of weights  $\{\mathbf{S}_p\}$  while integrating over  $\{\mathbf{W}_p\}$ . The fact that either set of factors may be determined in this way means that there is a duality between rows and columns imposed by this model that is similar, in principle, to the duality between factors and latent states for probabilistic PCA<sup>48</sup>.

If we let the noise distribution and prior distribution of  $\mathbf{W}$  both be Gaussian, then we can use standard Gaussian identities to derive the marginal density  $p(\mathbf{Y}|\mathbf{S}, \lambda)$  and the corresponding posterior density  $p(\mathbf{W}|\mathbf{Y}, \mathbf{S}, \lambda)$ . A simple starting assumption would be to let all elements of  $\mathbf{W}$  be independent standard normal (that is,  $\text{vec}(\mathbf{W}) \sim \mathcal{N}(0, \mathbf{I}_{rn})$  where  $\tilde{r} = \sum_p \text{rank}(\mathbf{B}_p)$ ). We, therefore, assume that the weights are a priori independent and that the noise variance is independent across both neurons and time. In principle, our framework supports

the application of more structured priors and noise covariances, but we will not explore more elaborate models in this paper. Further details are developed in Supplementary Math Note 2.

**Experimental details.** A detailed description of these data has been published previously<sup>1</sup>. Briefly, two adult male rhesus monkeys were trained to perform a context-dependent, two-alternative, forced-choice visual discrimination task. At the beginning of each trial, the monkeys were cued (Fig. 1a) to respond to either the motion or the color parts of the stimulus. After the context-cue presentation, two targets appear for 350 ms, followed by a 750-ms presentation of the stimulus. The stimulus was then followed by a randomized 300–1,500-ms delay, after which the monkey was cued to indicate its decision with a saccade to either of the two targets. The position of red and green targets was randomized on each trial.

Electrophysiological data were recorded from tungsten electrodes implanted in the arcuate sulcus in and around the FEF. Electrodes were lowered two at a time into adjacent grid holes and were advanced until at least one single unit could be isolated, although some trials yielded multi-unit activity. Data were recorded using the Multichannel Acquisition Processor Data Acquisition System (Plexon). All recorded units were included in the analysis. Spike sorting was conducted by clustering based on PCA using the Plexon Offline Sorter (Plexon). Each isolated cluster was functionally treated as a unit. Some clusters did not correspond to well-discriminated, single-unit activity and were therefore deemed multi-unit activity.

All analyses presented in this paper used spike counts binned at 50 ms (for model fitting and decoding) or 12.5 ms (for display of projections, jPCA and PSTHs). Analysis windows for both monkeys started 100 ms after stimulus onset and continued for 100 ms after stimulus onset<sup>1</sup>. Color coherence was transformed into position evidence based on the location of the red and green target. All data were analyzed with custom scripts written in MATLAB (MathWorks).

**Model structure.** Inclusion of linear terms for color, motion, choice and context were substantiated by previous work<sup>1</sup>. Examination of the PSTHs revealed that stimulus encoding was asymmetric (eg, unit 2 in Fig. 1c), such that the encoding of the stimulus strength was stronger in one direction than the other. This suggested that the absolute value of the stimulus strengths should be jointly modeled with the linear encoding of the stimuli. Model fits using terms for the absolute value of the stimuli resulted in smaller AIC than model fits with only linear terms (monkey A:  $\text{AIC}_{\text{linear}} = 9.79 \times 10^7$ ,  $\text{AIC}_{\text{abs}} = 7.33 \times 10^7$ ; monkey F:  $\text{AIC}_{\text{linear}} = 8.065 \times 10^7$ ,  $\text{AIC}_{\text{abs}} = 8.0628 \times 10^7$ ).

In general, we suggest that investigators proceed with task variable inclusion in the same way that one would when performing traditional linear regression. This process should include careful consideration of the phenomenology of encoding in the population that they are studying and principled model selection metrics. For large numbers of putative task variables, investigators should consider model selection via sparse priors on regression coefficients.

**Cross-validated variance explained.** To assess the variance in the population responses that is explained by our method, we conducted four-fold cross-validation where, on each fold of cross-validation, we used a randomly selected sample of 75% of the trials as training data to estimate the parameters of the model. Using the remaining 25% of the trials as test data, we made PSTHs for every possible task variable contingency for correct trials (total of 144 conditions). The reported variance explained was averaged over the four cross-validation folds.

When assessing variance explained, the population PSTHs for each condition was averaged over all extraneous task variables. For example, to assess the variance explained by the motion subspaces, we averaged the PSTHs over all task variables except motion. We, therefore, had six sets of PSTHs for each neuron that was projected onto the motion subspace.

To determine if the variance that was explained by the estimated subspaces was greater than chance, we compared the observed variance explained to the distribution of variance explained obtained by random projections. As a surrogate null distribution, we generated 500 samples for each task variable of random projection weights from a normal distribution and calculated the explained variance for each sample. We then asked what the probability was of the observed explained variance being larger than the explained variance of the random projections for each neuron. We found that many neurons exceeded the 95% Bonferroni-corrected significance threshold across nearly all dimensions.

**seqPCA.** The seqPCA algorithm identifies an orthogonal basis on which variance of a  $D$ -dimensional trajectory is sequentially explained. The algorithm starts by calculating the variance explained by the first singular vector of a sequence of  $D \times t$  data matrices  $\mathbf{Y}_t$ , where  $t$  indicates the number of time points included in the data. As the number of data points increases, the first singular vector explains a larger proportion of the variance,  $p_{1,t}$ , until trajectories change direction, after which  $p_{1,t}$  decreases. The  $t$  at which  $p_{1,t}$  reaches its peak is considered a transition time, and the left singular vector at this time is considered the first seqPC. Variability explained by this axis is subtracted from the data, and the procedure is repeated to identify the second seqPC and so on. For details, see Supplementary Math Note 9.

The seqPCA algorithm displays some sensitivity to noise by making peaks in  $p_{1,t}$  difficult to identify. However, moderate smoothing (Gaussian window, 50-ms width)

of the trajectories appeared to mitigate this effect. Greater robustness might be offered by translation of this algorithm into an optimization framework<sup>49</sup>. A related method has been developed for identification of sequential motifs of spike rasters<sup>50</sup>.

**Statistics.** No statistical methods were used to predetermine sample sizes, but our sample sizes are similar to those reported in previous publications<sup>1</sup>. Data collection and analysis were not performed blinded to the conditions of the experiments.

Tests for rotational dynamics by the method of Elsayed and Cunningham<sup>23</sup> depend on the distribution coming from a special form of tensor-variate normal distribution that corresponds to the maximum entropy distribution for tensor-valued data with independent tensor dimensions. Data distribution was assumed to be normal, but this was not formally tested.

Significance of the magnitude of inner products (Fig. 5 and Extended Data Fig. 5) was determined by a null distribution based on the positive half-Gaussian with zero mean and s.d.  $\sigma_0 = 1/n$ , where  $n$  is the number of neurons ( $n = 762$  for monkey A and  $n = 640$  for monkey F), and controlled by the positive false discovery rate<sup>45</sup>. We conducted the same procedure using standard deviations from bootstrap samples but found that the asymptotic formula ( $\sigma_0 = 1/n$ ) was slightly more conservative as the  $\sigma_0$  was a few percent larger. The null distribution was used following Kobak et al.<sup>27</sup> and validated by the one-sample Kolmogorov–Smirnov test with 1,000 permutations of the weight indices as samples from the surrogate null distribution ( $P = 0.78$  and  $D_n = 0.0207$ ).

Permutation tests for canonical correlations (Supplementary Math Note Fig. 3) were performed with 200 uniformly randomized permutations.

**Reporting Summary.** Further information on research design is available in the Nature Research Reporting Summary linked to this article.

## Data availability

Data are available for download at <https://www.ini.uzh.ch/en/research/groups/mante/data.html>.

## Code availability

Demo code for the mTDR method is available for MATLAB at <http://www.mikioaoi.com/samplecode/RDRdemo.zip>

## References

46. Brody, C. D., Hernández, A., Zainos, A. & Romo, R. Timing and neural encoding of somatosensory parametric working memory in macaque prefrontal cortex. *Cereb. Cortex* **13**, 1196–1207 (2003).
47. Roitman, J. D. & Shadlen, M. N. Response of neurons in the lateral intraparietal area during a combined visual discrimination reaction time task. *J. Neurosci.* **22**, 9475–9489 (2002).
48. Lawrence, N. Probabilistic non-linear principal component analysis with Gaussian process latent variable models. *J. Mach. Learn. Res.* **6**, 1783–1816 (2005).
49. Cunningham, J. P. & Ghahramani, Z. Linear dimensionality reduction: survey, insights, and generalizations. *J. Mach. Learn. Res.* **16**, 2859–2900 (2015).
50. Mackevicius, E. L. et al. Unsupervised discovery of temporal sequences in high-dimensional datasets, with applications to neuroscience. *eLife* **8**, e38471 (2019).

## Acknowledgements

This work was supported by grants from the McKnight Foundation, the Simons Collaboration on the Global Brain (SCGB AWD543027 to M.C.A. and J.W.P.), the National Institutes of Health BRAIN Initiative (R01EB026946 and NS104899 to J.W.P.), the National Science Foundation CAREER Award (IIS-1150186 to J.W.P.) and a U19 NIH-NINDS BRAIN Initiative Award (5U19NS104648 to M.C.A. and J.W.P.). V.M. was supported by the Swiss National Science Foundation (SNSF Professorship PP00P3-157539), the Simons Foundation (to W. T. Newsome and V.M., award 328189), the Swiss Primate Competence Center in Research, the Howard Hughes Medical Institute (through W. T. Newsome, investigator) and the DOD | USAF | AFMC | Air Force Research Laboratory: W. T. Newsome, agreement number FA9550-07-1-0537.

## Author contributions

M.C.A. and J.W.P. developed the model and performed data analysis. V.M. conceived and conducted the experiments and collected the data. All authors helped with the interpretation of analysis and writing of the paper.

## Competing interests

The authors declare no competing interests.

## Additional information

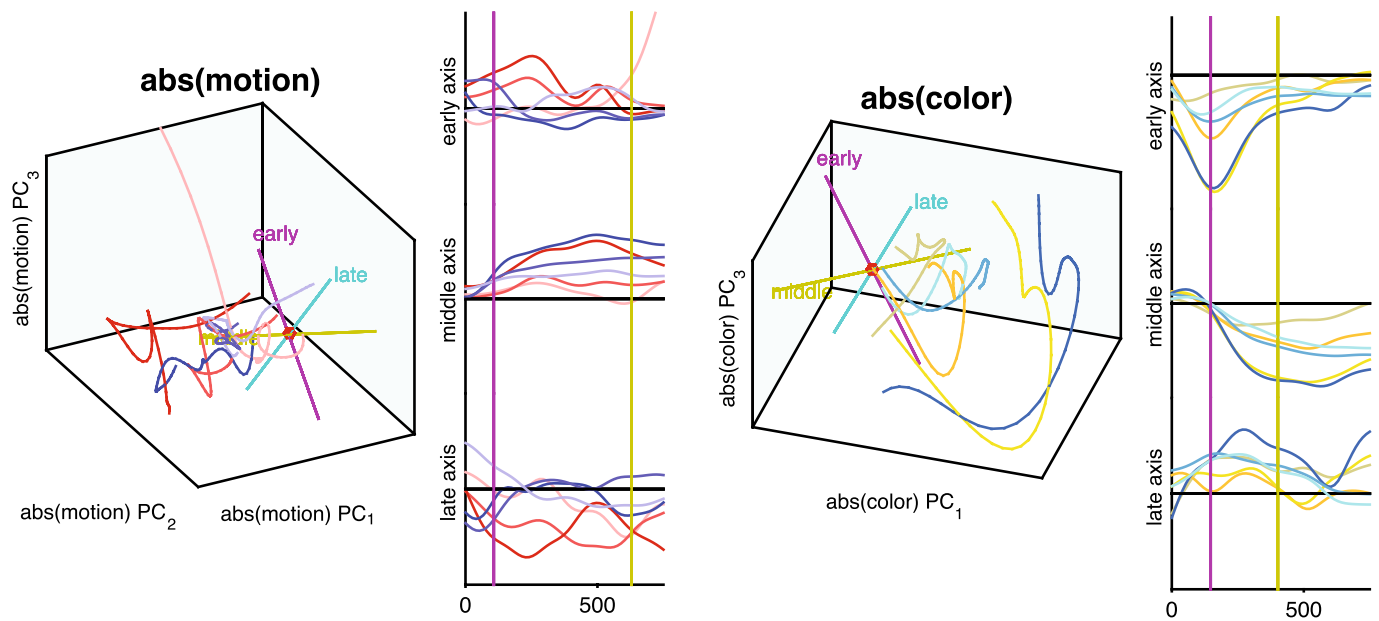
**Extended data** is available for this paper at <https://doi.org/10.1038/s41593-020-0696-5>.

**Supplementary information** is available for this paper at <https://doi.org/10.1038/s41593-020-0696-5>.

**Correspondence and requests for materials** should be addressed to M.C.A.

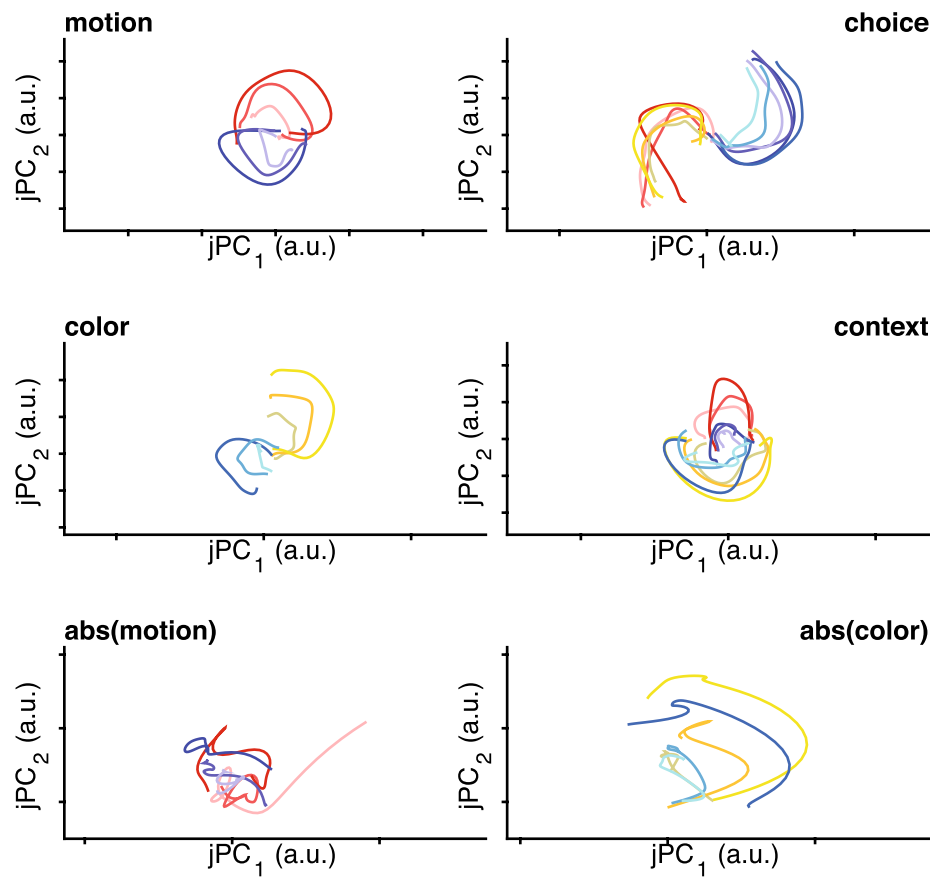
**Peer review information** *Nature Neuroscience* thanks John Cunningham, Camillo Padoa-Schioppa and the other, anonymous, reviewer(s) for their contribution to the peer review of this work.

**Reprints and permissions information** is available at [www.nature.com/reprints](http://www.nature.com/reprints).

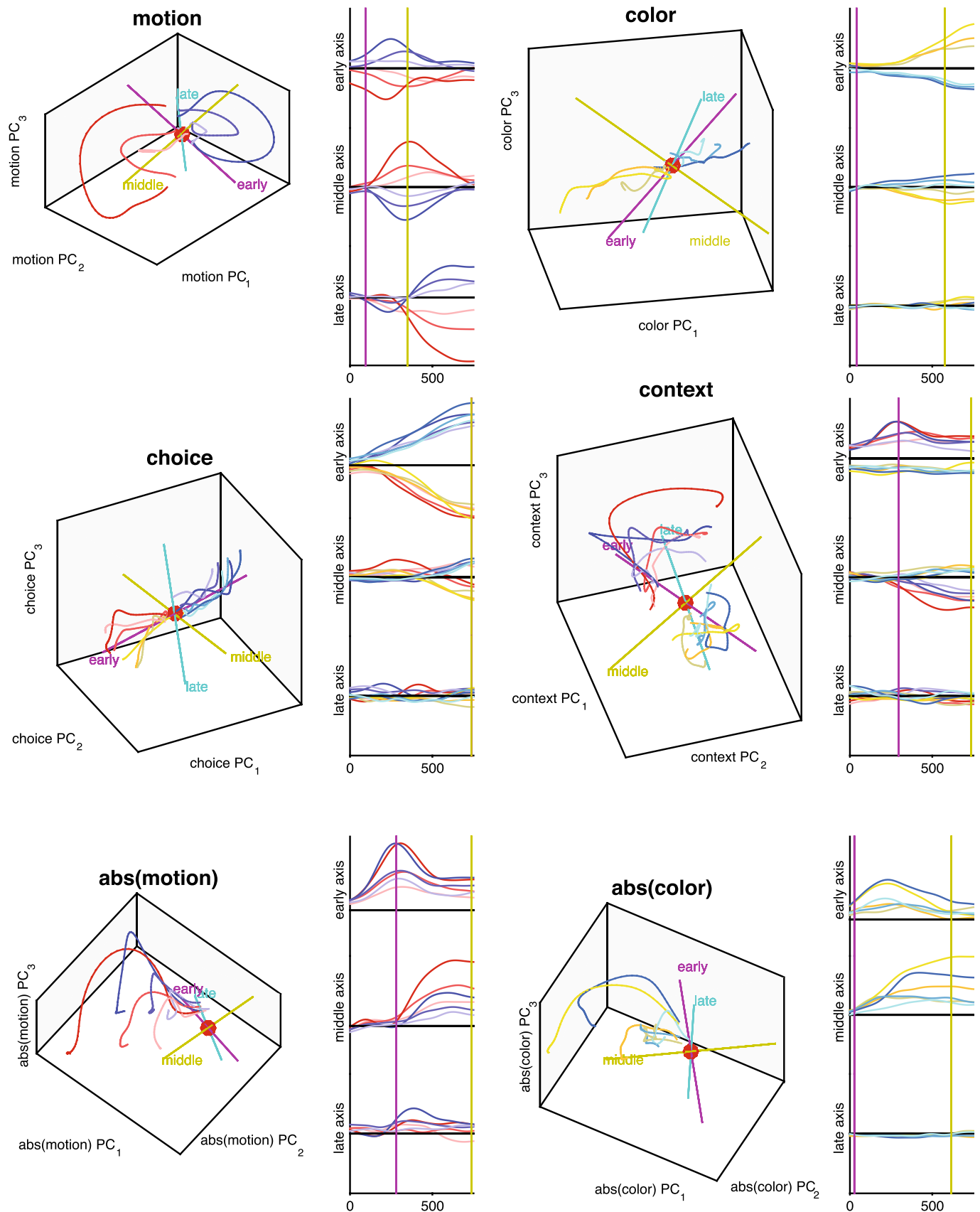


**Extended Data Fig. 1 | Projections of population PSTH's onto the first, second, and third PC-axes for monkey A.** **a**, The  $\text{abs(motion)}$  and **b**,  $\text{abs(color)}$  subspaces. Subspaces have been orthogonalized with respect to the first dimension of the choice subspace. The monkey gave the correct response for all trials used. Colored axes indicate dominant axes in the early, middle, and late periods of the stimulus epoch, as determined by the methods described in Supplementary section 9. Purple vertical lines indicate transition from the early to middle epochs. Yellow vertical lines indicate transition from the middle to late epochs as in Figure 4. Plotting colors are the same as those in Figure 4. Units of the ordinate are arbitrary but all axes are on the same scale.

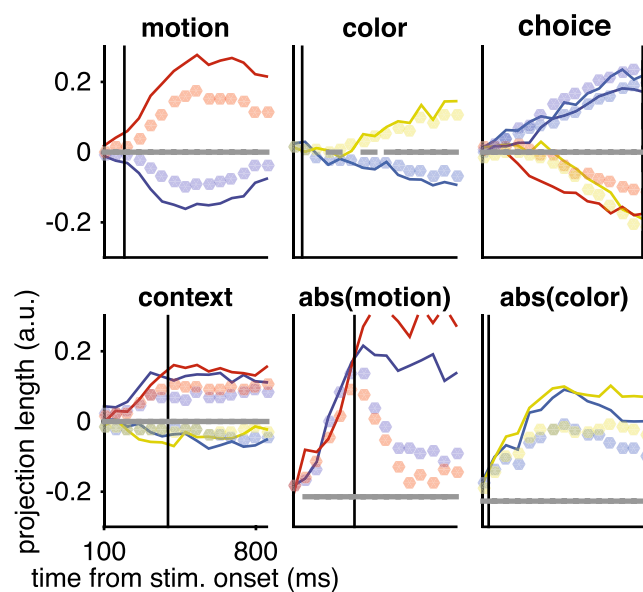




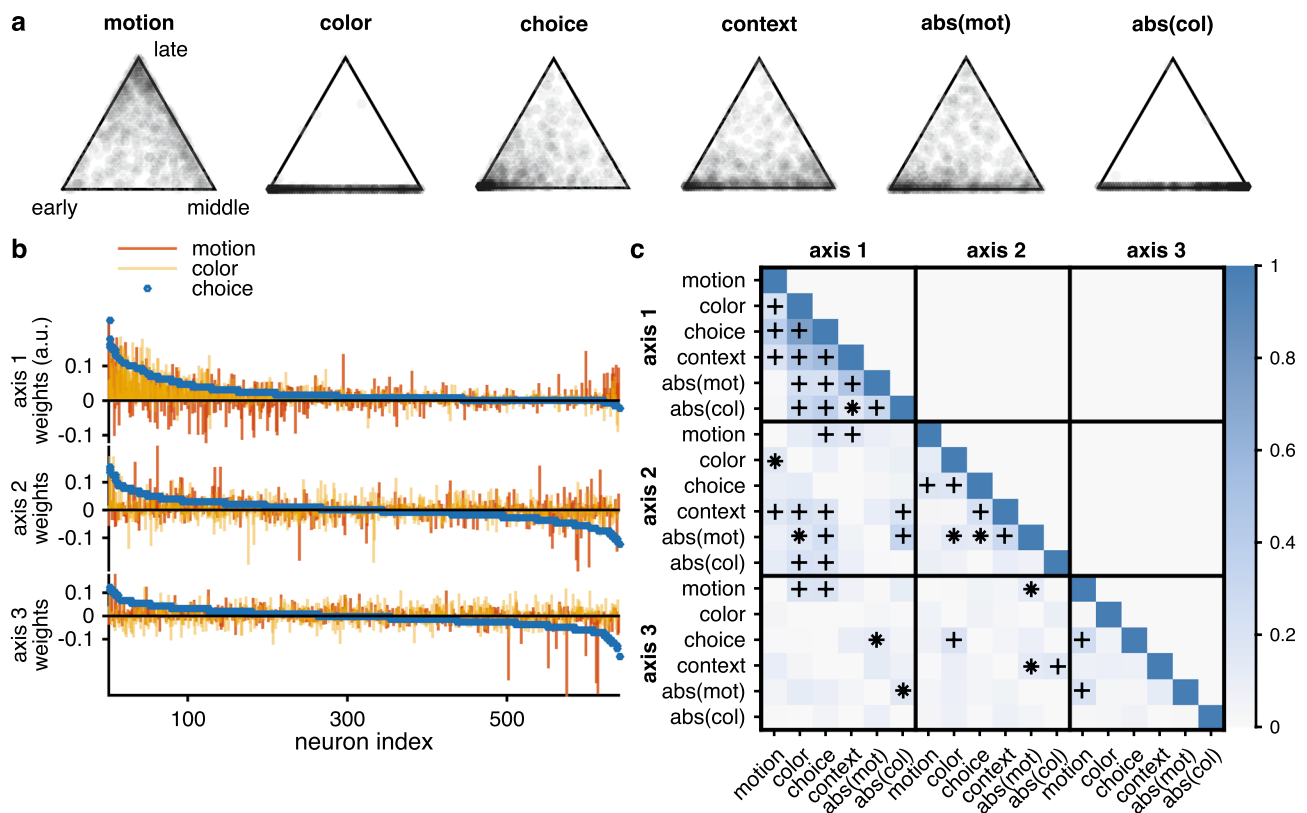
**Extended Data Fig. 2 | Projections of population PSTH's onto jPCA axes for monkey A.** Projections are onto the first two jPCA axes identified by the trajectories shown in Figure 4. The jPCA axes reveal strongly rotational dynamics for motion, color, choice, and context subspaces.



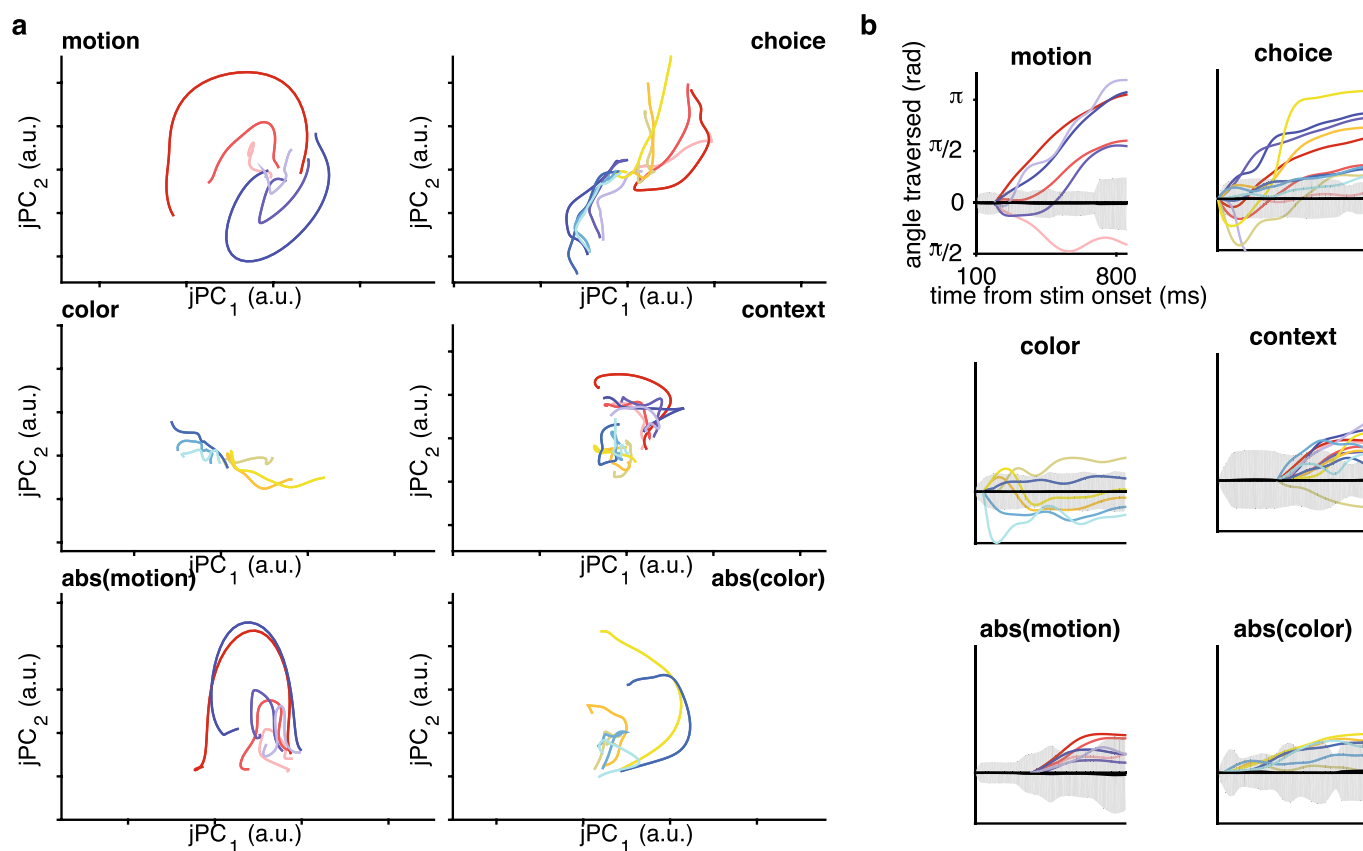
**Extended Data Fig. 3 |** Projections of population PSTH's for monkey F onto the first, second, and third PC-axes of all task variables subspaces. Plotting conventions and analyses are the same as those for Figure 4. Projected data is averaged over 2-folds of cross validated projections where a random sampling of half of the data was used to estimate parameters and the remaining half used to make projections.



**Extended Data Fig. 4 | Encoding strength of population pseudosamples for monkey F onto the first three axes of all task variables subspaces.** Plotting conventions and analyses are the same as those for Figure 4. Projected data is averaged over 2-folds of cross validated projections where pseudosamples were drawn from held-out trials. Grey bars at  $y = 0$  indicate time points where the mTDR projections had significantly stronger encoding across all stimulus levels than the 1D projections (left-tailed Wilcoxon signed-rank test,  $p\text{FDR}^{45}$  controlled at .01).

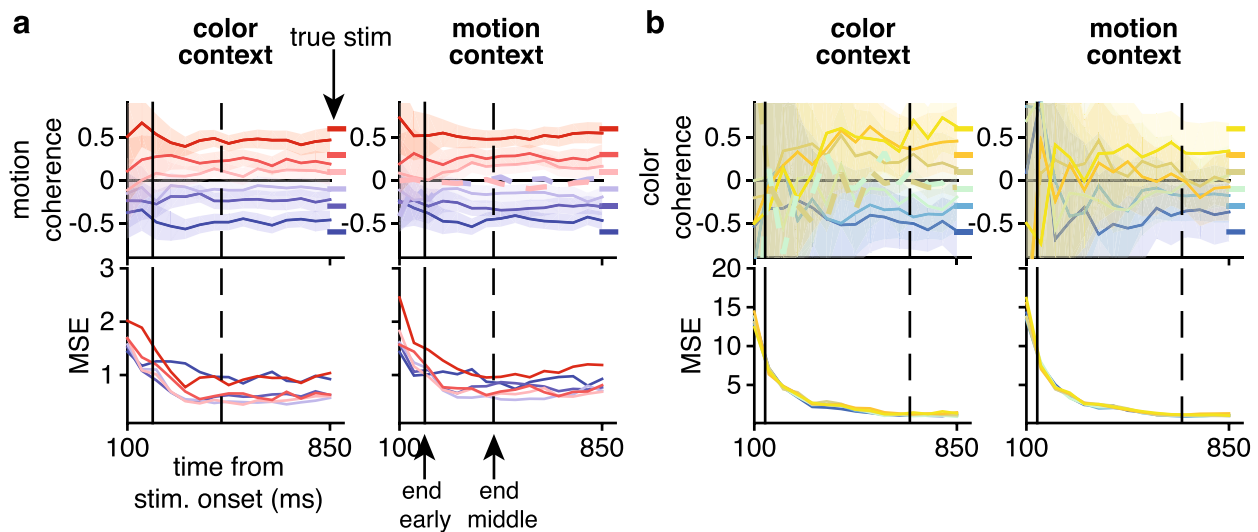


**Extended Data Fig. 5 | Distribution of variance among seqPCA axes. Monkey F.** Plotting conventions are the same as for Figure 5. **a**, Proportion of variance among seqPCA axes. Each marker corresponds to one neuron. The position of each neuron indicates the distribution of variance from PSTHs across corresponding early, middle, and late axes. e.g. a point that lies closer to the 'early' vertex of the motion plot has more of its motion-specific variance explained by the early axis while a point in the middle of the simplex has variance equally distributed across all axes. Darker regions indicate higher density of points. Colored dots correspond to cells displayed in Figure 3. **b**, Weights of the top (in terms of variance explained) 3 axes for all cells for motion, color, and choice subspaces. Cell indexes are sorted according to the choice weights from most positive to most negative. **c**, Magnitude of the Pearson correlation between top 3 subspace axes. The magnitude is used because the axes are only identifiable up to a sign. Markers indicate significant correlations controlled by the positive false discovery rate<sup>45</sup> (\*  $Q < .01$ , +  $Q < .01$ ). Null distribution is based on the positive half-Gaussian with zero-mean and standard deviation  $\sigma_0 = 1/n$ , where  $n = 640$  is the number of neurons. Significant correlations are most consistent between color-choice and motion-choice pairs.

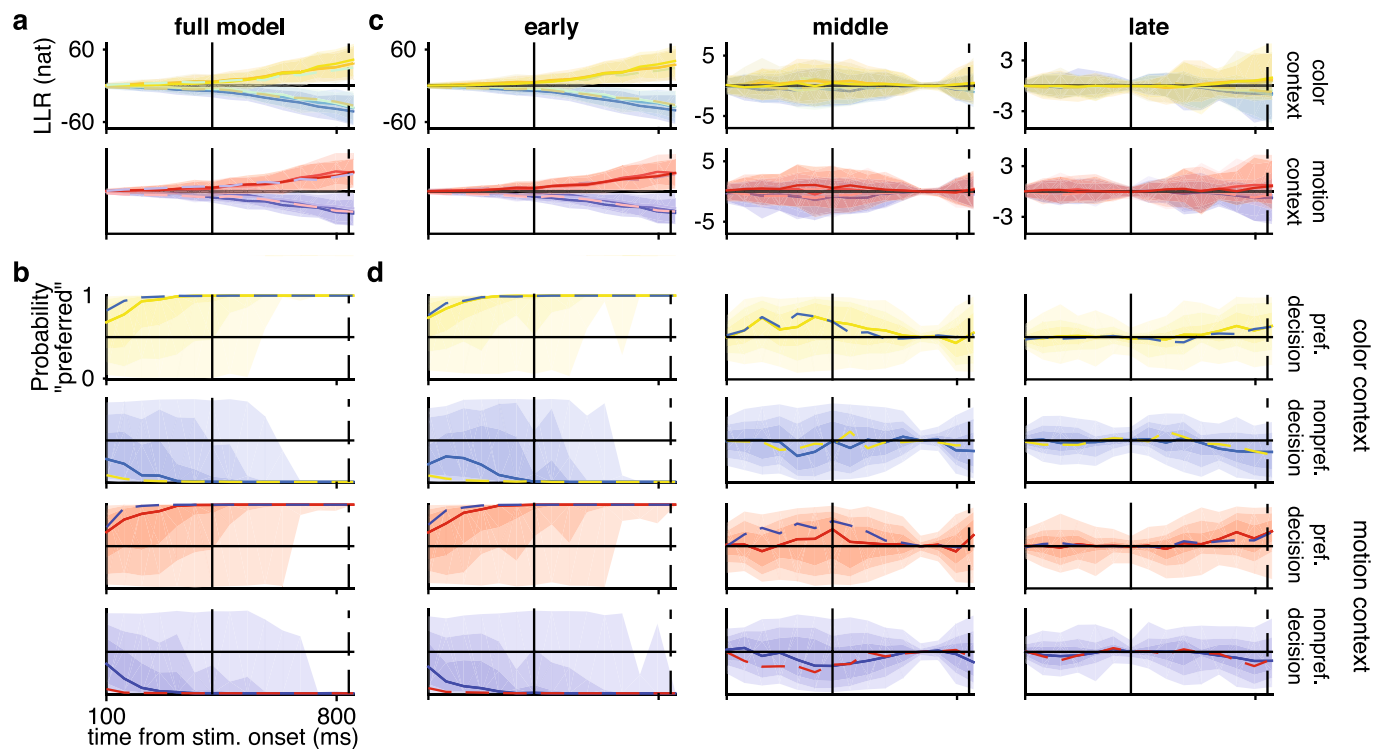


**Extended Data Fig. 6 | Rotational dynamics of subspace projections for Monkey F.** **a**, Projections of population PSTH's for monkey F onto the jPC-axes of all task variables subspaces. Plotting conventions and analyses are the same as those for Figure 4. Projected data is averaged over 2-folds of cross validated projections where a random sampling of half of the data was used to estimate parameters and the remaining half used to make projections. **b**, Angle of rotation over time for low-D trajectories of monkey F. Rotation angle traversed through rotational projection using jPCA. Angle was calculated starting from time when the projection transitions between the early and middle epochs. Coherent traversal across stimulus strengths that is consistent and monotonically increasing is an indication of rotation. Shaded areas are 95% confidence regions calculated using a maximum entropy method<sup>23</sup> ( $n = 100$  samples) under the null hypothesis of no population structure other than the empirical means and covariances across time, neurons, and task conditions.

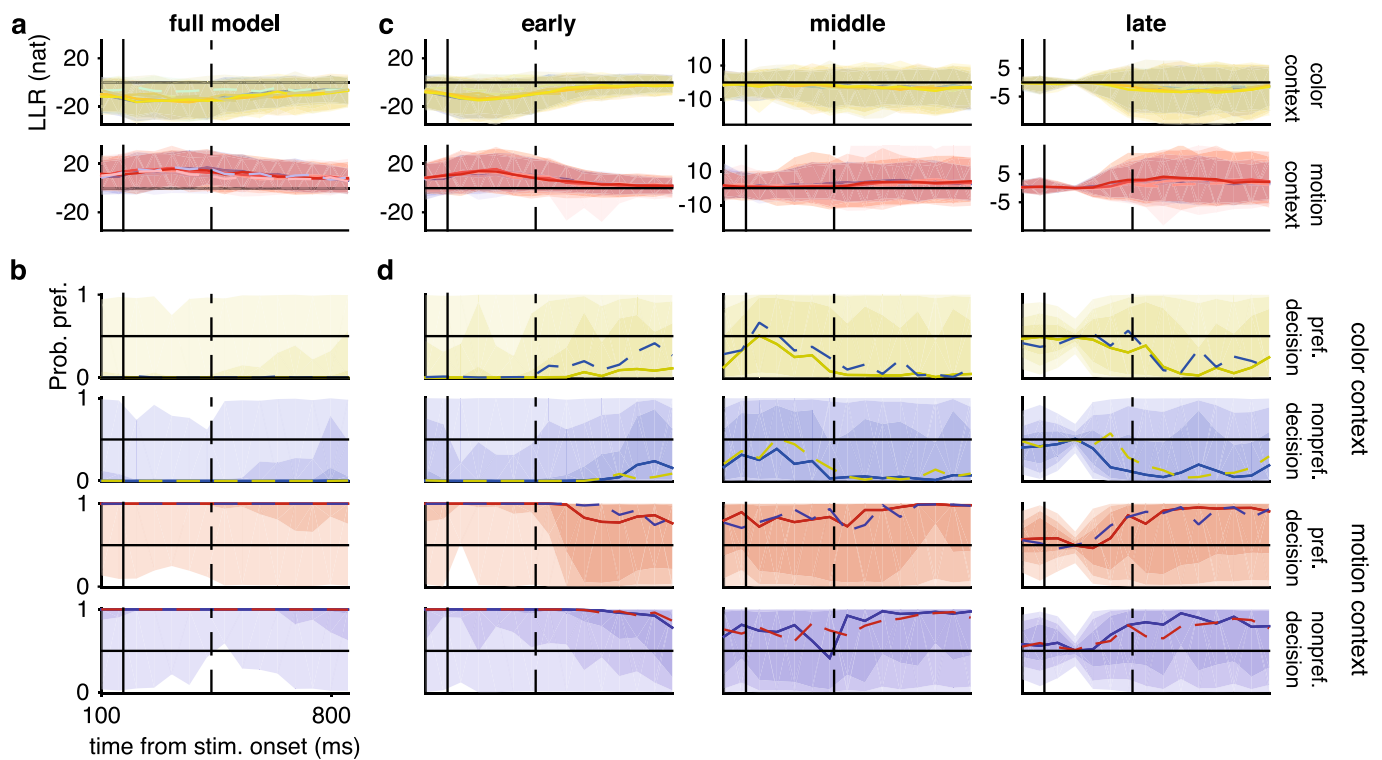




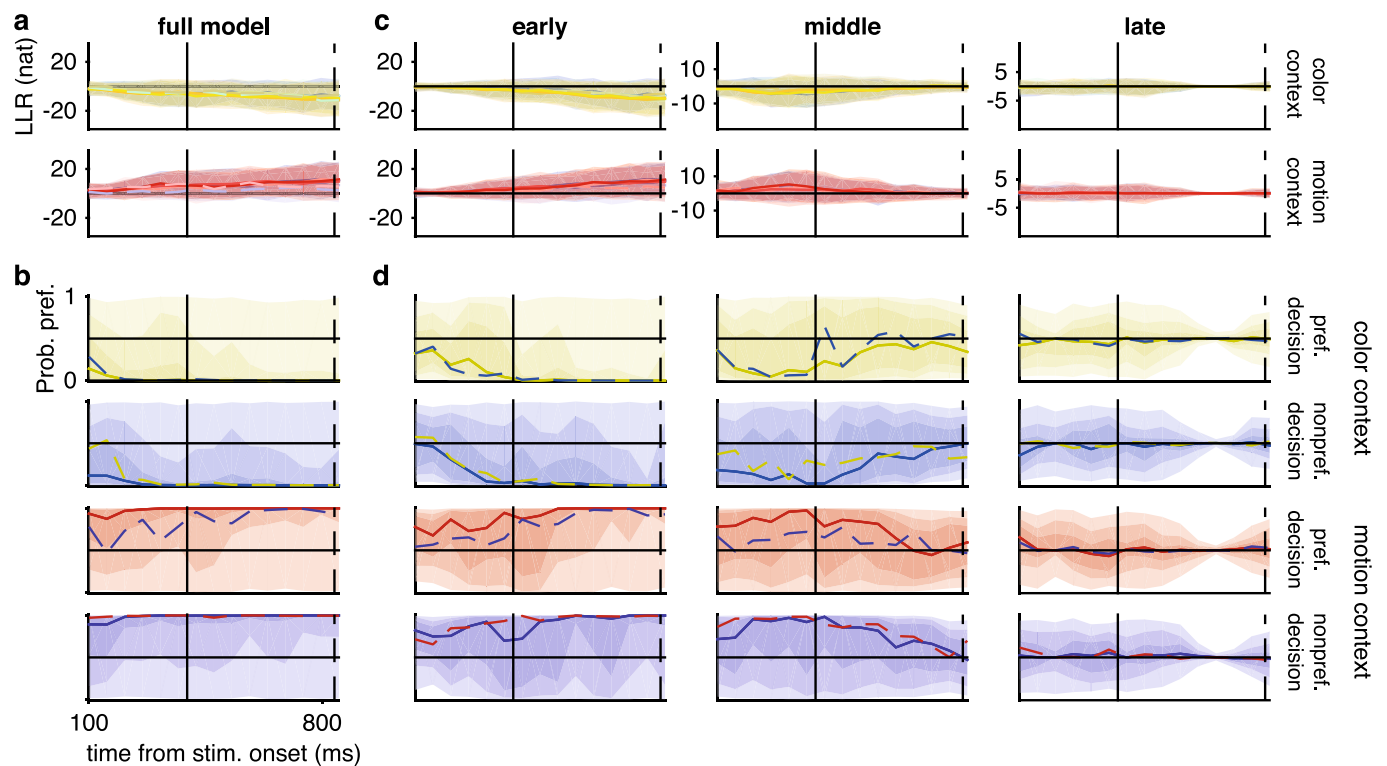
**Extended Data Fig. 7 | Instantaneous decoding of stimulus for monkey F.** Plotting conventions and analyses are the same as for Figure 6 **a**. Top: Decoded motion coherence by mTDR model in both contexts. Bottom: Mean squared error (MSE) over time of motion coherence decoding across stimulus levels and context. MSE decreases precipitously, and then stabilize around the time of the first transition. **b**, Same as **a**) for color coherence decoding. Shaded regions indicate 50% confidence intervals. Dashed lines indicate error trials from the corresponding context for the lowest stimulus strengths. 100 pseudotrials for each of 2-fold cross validation used for each analysis. Solid vertical lines indicate the time of early/middle axis transition for the corresponding stimulus subspace projections. Dashed vertical lines indicate the time of middle/late transition.



**Extended Data Fig. 8 | Instantaneous decoding of decision for monkey F.** Plotting conventions and analyses are the same as for Figure 6 **a**, Log-likelihood ratios (LLRs) in favor of a preferred choice using single pseudotrials from color - context (gold-blue, sorted by color coherence) and motion - context (red-violet, sorted by motion coherence) trials. Shaded regions indicate 95% quantile intervals for each stimulus strength. Solid lines indicate the median of correct trials. Dashed lines indicate median of error trials. **b**, Probability of a preferred choice based on corresponding LLRs combined over all stimulus strengths (see section 6.3 for details). Solid lines indicate median of correct trials. Dashed lines indicate median of error trials. Shaded regions indicate quantile coverage intervals of correct trials (light-to-dark: 95%, 75%, 50%). 100 pseudotrials for each of 2-fold cross validation folds used for all analyses. **c**, LLRs in favor of a preferred choice where the choice subspace has been restricted to only the early, middle, or late axes. **d**, Probability of a preferred choice based on LLRs from (**c**).



**Extended Data Fig. 9 | Instantaneous decoding of context for monkey A.** **a**, LLRs for monkey A in favor of the motion context using single pseudotrials, sorted by color coherence. Shaded regions indicate 95% quantile intervals for each stimulus strength. Solid lines indicate the median over correct trials. Dashed lines indicate median of error trials. **b**, Probability of the motion context based on corresponding LLRs combined over all stimulus strengths. Solid lines indicate median of correct trials. Dashed lines indicate median of error trials. Shaded regions indicate quantile intervals of correct trials (light-to-dark: 50%, 75%, 95%). Color conventions are the same as in Figure 4. 100 pseudotrials for each of 4-fold cross validation folds used for all analyses.



**Extended Data Fig. 10 | Instantaneous decoding of context for monkey F.** Plotting conventions are the same as in Extended Data 9. 100 pseudotrials for each of 2-fold cross validation folds used for all analyses.