

Neuron, Volume 109

Supplemental information

**The geometry of neuronal representations
during rule learning reveals complementary
roles of cingulate cortex and putamen**

Yarden Cohen, Elad Schneidman, and Rony Paz

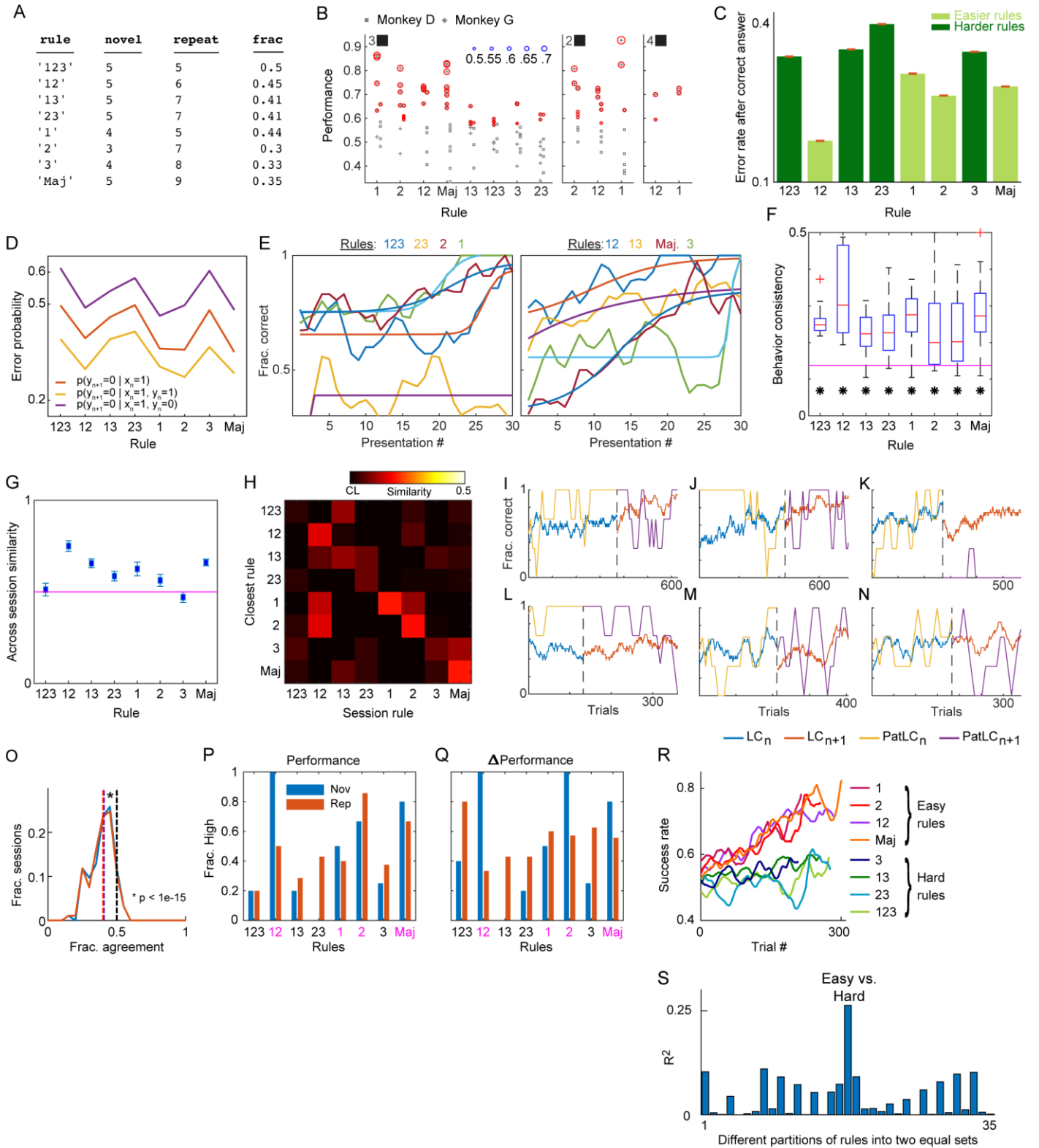


Fig. S1. Performance and rule-dependency (Related to main figure 2)

A. The number of sessions for each rule when it was different from the rule in the previous day-session (novel) or when it is the same rule (repeat).

B. Performance. Identical to Fig. 2E but with ordering the rules according to mean performance (left to right). Shown is the performance at the last quarter of each session (P_{end}), and the red circles show the confidence level. Both animals had at least one successful session in each rule, but four rules (left) had many more successful sessions than the other four (right).

C. Error probability after making a correct classification. For each rule shown is the probability of making an error on a pattern after choosing correctly in the previous presentation of that same pattern. Error bars depict SE (red). The strong dependency on the specific rule suggests the monkeys did not memorize specific patterns after a correct classification ($\chi^2(7)=88.97$, $p < 1e-15$, Kruskal-Wallis test).

D. Error probability in different stimulus-response associations show rule-dependent coupling. For each rule (x-axis) shown is the probability of making an error in classifying a pattern 'y' after correctly classifying a different pattern 'x' in the previous presentation of 'x' (red line). In yellow/purple shown is the same error probability conditioned also on correctly/incorrectly classifying 'y' in the previous presentation. The latter calculation shows the global effect of learning – fewer errors after correct choices. All probabilities have strong dependency on the rule being learned (Kruskal-Wallis test, $\chi^2 > 88$, $p < 1e-10$ for all), demonstrating a strong coupling between stimulus-response pairs as expected in rule-based learning.

E. Acquisition rates for 'salient' patterns. The most salient pattern, '000', was labeled 'right' and 'left' equally in the set of rules (4+4 rules). Shown are the learning curves (colored lines + sigmoidal fits), the proportion of correct labeling (y-axis, smoothed across days with a running window of 4 presentations) as a function of pattern presentations (x-axis). The two panels are for the 4 rules in which 000 is labeled 'left' (left panel) and for the 4 it was labeled 'right' (right panel). There was variability across rules in how the visually salient pattern was learned, showing that patterns that draw more attention are not learned faster or similarly across different rules, and providing more evidence against a simple memorization process (Kruskal-Wallis test, $24 < \chi^2(3,36) < 33$, $1e-7 < p < 1e-4$ for all).

F-H. Behavior consistency. **F.** Within-session consistency scores (y-axis, methods) for all sessions with a single rule (x-axis). All rules showed above-chance level consistency (t-test, $p < 0.05$ for all, chance-level in magenta). **G.** Between-session pairwise consistency score across sessions with the same rule. **H.** A classifier fitted to the last 25% of each session in which a certain rule was learned (x-axis) and compared to the other rules (y-axis, see methods). Color scale is from chance level (CL) to the maximal possible classification value (0.5).

I-N. Learning disrupts previously acquired stimulus-response associations. Blue and red lines show overall performance in consecutive sessions (LC_n, LC_{n+1}). The underlying rule is switched between the sessions (dashed line) but in these rules 4 out of 8 patterns did not change their correct stimulus-response association. The yellow and purple lines show the performance in one of these 4 patterns ($PatLC_n, PatLC_{n+1}$). As seen and in contrast to learning by independent stimulus-response pairing, rule-based learning implies that pattern-specific performance can deteriorate while the general performance improves (**A-D** for monkey G, **E,F** for monkey D). This happens even when the association-specific performance reaches 100% (**A,E**).

O. Win→Stay, Lose→Shift strategy poorly describes behavior. We simulated win→stay, lose→shift answer sequences, using the same pattern presentations and rules as in the actual sessions. Shown are the histograms of the overlap between the simulated sessions with the actual behavior (Red and blue for simulations starting with the left and right choices). Colored dashed lines show the median values, both smaller than chance level (0.5, black dashed line, Wilcoxon's signed ranked test, $p < 1e-15$).

P. The fraction of sessions in which performance was high, separating sessions with novel rules (blue) from sessions that repeat the rule from the previous day (red). Easy rules (magenta) are the highest 4 rules in de-novo learning.

Q. The fraction of sessions in which change-in-performance within a session was high, separating sessions with novel rules (blue) from sessions that repeat the rule from the previous day (red). Here as well, easy rules (magenta) are the highest 4 rules in de-novo learning.

R. Average learning curves also separate easy and hard rules.

S. Performance variance explained by partitioning rules into two equal sets of four rules. The Easy-Hard partition explains significantly more than all other partitions.

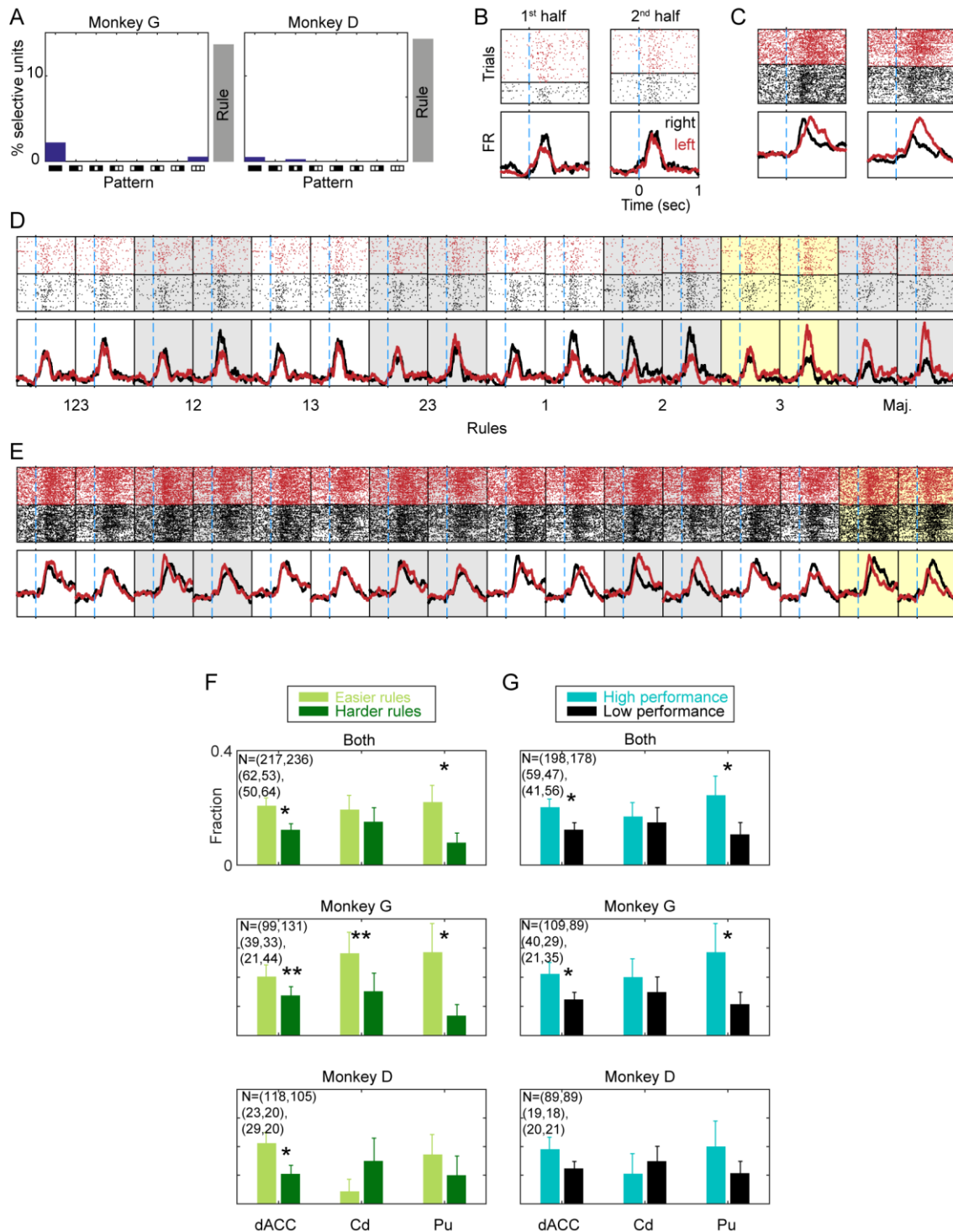


Fig. S2. Rule-representation in single-units (Related to main figure 3)

A. Pattern selective neurons. The percent of neurons that were selective for a single pattern (bars, y-axis) plotted for all patterns (x-axis). Gray bars on the right indicate the percent of rule correlated neurons. Left: monkey G. Right: monkey D. Very few neurons were selective for a single-patterns, suggesting against single stimulus-response learning.

B-E. Rasters and PSTHs from single neurons showing changes in representing the animals' choice, the correct rule, and all other rules. Spikes aligned to stimulus onset ($t=0$, x-axis, light blue dashed line). Trials (y-axis) are stacked according to the animal left or right choice (red, black dots in the top panel and curves in the bottom panel mark spike times and firing rates accordingly). Calculations are made separately in the first and second halves of the session (left - right panels).

B. A neuron that does not differentiate the animal's choice in the 2nd half of the session.

C. A neuron that learns to differentiate the animal's choice.

D. The same neuron (and data) as in panel **B**, separated by the labels of the 8 rules in the experiment (x-axis labels). Gray and white shadings distinguish the different rules and yellow shading marks the actual rule being learned.

E. similar to panel **D** but for the neuron in **C**.

F,G. Rule-representation differ in sessions of high/low performance, easy/hard rules in both monkeys.

Fraction of neurons that exhibited significant rule correlation during the last third of the session (Bars+SE), plotted for the 3 regions, and shown for both animals combined (top), separately for monkey G (middle) and monkey D (bottom).

F. Separating sessions by easy-hard rules. A-top (similar to Fig.3D): dACC: $z=2.43$, $p<0.01$, Pu: $z=2.12$, $p<0.02$; Monkey G: dACC: $z=1.305$, $p<0.1$, Cd: $z=1.3275$, $p<0.1$, Pu: $z=2.37$, $p<0.01$; Monkey D: dACC: $z=2.17$, $p<0.02$.

G. Separating sessions by performance, independent of the rule. Both: dACC: $z=2.046$, $p<0.021$, Pu: $z=1.79$, $p<0.04$; Monkey G: dACC: $z=1.85$, $p<0.04$, Pu: $z=1.97$, $p<0.03$.

Binomial z-test (*= $p<0.05$, **= $p<0.1$).

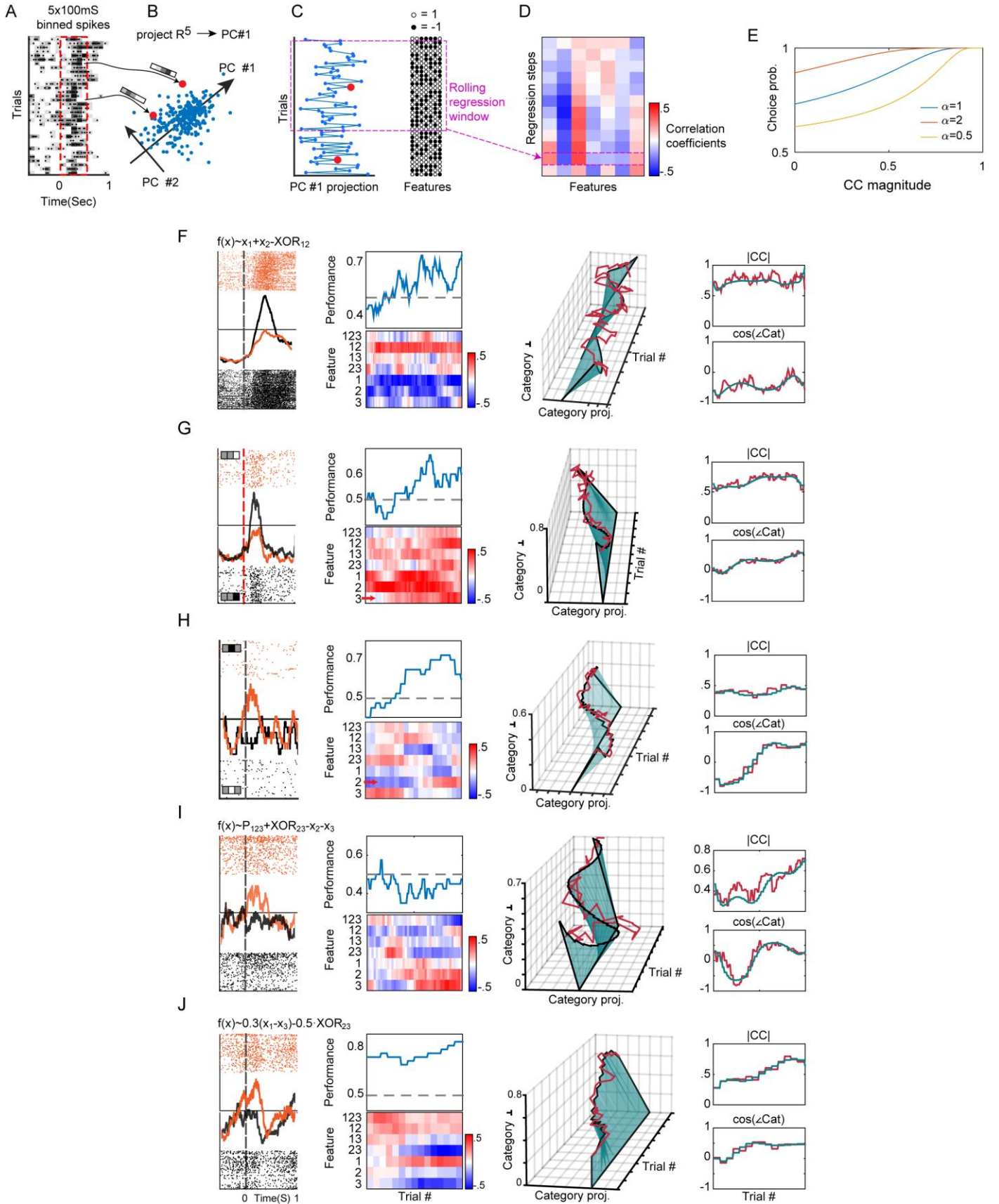


Fig .S3. Single neurons' dynamic feature representation (Related to main figure 4)

A. Spikes (dots) during 500 msec after the stimulus presentation (x-axis, dashed red frame) are counted in 5 x 100 msec bins (gray scale).

- B.** The 5-vectors of spike counts from all trials are projected on the principal component explaining the most variance across trials. (arrows and red dots show specific examples).
- C.** The resulting neural 1-d stimulus response, x-axis, is correlated across trials (y-axis) with stimulus visual features (7-vectors of black and white circles) in a rolling regression window (purple frame).
- D.** The resulting vector of 7 correlation coefficients spans the neuron's visual feature preference during that regression window and its dynamics across regression steps (*neural-vector*, y-axis).
- E.** Neural vector magnitude translates to behavioral confidence: confidence is defined as the choice probability of a decision based on neural activity (see methods 'Interpreting neural-vector magnitude as confidence'). As the magnitude of the features correlation vector (x-axis) increases, so does the choice probability (y-axis), indicating an increasing confidence.
- F-J. Examples of learning related neural dynamics.** Same format as in Fig.4, showing neurons with stable feature selectivity, high dimensional rotation, magnitude increase, and complex trajectories.

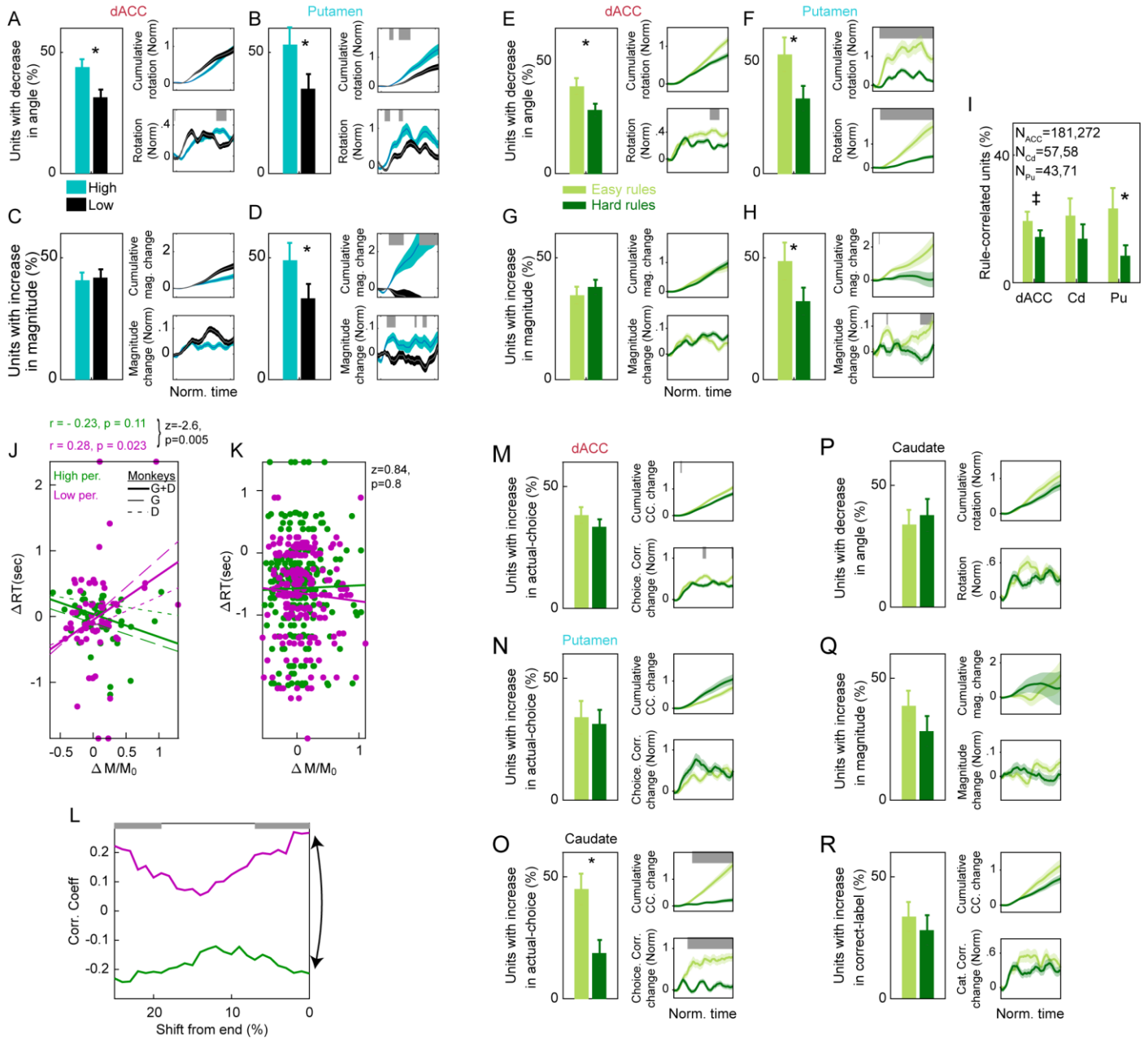


Fig .S4. Differential representation dynamics in the dACC, the Putamen, and the Caudate (Related to main figure 5)

A-D. Dividing sessions by performance. Replicating the results of main Fig. 5 but when dividing sessions into high and low performance independent of the rule type (i.e. instead of dividing session according to easy and hard rules, as in Fig. 5). High – Low performance was determined by above – below the median performance for each animal.

E-I. Redefining the Majority rule as hard for monkey D replicates the results. E-H. Same presentation as in Fig. 5 and replicating the results. **I.** Same presentation as in Fig. 2D and replicating the results.

J. Correlations between change in neural-vector magnitude and reaction-time are significantly different between sessions of high and low performance (Z-test of Fisher transformed Pearson correlation coefficients, $Z = -2.6, -2.05, -1.6, p = 0.005, p = 0.02, p = 0.055$ for both monkeys, monkey G, monkey D)

K. Same as in J for dACC neurons.

L. The correlation across neurons between change in reaction time and change in magnitude, shown for different offsets from the end of the session (x-axis, % of session duration). Green and purple lines show

the coefficients for high and low performance sessions. Gray bars mark significant differences (Fisher corrected measure, methods, $p < 0.05$). Arrow marks the contrast in panel J.

M-R. Neural representation in the Caudate reflect choice. **M,N.** The dACC (**M**) and Putamen (**N**) do not show a differential change in representing the actual choice. **O.** Caudate neurons show a rule-dependent (easy vs. hard) increase in representing the animals' actual-choice. **P-R.** Caudate neurons do not show a differential (easy vs. hard) decrease in angle-to-rule (**P**), or an increase in magnitude (**Q**), or in representation of the correct-label (**R**).

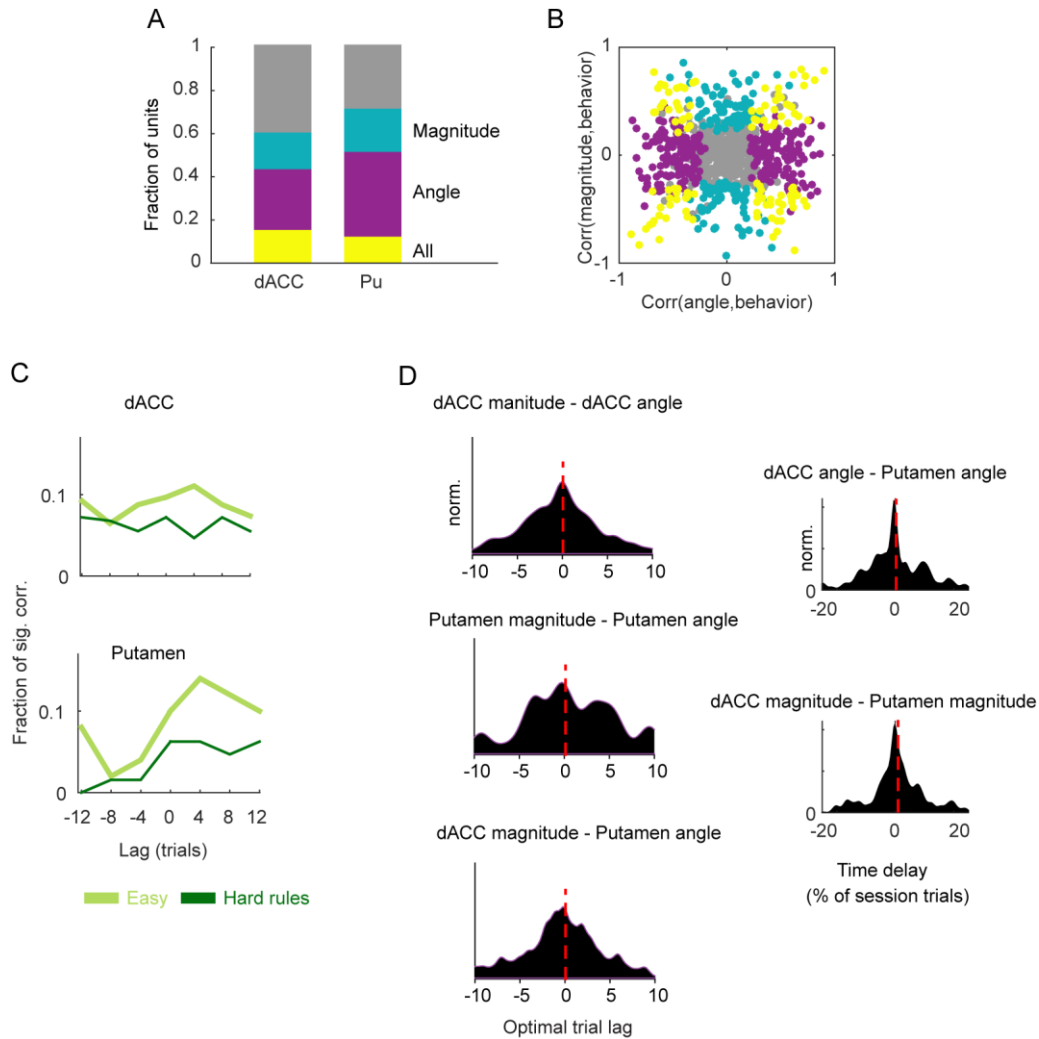


Fig. S5. Temporal alignment and similarity of dynamic geometric representations and learning curves (Related to main figure 6)

A. Enumeration of the significant (Pearson, $p < 0.05$) correlations to all variable combinations in the different regions (x-axis). The colored bars show the fraction of units with significant correlations between the learning curve and both angle-to-rule and correlations vector-magnitudes (yellow), only angle-to-rule (purple), or only vector-magnitudes (turquoise).

B. Comparing the magnitude correlation (y-axis) to the angle-to-rule correlation (x-axis) for all neurons. Color-coding as in panel (A)

C. Fraction of significant correlations between the change in rotation (the derivative of the angle-to-rule) and the behavioral curve, computed for different lags between the two measures. As shown in main Fig. 6E, more neurons followed the behavior (positive lags) in *easy* than in *hard* rules.

D. For all simultaneously recorded pairs, we computed the optimal lag between the vector-magnitude and the angle-to-rule, for all possible combinations of within and between regions. All options shown above were not different than zero (t-test, $p > 0.1$ for all). Only lags between vector-magnitude in Putamen and angle-to-rule in dACC were different than zero with the Putamen magnitude following dACC rotation (shown in Fig. 6G).

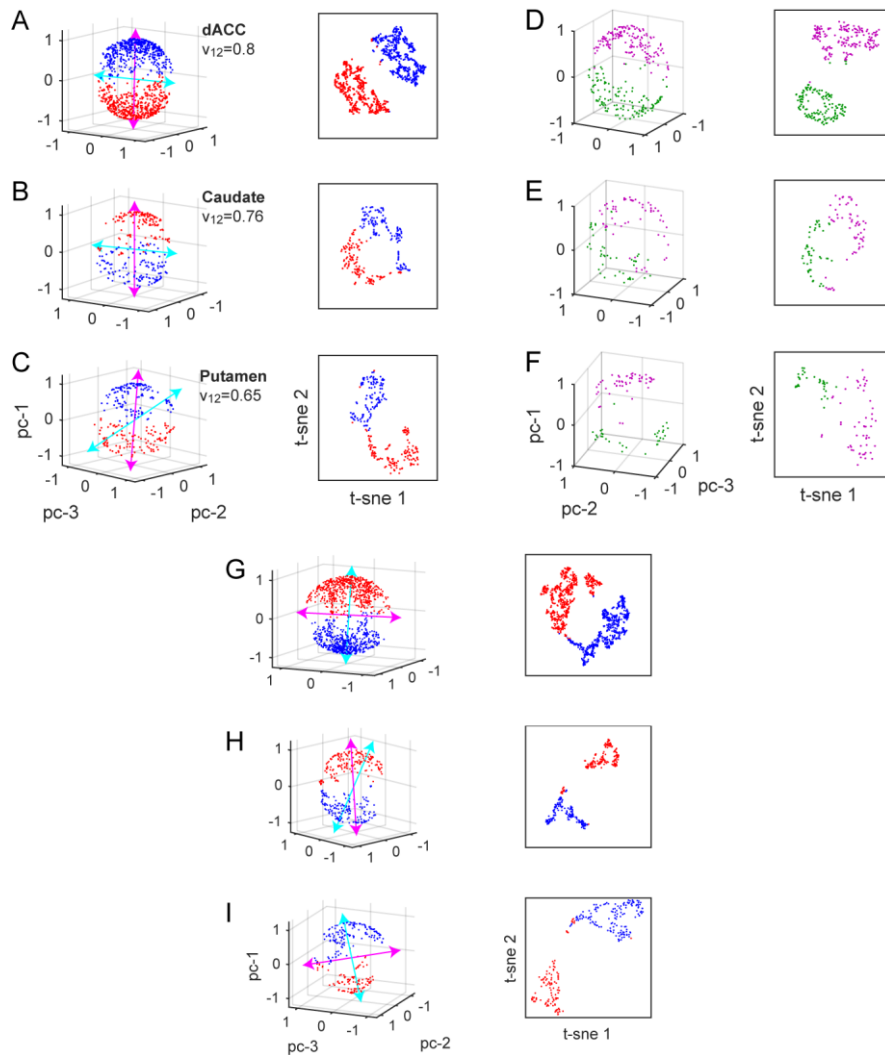


Fig. S6. Single-unit responses exhibit categorical properties (Related to methods section ‘Categorical properties of neural responses’)

Four task conditions defined by the correct category label and the animals’ answers (left or right) are used to cluster single neuron responses. Dots mark mean responses in 40 trial regression windows normalized to unit length. The left panels show responses from category-correlated regression windows, projected on their 3 main principal components. Magenta/cyan arrows show the direction of the category/answer variables. Right panels show the same data projected on 2 t-SNE components to reveal categorical separation.

A-C. Neurons in the dACC showed the lowest dimensionality (V_{12} is the fraction of variance captured by the first 2 PC’s).

D-F. Same as **A-C** but only for windows in which ANOVA tests showed significant dependence of the neural response on the task condition.

G-I. Same as **A-C** for regression windows with significant answer correlations. This reveals stronger categorical properties in the striatum and weaker in the dACC.

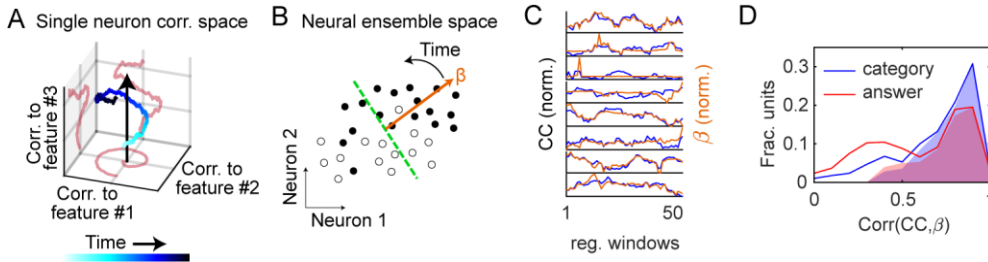


Fig. S7. Correspondence between feature-based correlations and population-level decoding of ensemble activity (Related to methods section ‘Behavior decoding from neural ensemble activity states’)

A. In our study the single neuron dynamics is a trajectory in the space of correlations to visual features (axes). A behavior policy is represented as a direction in this space (arrow). In successful learning the trajectory moves to increase the projection on the policy (magnitude increase and/or rotation towards the rule).

B. The behavior policy can be decoded from the activity of multiple neurons. The space of ensemble activity is defined by axes matching the activity of individual neurons and the response in each trial is a point in this space. A linear decoder, fitted to the behavior, defines a separating plane (dashed line) in this space perpendicular to the vector $\vec{\beta}$. To learn, the plane is rotated to achieve optimal separation by the desired policy (full vs. empty circles).

C. Example of 8 simultaneously recorded neurons. For each neuron (insets) the normalized category correlation (z-scored, blue) is overlaid on the β (z-scored, orange) for that neuron computed by a category classifier fitted in the same regression window (x-axis) to all 8 neurons. The traces expose the similarity of the two approaches and also points of disagreement.

D. Under assumptions of linear visual response of single-neurons and conditional independence between neurons, we find that the two approaches shown in A,B, can be related. The histograms show correlation coefficients between the single neuron category dynamics and the dynamics of the weight β fitted to neurons in decoding the category from the population activity. The shaded areas mark correlations with $p < 0.001$, amounting to 81% (or 62% for answers) of all neurons.