

1

2

3

4 Orthogonal neural encoding of targets and distractors
5 supports multivariate cognitive control

6
7 Harrison Ritz^{*1-3} & Amitai Shenhav^{1,2}

8

9 1. *Cognitive, Linguistic & Psychological Science, Brown University, Providence, RI, USA*
10 2. *Carney Institute for Brain Science, Brown University, Providence, RI, USA*
11 3. *Princeton Neuroscience Institute, Princeton University, Princeton, NJ, USA*

12

13 * *Corresponding author:* hritz@princeton.edu

14

15

1 Abstract

2 The complex challenges of our mental life require us to coordinate multiple forms of neural
3 information processing. Recent behavioral studies have found that people can coordinate
4 multiple forms of attention, but the underlying neural control process remains obscure. We
5 hypothesized that the brain implements multivariate control by independently monitoring
6 feature-specific difficulty and independently prioritizing feature-specific processing. During
7 fMRI, participants performed a parametric conflict task that separately tags target and distractor
8 processing. Consistent with feature-specific monitoring, univariate analyses revealed spatially
9 segregated encoding of target and distractor difficulty in dorsal anterior cingulate cortex.
10 Consistent with feature-specific attentional priority, a novel multivariate analysis (Encoding
11 Geometry Analysis) revealed overlapping, but orthogonal, representations of target and distractor
12 coherence in intraparietal sulcus. Coherence representations were mediated by control demands
13 and aligned with both performance and frontoparietal activity, consistent with top-down
14 attention. Together, these findings provide evidence for the neural geometry necessary to
15 coordinate multivariate cognitive control.

16
17 Keywords: cognitive control, attention, decision-making, fMRI

1 Introduction

2 We have remarkable flexibility in how we think and act. This flexibility is enabled by the array
3 of mental tools we can bring to bear on challenges to our goal pursuit^{1–6}. For example, someone
4 may respond to a mistake by becoming more cautious, enhancing task-relevant processing, or
5 suppressing task-irrelevant processing⁷, and previous work has shown that people
6 simultaneously deploy multiple such strategies at the same time in response to different task
7 demands^{3,8–10}. Flexibly coordinating multiple cognitive processes requires a control system that
8 can monitor multiple forms of task demands and deploy multiple forms of control (also referred
9 to as the necessity for *observability* and *controllability*;¹¹). These monitoring and regulation
10 processes are fundamental to control, and are thought to be underpinned by distinct cingulo-
11 opercular and frontoparietal neural systems^{12–19}. However, much is still unknown about how
12 multiple forms of control are represented across these domains.
13

14 Past research on the neural mechanisms of cognitive control has often sought to identify
15 representations that integrate over multiple different sources of task demands (i.e., represent
16 these different sources in *alignment*). For instance, previous studies has proposed that dorsal
17 anterior cingulate cortex (dACC) tracks integrative features like response conflict, effort, value,
18 error likelihood, and time-on-task^{20–27}. Because they integrate over different task features
19 instead of differentiating between them, these forms of ‘aligned encoding’ (Figure 1a) are ill-
20 suited for carrying out multidimensional control. Multidimensional cognitive control requires
21 independent representations that can track multiple sources of difficulty and regulate multiple
22 cognitive processes (e.g., prioritize multiple sources of information²⁸).
23

24 An alternative to aligned encoding – one that would allow the brain to separately control
25 multiple processes – is *independent* encoding, which can come in at least two forms. One way
26 the brain can have independent representations is by encoding different task features in spatially
27 segregated neural populations (‘segregated encoding’; Figure 1b). For example, past work has
28 shown that different subregions within dACC encode distinct task demands, including various
29 forms of errors and processing conflict^{29–34}. The brain can instead have independent
30 representations that are distributed across units within the same population, as has also been
31 observed in dACC^{35–37}. Within a shared population, independent encoding of information occurs
32 along a set of orthogonal dimensions or *subspaces* (Figure 1c, ‘subspace encoding’;^{38–41}).
33 Despite this exciting recent work, it remains unclear to what extent different components of the
34 cognitive control system leverage these aligned, segregated, or orthogonal encoding strategies
35 for monitoring multiple task demands and prioritizing multiple sources of information.
36

37 To gain new insight into the representations supporting cognitive control, we drew upon two key
38 innovations. First, we leveraged an experimental paradigm we developed to tag multiple control
39 processes¹⁰. Building on prior work^{3,30,41,42}, this task incorporates elements of perceptual
40 decision-making (discrimination of a target feature) and inhibitory control (overcoming a salient
41 and prepotent distractor). We have previously shown that we can separately tag target and
42 distractor processing from participants’ performance on this task, and that target and distractor
43 processing are independently controlled. For example, participants adjust target and distractor
44 sensitivity in response to distinct task demands (e.g., previous conflict or incentives;¹⁰). In
45 conjunction with this process-tagging approach, our second innovation was to develop a novel

multivariate fMRI analysis for measuring relationships between feature encoding (i.e., *encoding geometry*). Extending recent statistical approaches in systems neuroscience^{35,43,44}, we combined the strengths of multivariate encoding analyses and representation similarity analyses into a method we call ‘Encoding Geometry Analysis’ (EGA). We used EGA to characterize whether putative markers of monitoring and prioritization leverage independent representations for targets and distractors.

In brief, we found that key nodes within the cognitive control network use orthogonal representations of target and distractor information to support cognitive control. In the dorsal anterior cingulate cortex (dACC), encoding of target and distractor difficulty was spatially segregated and arranged along a rostrocaudal gradient. By contrast, in the intraparietal sulcus (IPS), encoding of target and distractor coherence was encoded along orthogonal neural subspaces. These regional distinctions are consistent with hypothesized roles in planning and implementing (multivariate) attentional policies^{12,17}. Furthermore, we found that coherence encoding depended on control demands, and was aligned with both task performance and frontoparietal activity, consistent with these coherence representations playing a critical role in cognitive control (e.g., feature prioritization). Together, these results suggest that cognitive control uses representational formats that allow the brain to monitor and control multiple streams of information processing.

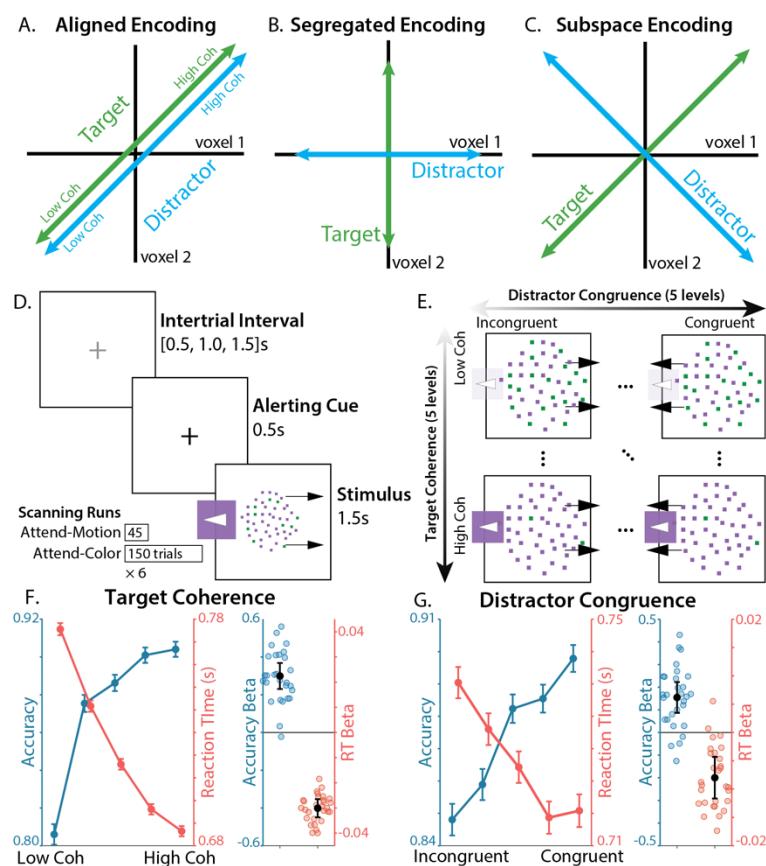
Results

Task overview

Twenty-nine human participants performed the Parametric Attentional Control Task (PACT¹⁰) during fMRI. On each trial, participants responded to an array of colored moving dots (colored random dot kinematogram; Figure 1d). In the critical condition (Attend-Color), participants respond with a left/right keypress based on which of two colors were in the majority. In alternating scanner runs, participants instead responded based on motion (Attend-Motion), which was designed to be less control-demanding due to the (Simon-like) congruence between motion direction and response hand^{3,10}. Across trials, we independently and parametrically manipulated target and distractor information across five levels of target coherence (e.g., percentage of dots in the majority color, regardless of which color) and distractor congruence (e.g., percentage of dots moving either in the congruent or incongruent direction relative to the correct color response; Figure 1e). This task allowed us to ‘tag’ participants’ sensitivity to each dimension by measuring behavioral and neural responses to independently manipulated target and distractor features. Unlike a similar task used to study post-error adjustments³, our parametric manipulation of target and distractor coherence allows us to better measure feature-specific representations. Unlike similar tasks used to study contextual decision-making^{30,41,45}, this task pits more control-demanding responses (towards color) against more automatic responses (towards motion), allowing comparisons between Attend-Color and Attend-Motion tasks to isolate the contributions of cognitive control^{46,47}.

1 Behavior

2 Participants had overall good performance on the task, with a high level of accuracy (median
 3 Accuracy = 89%, IQR = [84% - 92%]), and a low rate of missed responses (median lapse rate =
 4 2%, IQR = [0% - 5%]). We used mixed effects regressions to characterize how target coherence
 5 and distractor congruence influenced participants' accuracy and log-transformed correct reaction
 6 times. Replicating previous behavioral findings using this task, participants were sensitive to
 7 both target and distractor information ¹⁰. When target coherence was weaker, participants
 8 responded slower ($t_{(27.6)} = 16.1, p = 1.60 \times 10^{-15}$) and less accurately ($t_{(28)} = -8.90, p = 1.19 \times 10^{-9}$;
 9 Figure 1f). When distractors were more incongruent, participants also responded slower ($t_{(28.8)} =$
 10 5.09, $p = 2.15 \times 10^{-5}$) and less accurately ($t_{(28)} = -4.66, p = 6.99 \times 10^{-5}$; Figure 1g). Also
 11 replicating prior findings with this task, interactions between targets and distractors were not
 12 significant for reaction time ($t_{(28.2)} = 0.143, p = .887$) and had a weak influence on accuracy ($t_{(28)} =$
 13 2.36, $p = .0257$), with model omitting target-distractor interactions providing a better
 14 complexity-penalized fit (RT Δ AIC = 17.7, Accuracy Δ AIC = 1.38).
 15



16
 17 **Figure 1. Task and Behavior.** **A-C)** Three hypothesized encoding schemes. **A)** In *aligned encoding* features are
 18 represented similarly, e.g., encode performance variables like error likelihood or time-on-task. **B)** In *segregated*
 19 encoding features are encoded independently, in distinct voxel populations (i.e., voxel-level pure selectivity ⁴⁰). **C)**
 20 In *subspace encoding*, features are encoding independently, in overlapping voxel populations (i.e., voxel-level
 21 mixed selectivity). **D)** Participants responded to a color-motion random dot kinematogram (RDK) with a button
 22 press. Participants either responded to the left/right motion direction of the RDK (Attend-Motion runs) or based on
 23 the majority color (Attend-Color runs; critical condition). **E)** We parametrically and independently manipulated
 24 target coherence (% of dots in the majority color) and distractor congruence (motion coherence signed relative to the

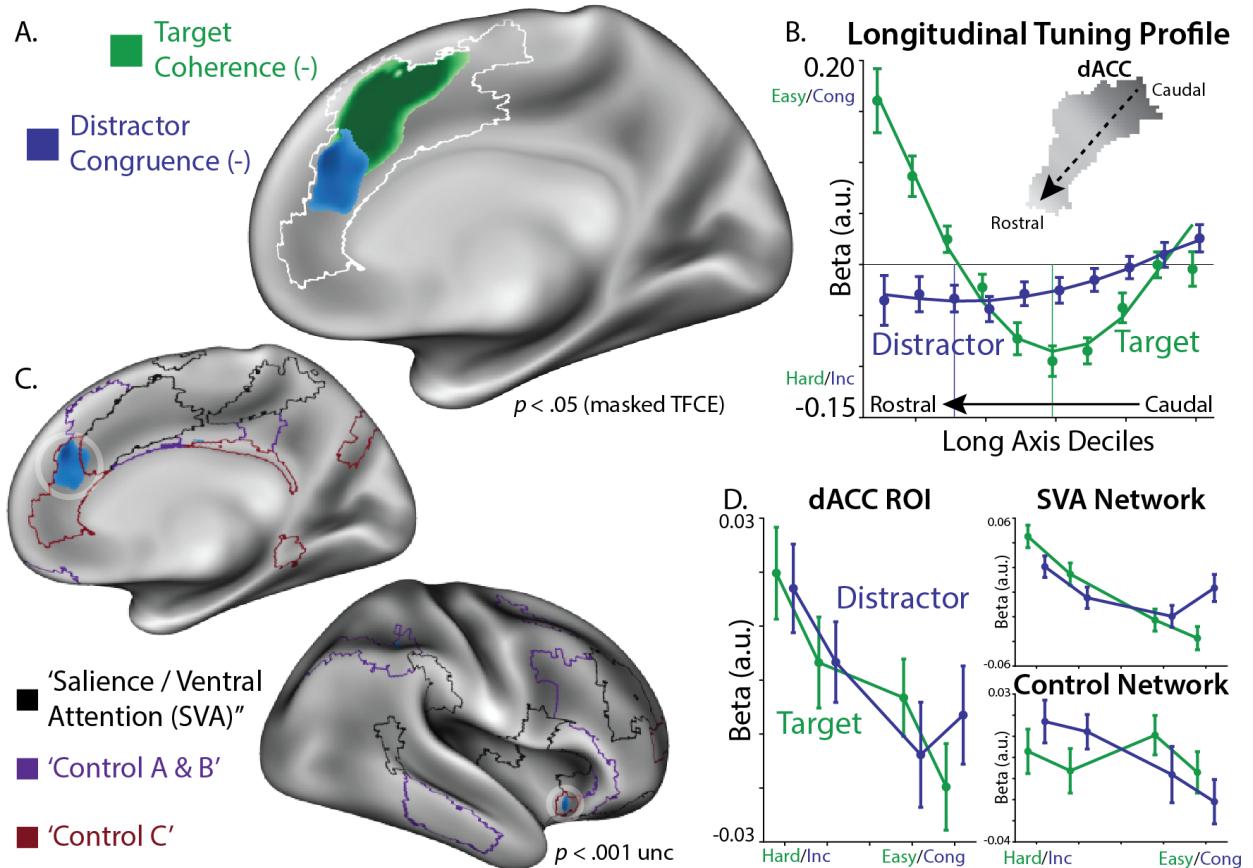
1 target response). **F**) Participants were faster and more accurate when the target was more coherent. **G**) Participants
2 were faster and more accurate when the distractor was more congruent with the target. Error bars on line plots
3 reflect within-participant SEM, error bars for regression fixed-effect betas reflect 95% CI.

4 Segregated encoding of target and distractor difficulty in dACC

5 Past work has separately shown that the dACC tracks task demands related to perceptual
6 discrimination (induced in our task when target information is weaker) and related to the need to
7 suppress a salient distractor (induced in our task when distractor information is more strongly
8 incongruent with the target^{12,30–32,48}). Our task allowed us to test whether these two sources of
9 increasing control demand are tracked within common regions of dACC (reflecting an
10 aggregated representation of multiple sources of task demands), or whether they are tracked by
11 separate regions (potentially reflecting a specialized representation according to the nature of the
12 demands).

13 Targeting a large region of dACC – a conjunction of a cortical parcellation with a meta-analytic
14 mask for ‘cognitive control’ (see ‘fMRI univariate analyses’ in Methods) – we found spatially
15 distinct signatures of target difficulty and distractor congruence within dACC. In caudal dACC,
16 we found significant clusters encoding the parametric effect of target difficulty (Figure 2a;
17 negative effect of target coherence in green), and in more rostral dACC we found clusters
18 encoding parametric distractor incongruence (negative effect of distractor congruence in blue).
19 Supporting this dissociation, the spatial patterns of target and distractor regression weights were
20 uncorrelated across dACC voxels ($t_{(28,0)} = 1.32$, $p = .197$, logBF = -0.363). These analyses
21 control for omission errors, and additionally controlling for commission errors produced the
22 same whole-brain pattern at a reduced threshold (see Supplementary Figure 1). While distractor
23 congruence was marginally significant when correcting for multiple comparisons across the
24 entire brain (one-sided $p = .08$, whole-brain TFCE), extensive previous research predicts
25 congruence effects in this ROI and in this direction^{12,20,48}, suggesting that a whole-brain
26 corrected estimate is overly conservative.
27

28



1 **Figure 2. Distinct coding of target and distractor difficulty in dACC.** A) We looked for linear target coherence and
2 distractor congruence signals within an a priori dACC mask (white outline; overlapping Kong22 parcels and medial
3 ‘cognitive control’ Neurosynth mask). We found that voxels in the most caudal dACC reflected target difficulty
4 (green), more rostral voxels reflected distractor incongruence (blue). Statistical tests are corrected using non-
5 parametric threshold-free cluster enhancement. B) We extracted the long axis of the dACC using a PCA of the voxel
6 coordinates. We plotted the target coherence (green) and distractor congruence (blue) along the deciles of this long
7 axis. Fit lines are the quantized predictions from a second-order polynomial regression. We used these regression
8 betas to estimate the minima for target and distractor tuning (i.e., location of strongest difficulty effects), finding that
9 the target difficulty peak (vertical green line) was more caudal than the distractor incongruence peak (vertical blue
10 line). C) Plotting the uncorrected whole-brain response, distractor incongruence responses (blue) were strongest
11 within the ‘Control C’ sub-network (red), both in dACC and anterior insula. D) BOLD responses across levels of
12 target coherence and distractor congruence, plotted within the whole dACC ROI (left), or the ‘Salience/Ventral
13 Attention (SVA)’ network and ‘Control’ network parcels within the dACC ROI (right). GLMs: A-C: Feature UV, D:
14 Difficulty Levels, see Table 2.

16 To further quantify how feature encoding changed along the longitudinal axis of dACC, we used
17 principal component analysis to extract the axis position of dACC voxels (see ‘dACC
18 longitudinal axis analyses’ in Methods), and then regressed target and distractor beta weights
19 onto these axis scores. We found that targets had stronger difficulty coding in more caudal
20 voxels ($t_{(27.9)} = 3.74, p = .000840$), with a quadratic trend ($t_{(26.5)} = 4.48, p = .000129$; Figure 2b).
21 In line with previous work on both perceptual and value-based decision-making^{30,49–52}, we found
22 that signatures of target discrimination difficulty (negative correlation with target coherence) in
23 caudal dACC were paralleled by signals of target discrimination ease (positive correlation with
24 target coherence) within the rostral-most extent of our dACC ROI (Supplementary Figure 2). In
25 contrast to targets, distractors had stronger incongruence coding in more rostral voxel ($t_{(28.0)} = -$

1 3.26, $p = .00294$), without a significant quadratic trend. We used participants' random effects
2 terms to estimate the gradient location where target and distractor coding were at their most
3 negative, finding that the target minimum was significantly more caudal than the distractor
4 minimum (signed-rank test, $z_{(28)} = 2.41, p = .0159$). Target and distractor minima were
5 uncorrelated across subjects ($r_{(27)} = .0282, p = .880, \log BF = -0.839$), again consistent with
6 independent encoding of targets and distractors.

7
8 As additional evidence that target-related and distractor-related demands have a dissociable
9 encoding profile, we found that the crossover between target and distractor encoding in dACC
10 occurred at the boundary between two well-characterized functional networks^{53–55}. Whereas
11 distractor-related demands were more strongly encoded rostrally in the Control Network
12 (particularly within regions of dACC and insula corresponding to the 'Control C' Sub-Network;
13^{54,56}), target-related demands were more strongly encoded caudally within the 'Salience / Ventral
14 Attention (SVA)' Network (Figure 2C-D). Including network membership alongside long axis
15 location predicted target and distractor encoding better than models with either network
16 membership or axis location alone ($\Delta BIC > 1675$).

17 Subspace encoding of target and distractor coherence in intraparietal 18 sulcus

19 We found that dACC appeared to dissociably encode target and distractor difficulty through
20 spatially segregated encoding, consistent with a role in monitoring different task demands and/or
21 specifying different control signals¹². To identify neural mechanisms for the implementation of
22 this control through the prioritization targets versus distractors, we next tested for regions that
23 encode target and distractor coherence (the amount of information in a feature, regardless of
24 which response it supports). Based on previous research, we might expect to find this form of
25 selective attention in posterior parietal cortex^{17,57,58}. We explored whether target and distractor
26 coherence share a common neural code (e.g., as a global index of spatial salience), compared to
27 where these features are encoded distinctly (e.g., as separate targets of control).

28 An initial whole-brain univariate analysis showed that overlapping regions throughout occipital,
29 parietal, and prefrontal cortices track the feature coherence (proportion of dots in the majority
30 category) for both targets and distractors (Figure 3a; conjunction in orange). These regions
31 showed elevated responses to lower target coherence and higher distractor coherence, potentially
32 reflecting the relevance of each feature for task performance. Note that in contrast to distractor
33 congruence, distractor *coherence* had an inconsistent relationship with task performance (RT:
34 $t_{(27.0)} = 2.08, p = .048$; Accuracy: $t_{(28)} = -0.845, p = .406$), suggesting that these neural responses
35 are unlikely to reflect task difficulty per se.

36
37 While these univariate activations point towards widespread and coarsely overlapping encoding
38 of the feature coherence (potentially consistent with aligned encoding; Figure 1a), they lack
39 information about how these features are encoded at finer spatial scales. To interrogate the
40 relationship between target and distractor encoding, we developed a multivariate analysis that
41 combines multivariate encoding analyses with pattern similarity analyses, which we term
42 Encoding Geometry Analysis (EGA). Whereas pattern similarity analyses typically quantify
43 relationships between representations of specific stimuli or responses (e.g., whether they could

1 be classified,⁵⁹), EGA characterizes relationships between encoding subspaces (patterns of
2 contrast weights) across different task features, consistent with recent analyses trends in systems
3 neuroscience^{35,36,43,60–62}. A stronger correlation between encoding subspaces (either positive or
4 negative) indicates that features are similarly encoded (i.e., that their representations are aligned
5 and thus confusable by a linear decoder; Figure 1a), whereas weak correlation indicate that these
6 representations are orthogonal (and thus distinguishable by a linear decoder;⁵⁹). In contrast to
7 standard pattern similarity, the sign of these relationships is interpretable in EGA, reflecting how
8 features are coded relative to one another. Compared to standard encoding analyses, EGA is less
9 sensitive to noise (Supplementary Figure 3). We estimated this encoding alignment within each
10 parcel, correlating unsmoothed and spatially pre-whitened patterns of parametric regression betas
11 across scanner runs to minimize spatiotemporal autocorrelation^{63–65}. This cross-validated
12 similarity further allowed us to anchor our analysis on the measurement reliability of encoding
13 profiles (i.e., the self-correlation of encoding patterns across cross-validation folds^{66,67}).

14

15 Focusing on regions that encoded both target and distractor information (parcels where both
16 group-level $p < .001$), EGA revealed clear dissociations between regions that represent these
17 features in alignment versus orthogonally. Within visual cortex and the superior parietal lobule
18 (SPL), target and distractor representations demonstrated significant negative correlations
19 (Figure 3b, red), reflecting (negatively) aligned encoding. In contrast, early visual cortex and
20 intraparietal sulcus (IPS; see Figure 3c for anatomical boundaries) demonstrated target-distractor
21 correlations near zero (Figure 3b, black), suggesting encoding along orthogonal subspaces.

22

23 To bolster our interpretation of the latter findings as reflecting orthogonal (i.e., uncorrelated)
24 representations rather than merely small but non-significant correlations, we employed Bayesian
25 t-tests at the group level to estimate the relative (log-10) likelihood that these encoding
26 dimensions were orthogonal or correlated. Consistent with our previous analyses, we found
27 strong evidence for correlation (positive log bayes factors) in more medial regions of occipital
28 and posterior parietal cortex (e.g., SPL), and strong evidence for orthogonality (negative log
29 bayes factors) in more lateral regions of occipital and posterior parietal cortex (e.g., IPS; Figure
30 3D). Control analyses confirmed that coherence orthogonality was not due to encoding
31 reliability, as a similar topography was observed with disattenuated correlations (normalizing
32 correlations by their reliability; see Supplementary Figure 4). Further supporting these results,
33 our Bayes factor analyses were robust to the choice of priors (see Supplementary Figure 5).

34

35

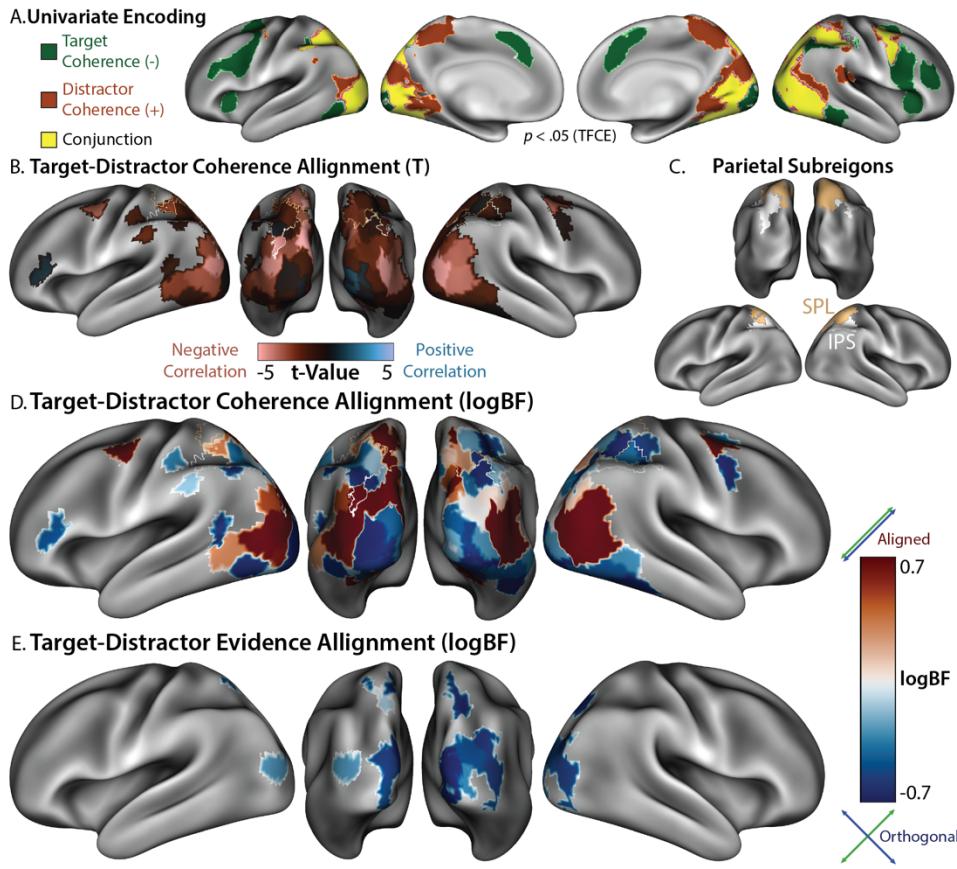


Figure 3. Encoding Geometry Analysis (EGA) dissociates target and distractor encoding. **A)** Parametric univariate responses to weak target coherence (green; percentage of dots in majority color), strong distractor coherence (orange; percentage of dots with coherent motion), and their conjunction (yellow). Statistical tests are corrected for multiple comparisons using non-parametric threshold-free cluster enhancement (TFCE). **B)** Encoding alignment within parcels in which target and distractor encoding was jointly reliable (both $p < .001$ uncorrected). Representations were negatively correlated within Superior Parietal Lobule (SPL in gold; Kong22 labels), and uncorrelated within Intraparietal Sulcus (IPS in white; Kong22 labels). **C)** Anatomical labels for parietal regions, based on the labels in the Kong22 parcellation. **D)** Bayesian analyses provide explicit evidence for orthogonality within IPS (i.e., negative BF; theoretical minima: -0.71). **E)** Coherence coded in terms of evidence (i.e., supporting a left vs right choice). Target and distractor evidence encoding overlapped in visual cortex and SPL and was represented orthogonally. GLMs: A: Feature UV, B-E: Feature MV, see Table 2.

While our analyses support independent encoding of targets and distractors within the same parcel, we further explored whether feature information is reflected in overlapping voxels (i.e., voxel-level mixed selectivity⁴⁰). Simulations revealed that the alignment between absolute encoding weights can differentiate between pure and mixed selectivity, and parietal coherence representations bore this signature of voxel-level mixed selectivity (Supplementary Figure 6), consistent with the subspace encoding hypothesis.

These results have focused on the coherence of different features regardless of the response they support, demonstrating that SPL exhibits aligned representations of target and distractor coherence. Past decision-making research has separately demonstrated that SPL tracks the amount of evidence supporting specific response^{42,68,69}, which we found was also true for our task. In addition to encoding target and distractor coherence, SPL and visual cortex also tracked target and distractor ‘evidence’ (proportion of dots supporting a rightward vs leftward response;

Figure 3e). EGA revealed orthogonal evidence representations between targets and distractors, in the same areas with aligned coherence representations (compare Figure 3d and 3e), consistent with previous observations of multiple decision-related signals in SPL⁶⁸.

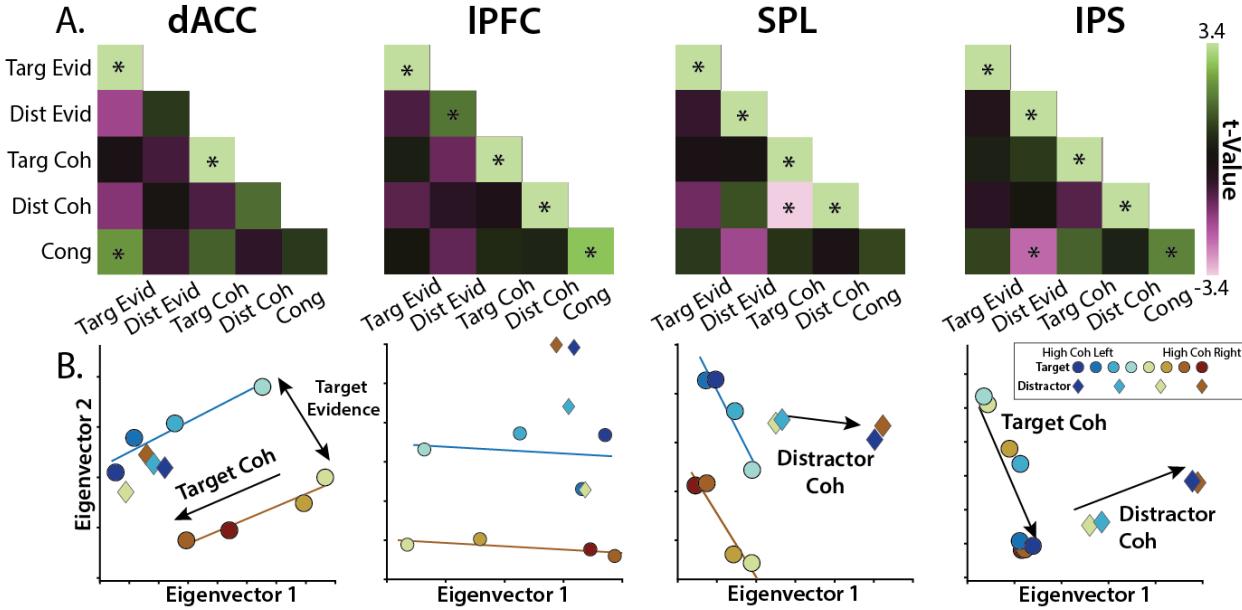


Figure 4. Region-specific feature encoding. **A)** Similarity matrices for dACC, IPFC, SPL, and IPS, correlating feature evidence ('Evid'), feature coherence ('Coh'), and feature congruence ('Cong'). Encoding strength on diagonal (right-tailed p -value), encoding alignment on off-diagonal (two-tailed p -value). **B)** Classical MDS embedding of target (circle) and distractor (diamond) representations at different levels of evidence. Colors denote coherence. GLMs: A: Feature MV, B: Evidence Levels, see Table 2.

We complemented our whole-brain analyses with ROI analyses in areas exhibiting reliable encoding of key variables, focusing on core frontal regions linked with cognitive control (dACC and lateral PFC [IPFC]), and parietal regions linked with decision-making and attention (SPL and IPS; ^{12,15}). Consistent with our analyses above, we found that target and distractor coherence encoding was aligned in SPL, but not in IPS (Figure 4a, compare to Figure 3d), whereas SPL encoded target and distractor evidence. Directly comparing these regions (see Table 1), we found stronger encoding of target evidence in SPL, stronger encoding of target coherence in IPS, and stronger alignment between target-distractor coherence alignment in SPL. Unlike our univariate results, we did not find distractor congruence encoding in dACC (though this was found in IPFC and IPS). Instead, dACC showed multivariate encoding of target coherence and evidence.

Task Feature	SPL	IPS	SPL – IPS
Target Evidence	$t_{(28)} = 10.39, p = 4.07 \times 10^{-11}, \text{logBF} = 8.37$	$t_{(28)} = 5.82, p = 3.00 \times 10^{-6}, \text{logBF} = 3.81$	$t_{(28)} = 3.89, p = .000562, \text{logBF} = 1.75$
Distractor Evidence	$t_{(28)} = 4.42, p = 1.34 \times 10^{-4}, \text{logBF} = 2.30$	$t_{(28)} = 3.62, p = 0.0012, \text{logBF} = 1.47$	$t_{(28)} = 0.896, p = .378, \text{logBF} = -0.545$
Target-Distractor Evidence Alignment	$t_{(28)} = -0.703, p = .488, \text{logBF} = -0.606$	$t_{(28)} = -0.436, p = 0.666, \text{logBF} = -0.667$	$t_{(28)} = -0.145, p = .886, \text{logBF} = -0.701$
Target Coherence	$t_{(28)} = 5.82, p = 2.94 \times 10^{-6}, \text{logBF} = 3.82$	$t_{(28)} = 7.73, p = 2.00 \times 10^{-8}, \text{logBF} = 5.83$	$t_{(28)} = -3.89, p = 9.36 \times 10^{-9}, \text{logBF} = 6.14$

Distractor Coherence	$t_{(28)} = 8.88, p = 1.25 \times 10^{-9}$, logBF = 6.97	$t_{(28)} = 8.53, p = 2.80 \times 10^{-9}$, logBF = 6.63	$t_{(28)} = 1.40, p = .170$, logBF = -0.320
Target-Distractor Coherence Alignment	$t_{(28)} = -4.75, p = 5.50 \times 10^{-5}$, logBF = 2.65	$t_{(28)} = -1.06, p = 0.294$, logBF = -0.479	$t_{(28)} = -2.99, p = .00580$, logBF = 0.861

Table 1. Feature encoding contrasted across parietal cortex. Encoding of feature evidence and coherence within SPL, within IPS, and contrasted between SPL and IPS. Note the stronger target evidence encoding in SPL, stronger target coherence encoding in IPS, and stronger target-distractor coherence alignment in SPL.

To further characterize how feature coherence and evidence are encoded across these regions, we performed multidimensional scaling over each regions task representations (Figure 4b; ^{64,70}). Briefly, this method allows us to visualize – in a non-parametric manner – the relationships between representations of different feature levels (e.g., levels of target coherence), by estimating each feature level separately within a GLM and then using singular value decomposition to project these patterns into a 2D space (see Methods for additional details). We found that coherence and evidence axes naturally emerge in the top two principal components in this analysis within dACC, SPL, and IPS. Coherence axes (light to dark shading) are parallel between left (blue) and right (brown) responses, suggesting a response-independent encoding. In these components, evidence encoding appeared to be binary, in contrast to parametric coherence encoding (we found similar whole-brain encoding maps for binary-coded evidence; see Supplementary Figure 7). Critically, whereas coherence encoding axes within SPL were aligned between targets (circles) and distractors (diamonds; confirming aligned encoding), in IPS these representations form perpendicular lines (confirming orthogonal encoding). When we visualized higher dimensions, we found that IPS did appear to have weak encoding alignment between target and distractor coherence in higher dimensions (Supplementary Figure 8). Nevertheless, the orthogonal encoding in the first two principal components is sufficient for a downstream region to have an independent read-out of feature-specific coherence. These analyses both help to visualize cross-region dissociations in encoding profiles and validate that task features are encoded in a monotonic fashion.

Finally, to explore the divisions between SVA and Control networks evident in the univariate analyses, we split up our two prefrontal ROIs by their network membership (Supplementary Figure 9). In dACC, we found that SVA parcels tended to have stronger feature encoding than Control parcels. Interestingly, in these SVA parcels several features were aligned with the target evidence dimension, consistent with recent human electrophysiology findings ³⁵. In IPFC, we found that Control parcels, but not SVA parcels, encoded distractor congruence (Control: $t_{(28)} = 3.60$, two-tailed $p = .0012$, logBF = 1.45; SVA: $t_{(28)} = 0.57, p = .57$, logBF = -0.64; Control – SVA: $t_{(28)} = 3.27, p = .0029$, logBF = 1.12). This distractor congruence encoding was present in IPFC both in ‘Control A/B’ parcels ($t_{(28)} = 3.66, p = .001$) and marginally in ‘Control C’ parcels ($t_{(28)} = 1.86, p = .073$). This network-selective encoding of congruence is consistent with the univariate results in dACC (see Figure 2).

Control demands dissociate coherence and evidence encoding

Our findings thus far demonstrate two sets of dissociations within and across brain regions. In dACC, we find that distinct regions encode the control demands related to discriminating targets

1 (caudal dACC) versus overcoming distractor incongruence (rostral dACC). In posterior parietal
2 cortex, we find that overlapping regions track the coherence of these two stimulus features, but
3 that distinct regions represent these features in alignment (SPL) versus orthogonally (IPS). While
4 these findings suggest that this set of regions was involved in translating between feature
5 information and goal-directed responding, they only focus on the information that was presented
6 to the participant on a given trial. To provide a more direct link between feature-specific
7 encoding and control, we examined how the encoding of feature coherence differed between
8 matched task that placed stronger or weaker demands on cognitive control. So far, our analyses
9 have focused on conditions in which participants needed to respond to the color feature while
10 ignoring the motion feature (Attend-Color task), but on alternating scanner runs participants
11 instead responded to the motion dimension and ignored the color dimension (Attend-Motion
12 task). These tasks were matched in their visual properties (identical stimuli) and motor outputs
13 (left/right responses), but critically differed in their control demands. Attend-Motion was
14 designed to be much easier than Attend-Color, as the left/right motion directions are compatible
15 with the left/right response directions (i.e., Simon facilitation; ^{3,10}). Comparing these tasks allows
16 us to disambiguate bottom-up attentional salience from the top-down contributions to attentional
17 priority ^{47,71–73}.
18

19 Consistent with previous work ¹⁰, performance on the Attend-Motion task was better overall
20 (mean RT: 565ms vs 725ms, sign-rank $p = 2.56 \times 10^{-6}$; mean Accuracy: 93.7% vs 87.5%, sign-
21 rank $p = .000318$). Unlike the Attend-Color task, performance was not impaired by distractor
22 incongruence (i.e., color distractors; RT: $t_{(28)} = -1.39$, $p = .176$; Accuracy: $t_{(28)} = 0.674$, $p = .506$).
23 To investigate these task-dependent feature representations, we fit a GLM that included both
24 tasks. To control for performance differences across tasks, we only analyzed accurate trials and
25 included trial-wise RT as a nuisance covariate, concatenating RT across tasks.
26

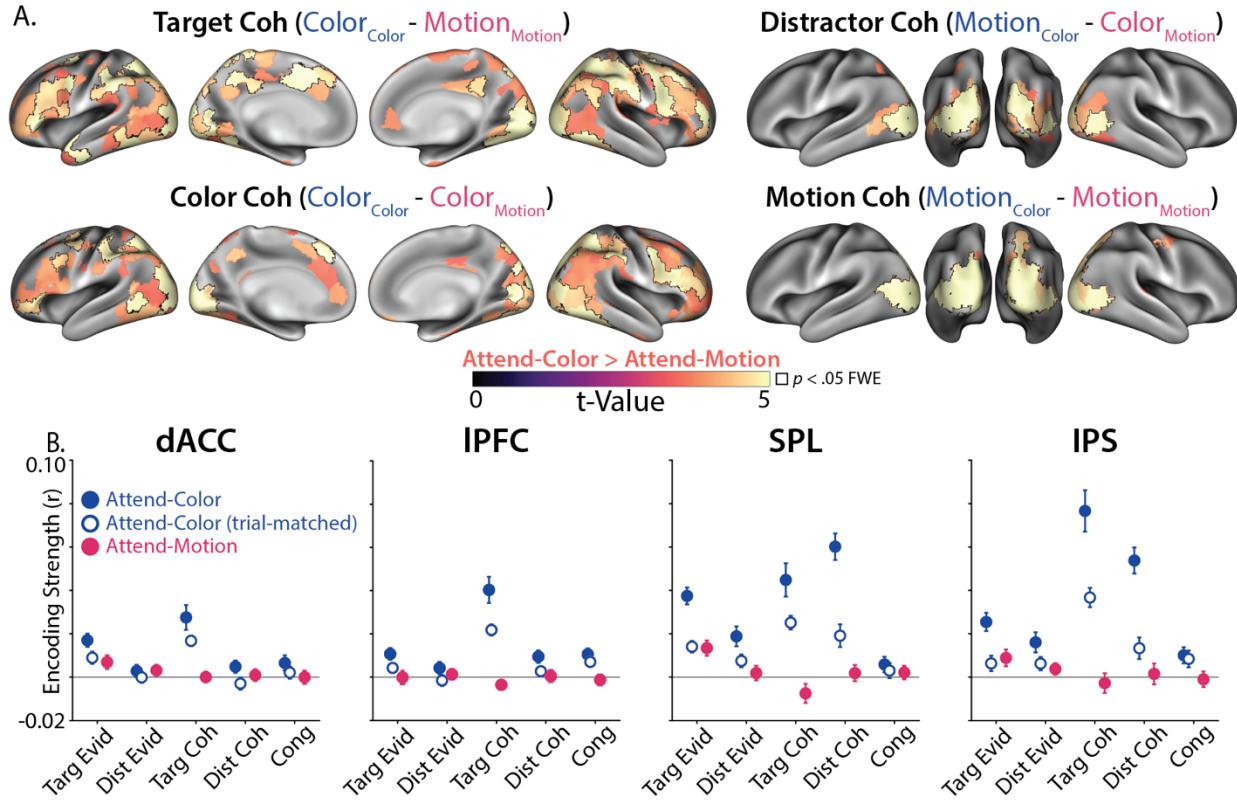


Figure 5. Task-dependent encoding strength. **A)** Across cortex, feature coherence encoding was stronger during Attend-Color than Attend-Motion, matched for the same number of trials. Attend-Color had stronger encoding when comparing target coherence (top left), distractor coherence (top right), color coherence (bottom left) and motion coherence (bottom right). Parcels are thresholded at $p < .001$ (uncorrected); outlined parcels are significant at $p < .05$ I (max-statistic randomization test across all parcels). Condition labels in title parentheses are coded ‘Feature_{Task}’. **B)** Target and distractor coherence information was encoded more strongly during Attend-Color than Attend-Motion in dACC, IPFC, SPL and IPS. Attend-Color encoding plotted from the whole sample (blue fill) and a trial-matched sample (first 45 trials of each run; white fill). In Attend-Motion runs, only target evidence was significantly encoded (magenta). **C)** Target and distractor coherence was not reliably encoded during the Attend-Motion task (liberally thresholded at $p < .01$ uncorrected). GLM: Between-Task, see Table 2.

Whereas the encoding of both color and motion coherence was widespread during the Attend-Color task (Figure 3), coherence encoding was consistently weaker during the less demanding Attend-Motion task (Figure 5A). Coherence encoding was weaker during Attend-Motion whether classifying according to goal-relevance (comparing targets or distractors) or the features themselves (comparing motion or color). Task-relevant ROIs revealed that coherence encoding was effectively absent during the easy Attend-Motion task (Figure 5B), suggesting that they depend on the control demands of the Attend-Color task^{47,74}.

In contrast to these stark task-related differences in coherence encoding, we found that neural encoding of the target evidence (color evidence in the Attend-Color task and motion evidence in the Attend-Motion task) was preserved across tasks, including within dACC, IPFC, SPL, and IPS (Figure 5B). Consistent with previous experiments examining context-dependent decision-making^{36,41,42,45,73,75,76}, we found stronger target evidence encoding relative to distractor evidence encoding, in our case in the evidence-encoding SPL (Attend-Color: $t_{(28)} = 4.26$, one-tailed $p = 0.0001$; Attend-Motion: $t_{(28)} = 2.37$, one-tailed $p = 0.0124$). We also found that target

1 evidence encoding during Attend-Motion was aligned with Attend-Color, both for *motion*
2 evidence encoding ('stimulus axis'; SPL: one-tailed $p = .0236$, IPS: one-tailed $p = .0166$) and
3 *target* evidence encoding ('decision axis'; SPL: one-tailed $p = 1.29 \times 10^{-6}$; IPS: one-tailed $p =$
4 $.0005$), again in agreement with these previous experiments. Whereas our experiment replicates
5 previous observations of the neural representations supporting contextual decision-making, we
6 now extended these findings to understand how putative attention signals (i.e., feature
7 coherence) are encoded in response to the asymmetric inference that is characteristic of cognitive
8 control⁷⁷.

9 **Aligned encoding dimensions for feature coherence and task
10 performance**

11 Feature coherence encoding (i.e., feature strength, regardless of response or congruence) depends
12 on task demands, consistent with a role in cognitive control. To further understand this
13 relationship between coherence encoding and control, we next explored how coherence encoding
14 was related to task performance. We tested this question by determining whether feature
15 coherence representations were aligned with performance representations (i.e., alignment
16 between stimulus and behavioral subspaces⁷⁸). Specifically, we included trial-level reaction time
17 and accuracy in our first-level GLMs. Encoding of performance was itself highly robust: most
18 parcels encoded reaction time and accuracy, with the strongest encoding in cognitive control
19 regions (Supplementary Figure 10). Across cortex, reaction time and accuracy were negatively
20 correlated, again most prominently across the cognitive control network. To explore the
21 behavioral relevance of coherence representations, we tested whether coherence encoding was
22 aligned with the voxel patterns encoding task performance.

23
24
25

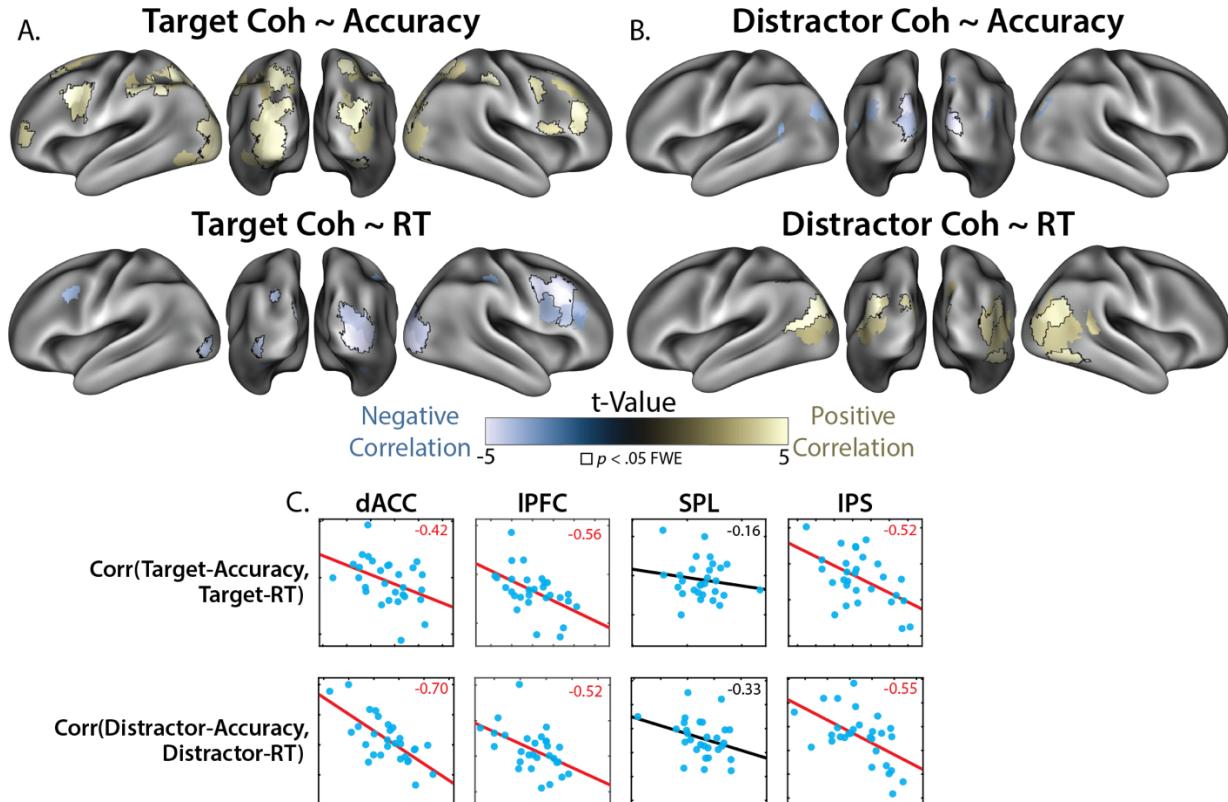


Figure 6. Alignment between feature and performance encoding. A) Alignment between encoding of target coherence and performance (top row: Accuracy, bottom row: RT). B) Alignment between encoding of distractor coherence and performance (top row: Accuracy, bottom row: RT). Across A and B, parcels are thresholded at $p < .001$ (uncorrected, in jointly reliable parcels), and outlined parcels are significant at $p < .05$ I (max-statistic randomization test across jointly reliable parcels). C) Individual differences in feature-RT alignment correlated with feature-accuracy alignment across regions (correlation values in top right; $p < .05$ in red). See Supplementary Table 1 for partial correlations controlling for reliability. GLM: Performance, see Table 2.

We found that the encoding of target and distractor coherence was aligned with performance across frontoparietal and visual regions (Figure 6a-b). If a regions' encoding of target coherence reflects how sensitive the participant was to target information on that trial (e.g., due to top-down priority), we would expect target encoding to be positively aligned with performance on a given trial, such that stronger target coherence encoding is associated with better performance and weaker target coherence encoding is associated with poorer performance. We would also expect distractor encoding to demonstrate the opposite pattern – stronger encoding associated with poorer performance and weaker encoding associated with better performance. We found evidence for both patterns of feature-performance alignment across visual and frontoparietal cortex: target encoding was aligned with better performance (faster RTs and higher accuracy; Figure 6a), whereas distractor encoding was aligned with worse performance (slower RTs and lower accuracy; Figure 6b).

Next, we examined whether performance-coherence alignment reflected individual differences in participants' task performance in our main task-related ROIs (see Figures 3-4). In particular, we tested whether the alignment between features and behavior reflects specific relationships with speed or accuracy, or whether they reflected overall increases in evidence accumulation (e.g., faster responding and higher accuracy). Within each ROI, we correlated feature-RT alignment

1 with feature-accuracy alignment across subjects. We found that in dACC and IPS, participants
2 showed the negative correlation between performance alignment measures predicted by an
3 increase in processing speed (Figure 6c). People with stronger alignment between target
4 coherence and shorter RTs tended to have stronger alignment between target coherence and
5 higher accuracy, with the opposite found for distractors. While these between-participant
6 correlations were present within targets and distractors, we did not find any significant
7 correlations across features (between-feature: all $p > .10$), again consistent with feature-specific
8 processing. These analyses were qualitatively similar after partialing out the reliability of
9 coherence and performance encoding, albeit with dACC and IPFC now showing marginal
10 correlations for target coherence (see Supplementary Table 1). While between-participant
11 analyses using small sample sizes warrant a note of caution, these findings are consistent across
12 features and regions. In conjunction with our within-participant evidence that feature coherence
13 representations are aligned with performance efficiency, these findings support a role for
14 coherence encoding in adaptive control.

15 Coherence encoding aligns with frontoparietal activity

16 Across frontal, parietal, and visual cortex, encoding of target and distractor coherence depended
17 on task demands and was aligned with performance. Since this widespread encoding of task
18 information likely reflects distributed network involvement in cognitive control ^{77,79,80}, we sought
19 to understand how frontal and parietal systems interact. We focused our analyses on IPS and
20 lateral PFC (IPFC), linking the core parietal site of orthogonal coherence encoding (IPS) to an
21 prefrontal site previous work suggests provides top-down feedback during cognitive control
22 ^{58,79,81,82}. Previous work has found that IPS attentional biases lower-level stimulus encoding in
23 visual cortices ^{83,84}, and that IPS mediates directed connectivity between IPFC and visual cortex
24 during perceptual decision-making ⁴². Here, we extended these experiments to test how IPS
25 mediates the relationship between prefrontal feedback and stimulus encoding.

26 To investigate these putative cortical interactions, we developed a novel multivariate
27 connectivity analysis to test whether coherence encoding was aligned with prefrontal activity,
28 and whether this IPFC-coherence alignment was mediated by IPS. We first estimated the voxel-
29 averaged residual timeseries in IPFC (SPM12's eigenvariate), and then included this residual
30 timeseries alongside task predictors in a whole-brain regression analysis (Supplementary Figure
31 11). This analysis can be schematized as:

$$\beta_{seed} = GLM(X, Y_{seed}) \quad (1)$$

$$e_{seed} = PCA(Y_{seed} - X\beta_{seed}) \quad (2)$$

$$\beta_{all} = GLM([X, e_{seed}], Y_{all}) \quad (3)$$

33
34 The GLM function performs regression using design matrix X and multivariate voxel timeseries
35 Y, and the PCA function extracts the first principal component of the residuals. Finally, we used
36 EGA to test whether there was alignment between patterns encoding IPFC functional
37 connectivity (i.e., betas from the residual timeseries predictor e_seed) and patterns encoding
38 target and distractor coherence. Note that while these results reflect functional connectivity, all
39 correlational measures are subject to potential confounding ⁸⁵.

40
41
42
43
44

1 We found that lPFC connectivity patterns were aligned with coherence-encoding patterns in
2 visual cortex (Figure 7A). Stronger prefrontal functional connectivity was aligned with weaker
3 target coherence and stronger distractor coherence, consistent with prefrontal recruitment during
4 difficult trials. Notably, IPS connectivity was also aligned with target and distractor coherence in
5 overlapping parcels, even when controlling for lPFC connectivity. These effects were liberally
6 thresholded for visualization, as significant direct and indirect effects are not necessary for
7 significant mediation⁸⁶.

8
9 Our critical test was whether IPS mediated the relationship between lPFC activity and coherence
10 encoding. We compared regression estimates between a model that only included lPFC residuals
11 ('solo' model) to a model that included both lPFC and IPS residuals ('both' model). Comparing
12 the strength of lPFC-coherence alignment with and without IPS is a test of whether parietal
13 cortex mediates lPFC-coherence alignment (MacKinnon et al., 2007). These models can be
14 schematized as:

$$\beta_{solo} = GLM([X, e_{lPFC}], Y_{all}) \quad (4)$$

$$\beta_{both} = GLM([X, e_{lPFC}, e_{IPS}], Y_{all}) \quad (5)$$

15
16 We found that this mediation was strongest in early visual cortex, where the alignment between
17 lPFC and feature coherence was reduced in a model that included IPS relative to a model without
18 IPS (Figure 7B). The negatively correlated target-lPFC relationship became more positive when
19 IPS was included (top), and the positively correlated distractor-PFC relationship became more
20 negative when IPS was included (bottom). Critically, we found that IPS reduced prefrontal-
21 coherence alignment in early visual cortex more than lPFC reduced parietal-coherence alignment
22 (Figure 7B inset; Supplementary Figure 12A-B), consistent with frontal-to-parietal directed
23 connectivity in previous research^{42,81}. Looking within color- and motion-sensitive parcels,
24 determined using task-free localizer runs (see Methods), we found this mediation was robustly
25 significant in color-sensitive cortex and marginally significant in motion-sensitive cortex. The
26 opposite relationship, lPFC mediation of IPS connectivity, appeared in higher-level visual cortex
27 for distractor coherence (Supplementary Figure 12C-D), though these effects were not reliable in
28 explicit contrasts and may reflect projections from both regions. Note that we did not see any
29 significant mediation of first-order target or distractor coherence encoding by IPS.
30
31
32
33

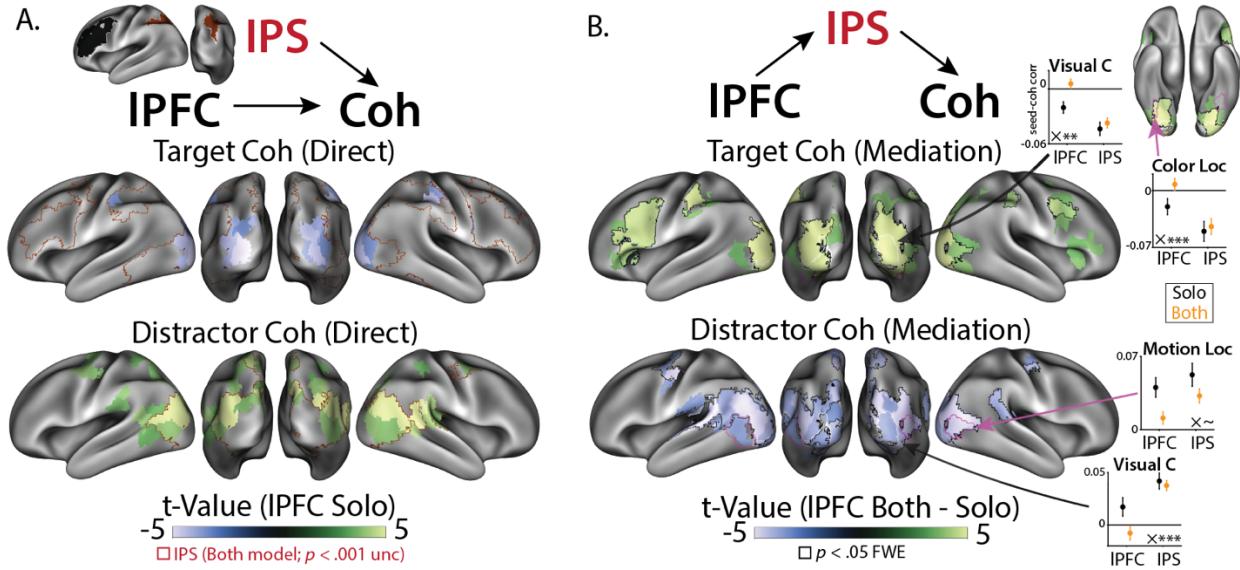


Figure 7. IPS mediates alignment between IPFC and feature encoding. **A)** Connectivity patterns from IPFC (color) and IPS (red outline) were aligned with target and distractor coherence patterns ($p < .001$ uncorrected, in jointly reliable parcels). IPS effects are outlined to show overlap, with all effects in a consistent direction to IPFC. **B)** IPFC-feature alignment contrasted between IPFC-only model ('Solo') and IPFC + IPS model ('Both'). Including IPS reduced the alignment between IPFC and feature encoding (compare the sign of the main effect in A to the contrast in B). Parcels are thresholded at $p < .001$ (uncorrected, jointly reliable parcels), and outlined parcels are significant at $p < .05$ I (max-statistic randomization test across jointly reliable parcels). Insets graphs: seed-coherence alignment in Solo models (black) and Both model (orange) across visual regions. 'Visual C' is defined by our parcellation⁵⁴; Color and Motion localizers are parcels near the peak response identified during feature localizer runs (see Methods). In general, IPFC alignment was more affected by IPS than IPS alignment was affected by IPFC (inset 'X': difference of differences; $\sim p < .10$, * $p < .01$, *** $p < .001$). GLM: Performance CX, see Table 2.

While we were primarily interested in alignment with IPFC due to previous work implicating these regions in top-down control (for reviews, see^{12,87}), for completeness we also examined how different subnetworks in both IPFC and dACC aligned with coherence encoding. In IPFC, we found that SVA and Control subnetworks had similar patterns of alignment (Supplementary Figure 13). In dACC we found that the SVA subnetwork had a qualitatively similar profile of coherence alignment as IPFC, but this alignment was absent in the Control subnetwork. Whereas this seed-coherence alignment was similar across IPFC and SVA dACC, unlike IPFC we found that SVA dACC failed to demonstrate strong evidence for mediation by IPS (Supplementary Figure 14).

A final set of analyses examined whether SPL and IPS demonstrated different patterns of task-related functional connectivity with other regions, given that we found that these regions differentially encoded evidence and coherence. When seeding our connectivity analyses with SPL activity, we found that SPL activity aligned with evidence encoding in bilateral motor cortex (Supplementary Figure 15). In contrast, IPS activity did not significantly align with evidence encoding, and this seed-evidence alignment in motor cortex was stronger for SPL than IPS, consistent with a putative role for SPL in response selection⁶⁸.

1

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

1 Discussion

2 In this experiment, we explored whether neural control systems use representations with the
3 same dimensionality as the processes they regulate^{2,5,11}. Inspired by behavioral evidence that
4 participants can independently control their sensitivity to targets and distractors¹⁰, we set out to
5 understand whether the neural correlates of monitoring and prioritization leverage independent
6 encoding for feature-selective control (Figures 1a-c). We found that key nodes of canonical
7 cognitive control networks had orthogonal neural representations of targets and distractors.

8 Within dACC, orthogonal representations of target and distractor difficulty arose from
9 segregated encoding along a rostrocaudal axis. Within IPS, orthogonal representations of target
10 and distractor coherence arose from orthogonal subspaces in overlapping voxels. Consistent with
11 a role in attentional priority, coherence representations depended on control demands, task
12 performance, and frontoparietal activity. Together, these results reveal a neural mechanism for
13 how cognitive control prioritizes multiple streams of information during decision-making.

14

15 Neurocomputational theories have proposed that dACC is involved in planning control across
16 multiple levels of abstraction^{12,88–90}. Past work has found that control abstraction is
17 hierarchically organized along dACC’s rostrocaudal axis, with more caudal dACC involved in
18 lower-level action control, and more rostral dACC involved in higher-level strategy control^{30–}
19 ^{32,34}, an organization that may reflect a more general hierarchy of abstraction within PFC^{31,91–93}.

20 Consistent with this account, we found that caudal dACC tracked the coherence of the target and
21 distractor dimensions, especially within the SVA network. In contrast, more rostral dACC
22 tracked incongruence between targets and distractors, especially within the Control network.
23 Speculatively, our results are consistent with caudal dACC tracking the first-order difficulty
24 arising from the relative salience of feature-specific information, and more rostral dACC
25 tracking the second-order difficulty arising from cross-feature (in)compatibility⁹², the latter of
26 which may require additional disengagement from distractor-dependent attentional capture.

27

28 Whereas dACC encoded feature difficulty (e.g., distractor incongruence), in parietal cortex we
29 found overlapping representations of feature coherence (e.g., distractor coherence). In SPL,
30 features had correlated coherence encoding (similarly representing low target coherence and high
31 distractor coherence), consistent with this region’s transient and non-selective role in attentional
32 control^{94–99}. In contrast, IPS had orthogonal representations of feature coherence, consistent with
33 selective prioritization of task-relevant information^{47,71–73,81,83,94–96,99,100}. While IPS primarily
34 encoded features orthogonally (i.e., in the largest components of our multidimensional scaling
35 analysis), the total coherence across features could also be read out at higher dimensions. The
36 ability of IPS to communicate both orthogonal and aligned coherence representations is
37 consistent with the diverse roles of IPS in attentional control.

38

39 Our previous work has demonstrated behavioral evidence for independent control over target and
40 distractor attentional priority in this task¹⁰, with different task variables selectively enhancing
41 target or distractor sensitivity (see also^{4,101}). Orthogonal feature representation in IPS may offer
42 a mechanism for this feature-selective control, consistent with theoretical accounts of IPS
43 implementing a priority map that combines stimulus- or value-dependent salience with goal-
44 dependent feedback from PFC^{17,57,58,80,102}.

45

1 In dACC, we found that target and distractor difficulty encoding was consistent with the
2 segregated encoding hypothesis, with features evoking univariate responses in distinct but
3 adjacent regions. Interestingly, we did not find corresponding encoding of distractor congruence
4 in our multivariate analyses within dACC, potentially reflecting the spatial smoothness of this
5 response. However, we did find multivariate encoding of distractor congruence in IPFC, and
6 multivariate encoding of target and distractor coherence in IPS. These multivariate profiles were
7 consistent with our subspace encoding hypothesis. The reasons for a mix of segregated and
8 subspace encoding across cortex is unclear, but this may speculatively reflect the segregation
9 across functional networks. Like in dACC, distractor congruence had stronger encoding within
10 the IPFC Control network, albeit without the feature segregation (IPFC Control parcels also
11 encoded target coherence in an orthogonal subspace). It is possible that these network
12 segregations help bind related control processes^{15,18,80}, a hypothesis that future experiments
13 should test with targeted paradigms (e.g., with subject-specific functional networks).

14
15 By comparing two different task goals (Attend-Color vs. Attend-Motion), our study was able to
16 test whether coherence representations reflect control-dependent prioritization of information
17 processing. Previous research has shown that these tasks differ dramatically in their control
18 demands¹⁰. As in previous work, task performance was much better in Attend-Motion runs than
19 Attend-Color runs, and participants were not sensitive to color distractors. Consistent with
20 previous work on context-dependent decision-making, target evidence had similarly strong
21 encoding across tasks, with generalizable encoding dimensions for choice and motion directions
22^{36,41,45}. In contrast to these putative decision representations, we found that coherence
23 representations disappeared in the easier Attend-Motion task. On its own, weaker encoding of
24 color distractors in Attend-Motion could be explained by the weaker bottom-up salience of the
25 color dimension. However, the stark drop in the encoding of target (motion) coherence in these
26 blocks cannot be similarly accounted for – these differences in target coherence encoding
27 showed the opposite relationship expected from salience: better encoding of low-salience color
28 targets (hard Attend-Color task) and weaker encoding of high-salience motion targets (easy
29 Attend-Motion task). Instead, this encoding profile is consistent with previous research finding
30 that feature decoding is stronger for more difficult tasks^{47,71,72,103} or when people are
31 incentivized to use cognitive control^{104,105}.

32
33 Critically, stimuli and responses were matched across tasks, helping to rule out alternative
34 accounts of coherence encoding based on ‘bottom-up’ stimulus salience, decision-making, or eye
35 movements. Difficulty-dependent coherence encoding may instead reflect the involvement of an
36 attention control system that can separately regulate target and distractor processing,
37 speculatively indexing the top-down ‘gain’ or ‘priority’ on these features^{17,58,102}. Supporting this
38 account, coherence representations in cognitive control regions like IPS were aligned with
39 performance representations, with target encoding strength aligned with better performance and
40 distractor encoding strength aligned with poorer performance. Individual difference in feature-
41 performance alignment was correlated across features, consistent with these representations
42 reflecting the underlying processes (e.g., priority) that give rise to behavior, rather than
43 performance monitoring or surprise (which would likely have the opposite relationship, e.g., high
44 target coherence aligned with poorer performance).

45

1 Classic models of prefrontal involvement in cognitive control^{77,82,106} propose that prefrontal
2 cortex biases information processing in sensory regions. In line with this macro-scale
3 organization, we found that coherence encoding in visual cortex was related to functional
4 connectivity with the frontoparietal network. In particular, coherence encoding in visual cortex
5 was aligned with patterns of functional connectivity to lateral prefrontal cortex, and this feature-
6 seed relationship was mediated by IPS. The results of this novel multivariate connectivity
7 analysis are consistent with previous research supporting a role for IPS in top-down control of
8 visual encoding^{83,84,107}, as well as a granger-causal PFC-IPS-visual pathway during a similar
9 decision-making task⁴². Here, we demonstrate stable ‘communication subspaces’ between visual
10 cortex and PFC^{108,109}, which can plausibly communicate feedback adjustments to feature gain.
11 With that said, while our interpretation of the direction of communication is therefore supported
12 by prior work, these connectivity methods are correlational⁸⁵, and cannot rule out the possibility
13 that our mediation findings reflect a bottom-up pattern of communication (e.g., visual-IPS-PFC).
14 The asymmetric mediation between regions (i.e., IPS mediates IPFC more than IPFC mediates
15 IPS; Supplementary Figure 12) rules out a range of potential confounders, and these regions
16 were selected based on the anatomical connectivity within the frontoparietal network, notably
17 through the superior longitudinal fasciculus¹¹⁰. Future research should use temporally precise
18 neuroimaging to account for directionality, causal manipulations to account for causality (e.g.,
19¹¹¹), and should explore the higher dimensional connectivity subspaces that link different regions
20^{103,109}. These considerations notwithstanding, our findings are consistent with IPS, a critical site
21 for orthogonal feature representations, playing a key role in linking prefrontal cortex with early
22 perceptual processing.

23
24 Collectively, our findings provide new insights into how the brain may control multiple streams
25 of information processing. While evidence for multivariate control has a long history in
26 attentional tracking^{28,112}, including parametric relationships between attentional load and IPS
27 activity^{113–117}, little is known about how the brain coordinates multiple control signals^{2,5}. Future
28 experiments should further elaborate on this frontoparietal control circuit, for instance by
29 interrogating how incentives influence different task representations^{104,105,118–120}, or how neural
30 and behavioral indices of control causally depend on perturbations of neural activity¹¹¹. Future
31 experiments should also use fast timescale neural recording technologies like (i)EEG or (OP-
32)MEG to better understand the within-trial dynamics of multivariate control^{10,121}. In sum, this
33 experiment provides new insights into the large-scale neural networks involved in multivariate
34 cognitive control, and points towards new avenues for developing a richer understanding of goal-
35 directed attention.

36

1 **Methods**

2 **Participants**

3 Twenty-nine individuals (17 females, Age: M = 21.2, SD = 3.4) participated in this experiment.
4 All participants had self-reported normal color vision and no history of neurological disorders.
5 Two participants missed one Attend-Color block (see below) due to a scanner removal, and one
6 participant missed a motion localizer due to a technical failure, but all participants were retained
7 for analysis. Participants provided informed consent, in accordance with Brown University's
8 institutional review board.

9 **Task**

10 The main task closely followed our previously reported behavioral experiment ¹⁰. On each trial,
11 participants saw a random dot kinematogram (RDK) against a black background. This RDK
12 consisted of colored dots that moved left or right, and participants responded to the stimulus with
13 button presses using their left or right thumbs.

14 In Attend-Color blocks (six blocks of 150 trials), participants responded depending on which
15 color was in the majority. Two colors were mapped to each response (four colors total), and dots
16 were a mixture of one color from each possible response. Dots colors were approximately
17 isolument (uncalibrated; RGB: [239, 143, 143], [191, 239, 143], [143, 239, 239], [191, 143,
18 239]), and we counterbalanced their assignment to responses across participants.

20 In Attend-Motion blocks (six blocks of 45 trials), participants responded based on the dot motion
21 instead of the dot color. Dot motion consisted of a mixture between dots moving coherently
22 (either left or right) and dots moving in a random direction. Attend-Motion blocks were shorter
23 because they acted to reinforce motion sensitivity and provide a test of stimulus-dependent
24 effects.

26 Critically, dots always had color and motion, and we varied the strength of color coherence
27 (percentage of dots in the majority) and motion coherence (percentage of dots moving
28 coherently) across trials. Our previous experiments have found that in Attend-Color blocks,
29 participants are still influenced by motion information, introducing a response conflict when
30 color and motion are associated with different responses ¹⁰. Target coherence (e.g., color
31 coherence during Attend-Color) was linearly spaced between 65% and 95% with 5 levels, and
32 distractor congruence (signed coherence relative to the target response) was linearly spaced
33 between -95% and 95% with 5 levels. In order to increase the salience of the motion dimension
34 relative to the color dimension, the display was large (~10 degrees of visual angle) and dots
35 moved quickly (~10 degrees of visual angle per second).

37 Participants had 1.5 seconds from the onset of the stimulus to make their response, and the RDK
38 stayed on the screen for this full duration to avoid confusing reaction time and visual stimulation
39 (the fixation cross changed from white to gray to register the response). The inter-trial interval
40 was uniformly sampled from 1.0, 1.5, or 2.0 seconds. This ITI was relatively short in order to

1 maximize the behavioral effect, and because efficiency simulations showed that it increased
2 power to detect parametric effects of target and distractor coherence (e.g., relative to a more
3 standard 5 second ITI). The fixation cross changed from gray to white for the last 0.5 seconds
4 before the stimulus to provide an alerting cue.

5 Procedure

6 Before the scanning session, participants provided consent and practiced the task in a mock MRI
7 scanner. First, participants learned to associate four colors with two button presses (two colors
8 for each response). After being instructed on the color-button mappings, participants practiced
9 the task with feedback (correct, error, or 1.5 second time-out). Errors or time-out feedback were
10 accompanied with a diagram of the color-button mappings. Participants performed 50 trials with
11 full color coherence, and then 50 trials with variable color coherence, all with 0% motion
12 coherence. Next, participants practiced the motion task. After being shown the motion mappings,
13 participants performed 50 trials with full motion coherence, and then 50 trials with variable
14 motion coherence, all with 0% color coherence. Finally, participants practiced 20 trials of the
15 Attend-Color task and 20 trials of Attend-Motion tasks with variable color and motion coherence
16 (same as scanner task).

17
18 Following the twelve blocks of the scanner task, participants underwent localizers for color and
19 motion, based on the tasks used in our previous experiments ³⁰. Both localizers were block
20 designs, alternating between 16 seconds of feature present and 16 seconds of feature absent for
21 seven cycles. For the color localizer, participants saw an aperture the same size as the task, either
22 filled with colored squares that were resampled every second during stimulus-on ('Mondrian
23 stimulus'), or luminance-matched gray squares that were similarly resampled during stimulus-
24 off. For the motion localizer, participants saw white dots that were moving with full coherence in
25 a different direction every second during stimulus-on, or still dots for stimulus-off. No responses
26 were required during the localizers.

27 MRI sequence

28 We scanned participants with a Siemens Prisma 3T MR system. We used the following sequence
29 parameters for our functional runs: field of view (FOV) = 211 mm × 211 mm (60 slices), voxel
30 size = 2.4 mm³, repetition time (TR) = 1.2 sec with interleaved multiband acquisitions
31 (acceleration factor 4), echo time (TE) = 33 ms, and flip angle (FA) = 62°. Slices were acquired
32 anterior to posterior, with an auto-aligned slice orientation tilted 15° relative to the AC/PC plane.
33 At the start of the imaging session, we collected a high-resolution structural MPRAGE with the
34 following sequence parameters: FOV = 205 mm × 205 mm (192 slices), voxel size = 0.8 mm³,
35 TR = 2.4 sec, TE1 = 1.86 ms, TE2 = 3.78 ms, TE3 = 5.7 ms, TE4 = 7.62, and FA = 7°. At the
36 end of the scan, we collected a field map for susceptibility distortion correction (TR = 588ms,
37 TE1 = 4.92 ms, TE2 = 7.38 ms, FA = 60°).

1 fMRI preprocessing

2 We preprocessed our structural and functional data using fMRIprep (v20.2.6; ¹²² based on the
3 Nipype platform ¹²³. We used FreeSurfer and ANTs to nonlinearly register structural T1w
4 images to the MNI152NLin6Asym template (resampling to 2mm). To preprocess functional T2w
5 images, we applied susceptibility distortion correction using fMRIprep, co-registered our
6 functional images to our T1w images using FreeSurfer, and slice-time corrected to the midpoint
7 of the acquisition using AFNI. We then registered our images into MNI152NLin6Asym space
8 using the transformation that ANTs computed for the T1w images, resampling our functional
9 images in a single step. For univariate analyses, we smoothed our functional images using a
10 Gaussian kernel (8mm FWHM, as dACC responses often have a large spatial extent). For
11 multivariate analyses, we worked in the unsmoothed template space (see below).

12 fMRI univariate analyses

13 We used SPM12 (v7771) for our univariate general linear model (GLM) analyses. Due to high
14 trial-to-trial collinearity from to our short ITIs, we performed all analyses across trials, rather
15 than extracting single-trial estimates. Our regression models used whole trials as events (i.e., a
16 1.5 second boxcar aligned to the stimulus onset). We parametrically modulated these events with
17 standardized trial-level predictors (e.g., linear-coded target coherence, or contrast-coded errors),
18 and then convolved these predictors with SPM's canonical HRF, concatenating our voxel
19 timeseries across runs. We included nuisance regressors to capture 1) run intercepts and 2) the
20 average timeseries across white matter and CSF (as segmented by fMIRPrep). To reduce the
21 influence of motion artifacts, we used robust weighted least-squares ^{124,125}, a procedure for
22 optimally down-weighting noisy TRs.

23 We estimated contrast maps at the subject-level, which we then used for one-sample t-tests at the
24 group-level. We controlled for family-wise error rate using threshold-free cluster enhancement
25 ¹²⁶, testing whether voxels have an unlikely degree of clustering under a randomized null
26 distribution (Implemented in PALM ¹²⁷; 10,000 randomizations). To improve the specificity of
27 our coverage (e.g., reducing white-matter contributions) and to facilitate our inference about
28 functional networks (see below), we limited these analyses to voxels within the Kong2022
29 whole-brain parcellation ^{54,55}. This parcellation assigns regional labels to parcels (e.g., whether
30 parcels are in 'SPL' or 'IPS'), which was used through-out to generate ROIs. Surface renders
31 were generated using surfplot ^{128,129}, projecting from MNI space to the Human Connectome
32 Project's fsLR space (164,000 vertices).

34 dACC longitudinal axis analyses

35 To characterize the spatial organization of target difficulty and distractor congruence signals in
36 dACC, we constructed an analysis mask that provided broad coverage across cingulate cortex
37 and preSMA. This mask was constructed by 1) getting a meta-analytic mask of cingulate
38 responses co-occurring with 'cognitive control' (Neurosynth uniformity test; ¹³⁰, and taking any
39 parcels from the whole-brain Schaefer parcellation (400 parcels; ^{54,55} that had a 50 voxel overlap
40 with this meta-analytic mask. We used this parcellation because it provided more selective gray

1 matter coverage than the Neurosynth mask alone and it allowed us to categorize voxels
2 membership in putative functional networks.
3
4 To characterize the spatial organization within dACC, we first performed PCA on the masked
5 voxel coordinates (y and z), getting a score for each voxel's position on the longitudinal axis of
6 this ROI. We then regressed voxel's gradient scores against their regression weights from a
7 model including linear target coherence and distractor congruence (both coded -1 to 1 across
8 difficulty levels). We used linear mixed effects analysis to partially pool across subjects and
9 accommodate within-subject correlations between voxels. Our model predicted gradient score
10 from the linear and quadratic expansions of the target and distractor betas (gradientScore ~ 1 +
11 target + target² + distractor + distractor² + (1 + target + target² + distractor + distractor² |
12 subject)). To characterize the network-dependent organization of target and distractor encoding,
13 we complexity-penalized fits between models that either 1) predicted target or distractor betas
14 from linear and quadratic expansions of gradient scores, or 2) predicted target/distractor betas
15 from dummy-coded network assignment from the Schaefer parcellation, comparing these models
16 against a model that used both network and gradient information.

17 Encoding Geometry Analysis (EGA)

18 We adapted functions from the pcm-toolbox and rsatoolbox packages for our multivariate
19 analyses^{65,131}. We first fit whole-brain GLMs without spatial smoothing, separately for each
20 scanner run. These GLMs estimated the parametric relationship between task variables and
21 BOLD response (e.g., linearly coded target coherence), with a pattern of these parametric betas
22 across voxels reflecting linear encoding subspace⁵⁹. Within each Schaefer parcel (n=400), we
23 spatially pre-whitened these beta maps, reducing noise correlations between voxels that can
24 inflate pattern similarity and reduce reliability⁶³. We then computed the cross-validated Pearson
25 correlation, estimating the similarity of whitened patterns across scanner runs. We used a
26 correlation metric to estimate the alignment between encoding subspaces, rather than distances
27 between condition patterns, to normalize biases and scaling across stimuli (e.g., greater
28 sensitivity to targets vs distractors) and across time (e.g., representational drift). Note that this
29 analysis approach is related to ‘Parallelism Scores’⁴³, but here we use parametric encoding
30 models and emphasize not only deviations from parallel/orthogonal, but also the direction of
31 alignment between features (e.g., Figures 5 and 7).
32

33 We computed subspace alignment between contrasts of interest within each participant, and then
34 tested these against zero at the group level. Since our correlations were less than $r = |0.5|$, we did
35 not transform correlations before analysis. We used a Bayesian t-test to test for orthogonality
36 (bayesFactor toolbox in MATLAB, based on¹³²). The Bayes factor from this t-test gives
37 evidence for either non-orthogonality (BF_{10} further from zero) or orthogonality (BF_{10} closer to
38 zero, often defined as the reciprocal BF_{01}). Using a standard prior (Cauchy, width = 0.707), our
39 strongest possible evidence for the orthogonality is $BF_{01} = 5.07$ or equivalently $\log BF = -0.705$
40 (i.e., the Bayes factor when $t_{(28)} = 0$).
41

42 Our measure of encoding strength was whether encoding subspaces were reliable across blocks
43 (i.e., whether within-feature encoding pattern correlations across runs were significantly above
44 zero at the group level). We used pattern reliability as a geometric proxy for how well a linear

1 encoder would predict held-out brain data, as reliability provides the similarity between the
2 cross-validated model and the best linear unbiased estimator of the within-sample data. We
3 confirmed through simulations that pattern reliability is a good proxy for the traditional encoding
4 metric of predicting held-out timeseries⁵⁹. However, we found that pattern reliability is more
5 powerful, due to it being much less sensitive to the magnitude of residual variance (these two
6 methods are identical in the noise-free case; see Supplementary Figure 3).

7
8 When looking at alignment between two subspaces across parcels, we first selected parcels that
9 significantly encoded both factors ('jointly reliable parcels', both $p < .001$ uncorrected). This
10 selection process acts as a thresholded version of classical correlation disattenuation^{66,67}, and we
11 confirmed through simulations this selection procedure does not increase type 1 error rate. We
12 corrected for multiple comparisons using non-parametric max-statistic randomization tests across
13 parcels¹³³. These randomization tests determine the likelihood of an observed effects under a
14 null distribution generated by randomizing the sign of alignment correlations across participants
15 and parcels 10,000 times. Within each randomization, we saved the max and min group-level
16 effect sizes across all parcels, estimating the strongest parcel-wise effect we'd expect if there
17 wasn't a systematic group-level effect.

18
19 Some of our first-level models had non-zero levels of multicollinearity, due to conditioning on
20 trials without omission errors or when including feature coherence and performance in the same
21 model. Multicollinearity was far below standard thresholds for concern, generally (much) less
22 than 5 for a standard threshold is 30 (ratio between largest and smallest singular values in the
23 design matrix, using MATLAB colintest;¹³⁴). However, we wanted to confirm that predictor
24 correlations wouldn't bias our estimates of encoding alignment. We simulated data from a
25 pattern component model¹³¹ in which two variables were orthogonal (generated by separate
26 variance components with no covariance), but were generated from a design matrix with
27 correlated predictors. These simulations confirmed that cross-validated similarity measures were
28 not biased by predictor correlations (Supplementary Figure 16).

29
30 To provide further validation for our parametric analyses, we estimated encoding profiles using
31 an analysis with fewer parametric assumptions. First, we fit a GLM with separate predictors for
32 levels of target and distractor evidence ('Evidence Levels' GLM in Table 2). Next, we estimated
33 a traditional (cross-validated) representational dissimilarity matrix across all feature levels.
34 Finally, we visualized these encoding profiles using classical multidimensional scaling
35 (eigenvalue decomposition; see Figure 4B and Supplementary Figure 8).

36 Multivariate Connectivity Analysis

37 To estimate what information is plausibly communicated between cortical areas, we measured
38 the alignment between multivariate connectivity patterns (i.e., the 'communication subspace'
39 with a seed region,¹⁰⁸) and local feature encoding patterns. First, we residualized our
40 Performance GLM (see Table 2) from a seed region's timeseries, and then extracted the
41 variance-weighted average timecourse (i.e., the leading eigenvariate from SPM12's volume of
42 interest function). We then re-estimated our Performance GLM, including the block-specific
43 seed timeseries as a covariate, and performed EGA between seed and coherence patterns (see
44 Equations 1-3). We found convergent results when we residualized a quadratic expansion of our

1 Performance GLM from our seed region, helping to confirm that connectivity alignment wasn't
2 due to underfitting. Note that our cross-validated EGA helps avoid false positives due to any
3 correlations in the design matrix (see above). We localized this connectivity analysis to color-
4 and motion-sensitive cortex by finding the bilateral Kong22 parcels that roughly covered the area
5 of strongest block-level contrast during our localizer runs. Note that these analyses reflect
6 'functional connectivity', which is susceptible to unmodelled confounders⁸⁵.
7

8 To estimate the mediation of IPFC connectivity by IPS, we compared models in which just IPFC
9 or just IPS were used for EGA against a model where both seeds were included as covariates in
10 the same model⁽⁸⁶⁾; see Equations 4-5). Our test of mediation was the group-level difference in
11 IPFC seed-coherence alignment before and after including IPS. While these analyses are
12 inherently cross-sectional (i.e., IPFC and IPS are measured at the same time), we supplemented
13 these analyses by showing that the mediating effect of IPS on IPFC was much larger than the
14 mediating effect of IPFC on IPS (see Figure 7B; Supplementary Figure 12). Unlike traditional
15 mediation analyses looking at the first-order change in regression estimates, our analysis looks at
16 the second-order change in the multivariate alignment between regression estimates, using the
17 same core rationale.
18

Model Name	Trial selection	Predictors (z-scored)
Feature UV	No omission errors; run-concatenated	target coherence, distractor coherence, target evidence, distractor evidence, distractor congruence; omission errors (run-concatenated)
Difficulty Levels	No omission errors; run-concatenated	Separate levels (1,2,4,5) of target coherence, separate levels (1,2,4,5) of distractor congruence; omission errors (run-concatenated)
Feature MV	No errors; run-separated	target coherence, distractor coherence, target evidence, distractor evidence, distractor congruence; errors (run-concatenated)
Evidence Levels	No errors; run-separated	Levels (1-5, 6-10) of target evidence, Levels (1,2,4,5) of distractor evidence; errors (run-concatenated)
Between-Task	No errors; run-separated	target coherence, distractor coherence, target evidence, distractor evidence, distractor congruence; errors (run-concatenated); reaction time (run-concatenated)
Performance	No omission errors; run-separated	target coherence, distractor coherence, target evidence, distractor evidence, distractor congruence, reaction time, accuracy; omission errors (run-concatenated)
Performance CX	No omission errors; run-separated	target coherence, distractor coherence, target evidence, distractor evidence, distractor congruence, reaction time, accuracy;

		omission errors (run-concatenated); seed timeseries (run-separated)
--	--	--

1 **Table 2.** *fMRI models*. First-level general linear models used for univariate and multivariate fMRI analyses.
 2 Coherence: percentage of dots supporting the same response ('unsigned coherence'). Evidence: % dots supporting a
 3 rightwards vs leftwards response ('signed coherence'). Distractor Congruence: % dots supporting the same response
 4 as the target dimension. All predictors were z-scored within their run. For difficulty and feature levels, we included
 5 each level as a separate predictor, with collinearity with the block intercept preventing all levels from being
 6 included. For Evidence Levels, targets had greater granularity due to distractors being coded relative to targets (5
 7 levels of congruence led to 5 levels of coherence). For Performance CX, seed timeseries were included as run-
 8 separated regressors (see Multivariate Connectivity Analysis in Methods).

9

10 **Acknowledgements:** This work was supported by NIH grant R01MH124849 (A.S.), NSF
 11 CAREER Award 2046111 (A.S.), NIH grant S10OD025181, and the C.V. Starr Postdoctoral
 12 Fellowship (H.R.). We are grateful to Joonhwa Kim for her assistance in data collection, and to
 13 Michael J. Frank, Matthew N. Nassar, Jonathan Cohen, Michael Esterman, Romy Frömer, Jörn
 14 Diedrichsen, Apoorva Bhandari, Debbie Yee, Sam Nastase, Caroline Jahn, and the Shenhav Lab
 15 for helpful discussions.

16 **Conflicts of Interest:** None

17 **Data Availability:** Data will be made available upon publication.

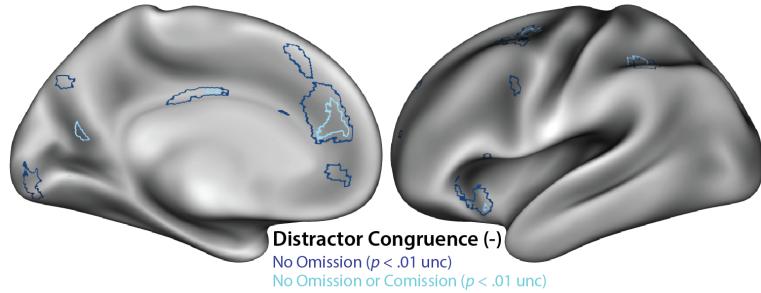
18 **Code Availability:** Code will be made available upon publication.

19

20

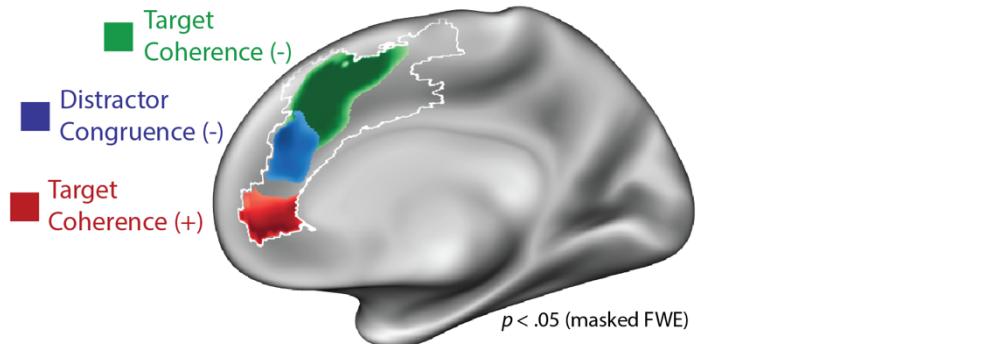
21

1 Supplementary Figures



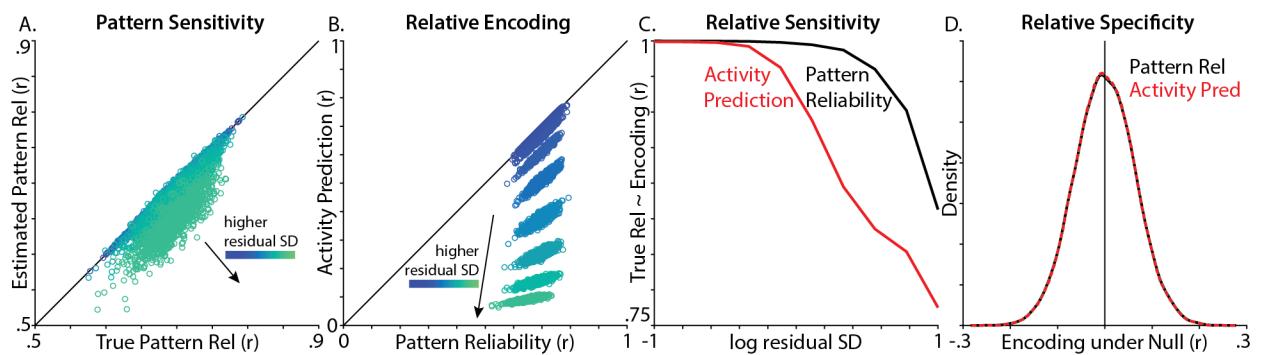
2

3 **Supplementary Figure 1.** *Error control analysis.* Distractor congruence effect when controlling for different types
4 of errors. Our primary analysis only analyzed trials without omission errors (navy), here plotted at a liberal
5 uncorrected threshold. When we analyze trials without omission errors and commission errors (cyan), we see a
6 consistent whole-brain topography, albeit at a lower statistical threshold. In both cases, relevant errors trials were
7 included as nuisance events.
8



9
10
11
12

Supplementary Figure 2. *Target ease.* Parametric effects of target coherence and distractor congruence, showing
the rostral effect of target ease (positive relationship with target coherence) in red.

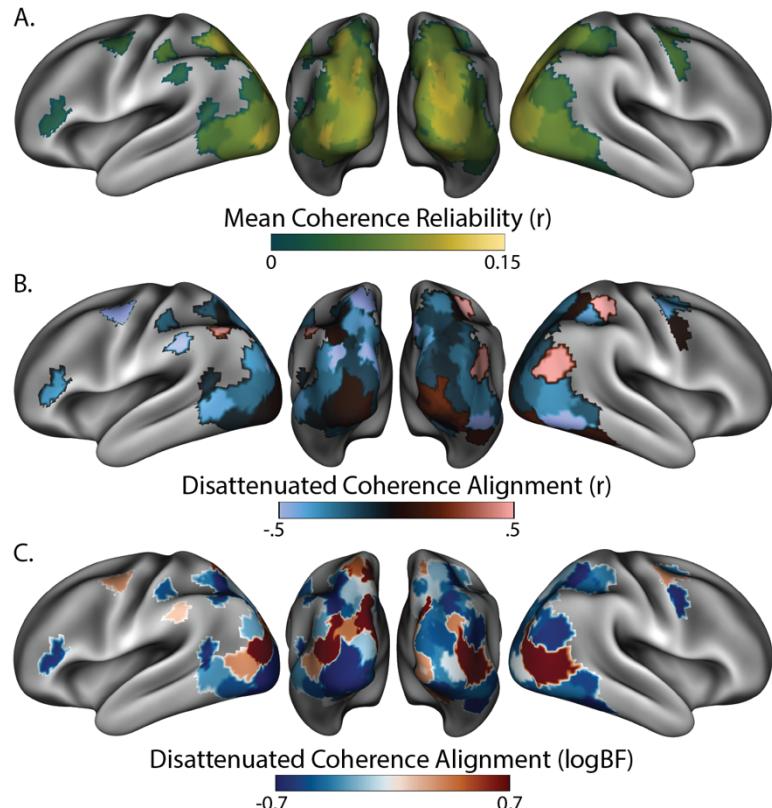


13
14
15
16
17
18
19
20
21

Supplementary Figure 3. *Encoding Geometry Analysis (EGA) validation.* We validated how well we could recover
the similarity between linear Gaussian models (training: $Y = XB + \Sigma$, test: $Y' = X'B' + \Sigma$). Y is the $[1000 \times 250]$
activity timeseries, X is the $[1000 \times 1]$ design matrix, B is the $[1 \times 250]$ encoding profile, and Σ reflects IID
Gaussian noise. In each of our 1000 simulations, we used two different methods to recover the similarity between
the true training encoding profile (B) and the true test encoding profile ($B' = B + \mathcal{N}(0, 1)$), based on noisy activity
timeseries ($Y = XB + \mathcal{N}(0, \sigma_Y)$; $Y' = X'B' + \mathcal{N}(0, \sigma_Y)$). The first method was *pattern reliability* (i.e., our EGA
method), correlating the encoding profile estimated during training ($\hat{B} = X^\dagger Y$, \dagger indicates pseudoinverse) with the
encoding profile estimated during test ($\hat{B}' = X'^\dagger Y'$). The second method was *activity prediction* (i.e., the traditional

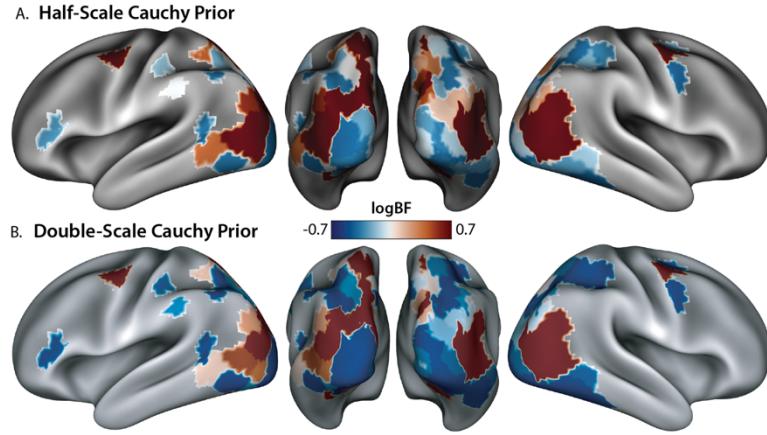
1 encoding approach), correlating the ground-truth test activity (Y') with the predicted test activity ($\hat{Y}' = X'\hat{B}$) after
 2 vectorizing both multivariate timeseries. To simulate the high measurement noise inherent to fMRI, we compared
 3 these methods under different levels of residual SD (σ_Y). **A)** Estimated pattern reliability tracked the true pattern
 4 reliability (i.e., the true correlation between B' and B) across the full range of residual SD, with some attenuation at
 5 high levels of noise **B)** Unlike pattern reliability, activity prediction became much poorer as residual SD increased.
 6 **C)** Correlating the true pattern reliability (correlation between B and B') and estimated encoding strength (i.e.,
 7 pattern reliability or activity prediction), we found pattern reliability was better correlated with the true reliability,
 8 particularly at higher levels of noise. **D)** Both methods had similar performance in the absence of a signal ($B'_{null} =$
 9 $\mathcal{N}(0,1)$).

10
 11

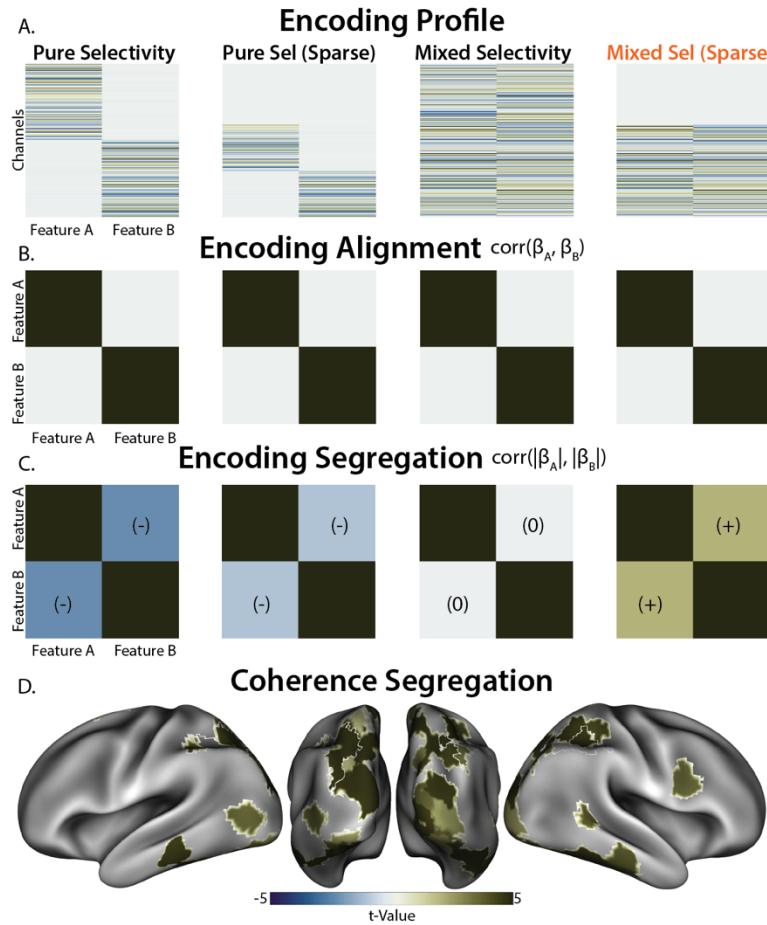


12
 13 **Supplementary Figure 4.** Reliability control analysis. **A)** Geometric mean of target and distractor coherence
 14 reliability ($\sqrt{r_{targ} \times r_{dist}}$), plotted in the reliability-thresholded parcels used in Figure 4. Reliability provides the
 15 theoretical upper bound on correlation strength. Median across participants, excluding participants with non-positive
 16 reliability. **B)** Target-distractor correlations, normalized by target-distractor reliability (i.e., disattenuated
 17 correlations) **C)** Log bayes factors for disattenuated target-distractor correlations. Compare to Figure 4C.
 18

1
2
3
4
5
6
7
8
9
10
11
12

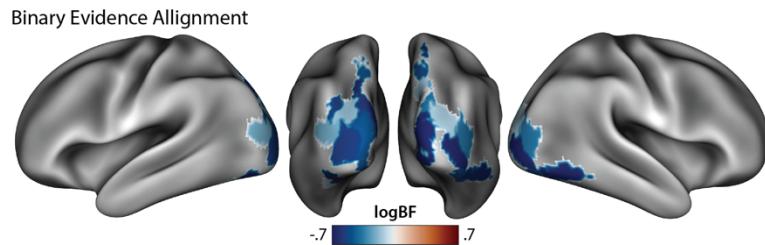


Supplementary Figure 5. Bayes factor prior control analysis. **A)** Log bayes factors for target-distractor coherence alignment using a narrower prior (one-half the default Cauchy scale = 0.35). Minimum logBF is -0.46 at $t_{(28)} = 0$. **B)** Same log bayes factor using a wider prior (double the default Cauchy scale = 1.41). Minimum logBF = -0.99 at $t_{(28)} = 0$. Across different prior parameterizations, note the similarity to Figure 4C.

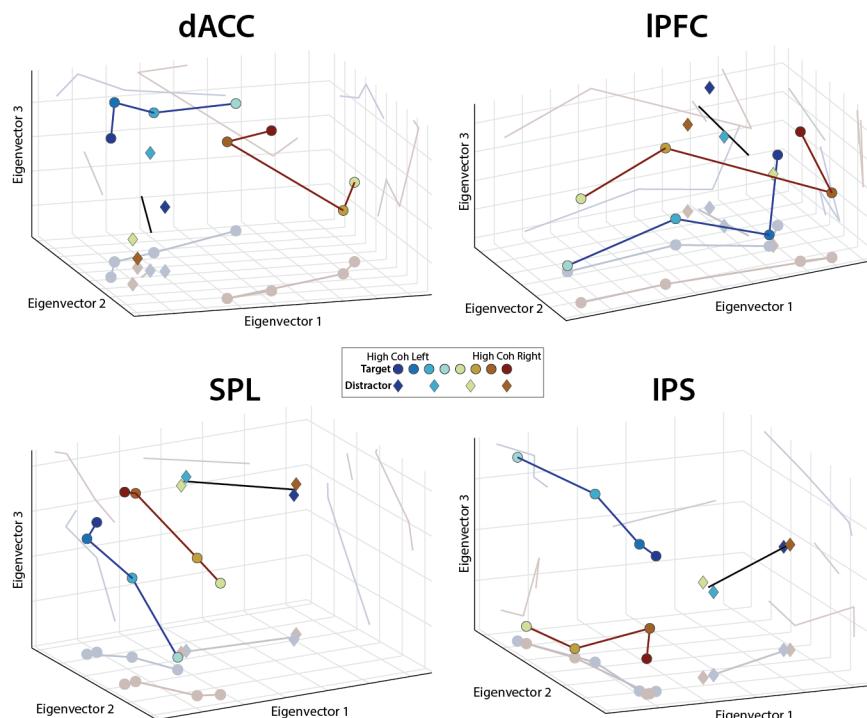


Supplementary Figure 6. Segregation Analysis. **A)** We used pattern component modelling¹³¹ to simulate different candidate encoding profiles. ‘Pure Selectivity’ reflects the segregated encoding hypothesis, with different voxels (rows) encoding different features (columns). ‘Mixed Selectivity’ reflects the orthogonal subspace hypothesis, with the same voxels encoding both features. ‘Sparse’ models include non-selective voxels. **B)** By design, all of these encoding profiles had the same orthogonal encoding alignment (uncorrelated encoding weights), highlighting that

1 this measure is unable to adjudicate between candidate encoding profiles. **C)** These models can be differentiated by
 2 correlating their absolute encoding weights, testing whether the sensitivity of a voxel to one feature is related to its
 3 sensitivity to the other feature, ignoring the direction of encoding. Pure selective encoding predicts a negative
 4 relationship, mixed selective encoding predicts no relationship, and sparse mixed selective encoding predicts a
 5 positive relationship. Similarity matrices averaged over 10,000 simulations. **D)** Correlating the absolute encoding
 6 weights, we found that the IPS profile was consistent with sparse mixed selective encoding.
 7

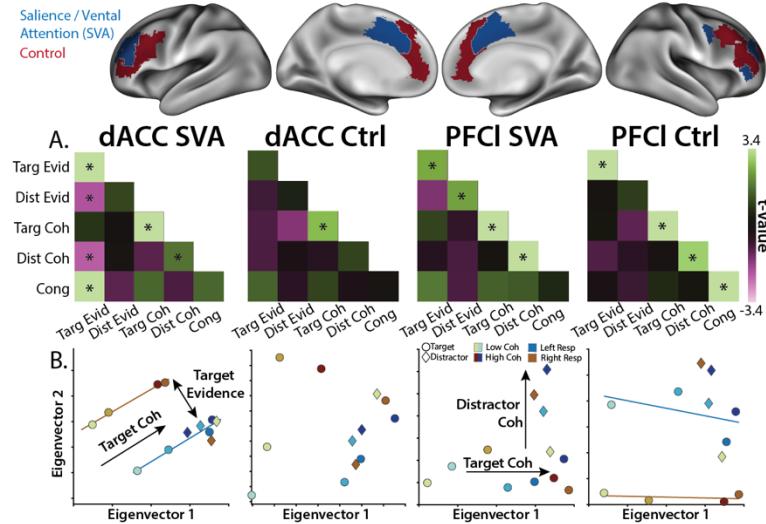


8 **Supplementary Figure 7.** *Binary evidence encoding control analysis.* Target-distractor response encoding
 9 alignment using binary evidence rather than coherence-modulated evidence. Note the similarity to Figure 4D.
 10

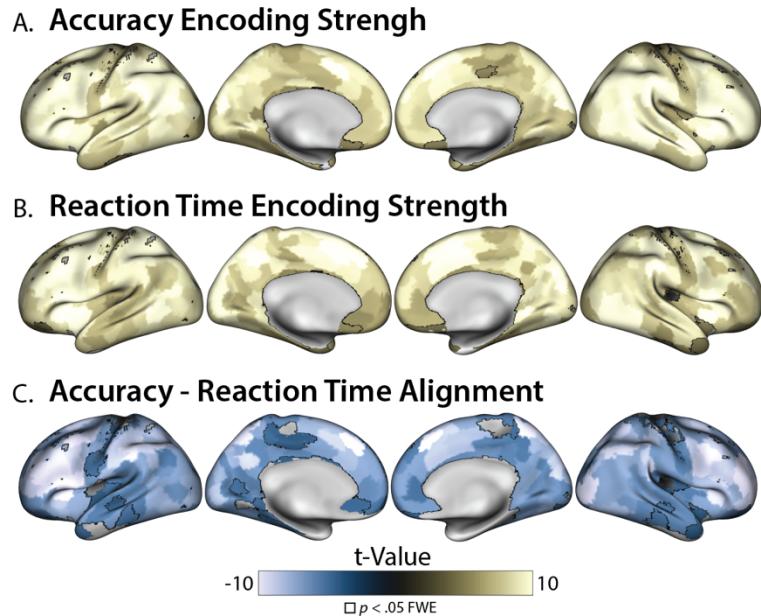


11 **Supplementary Figure 8.** *Multidimensional scaling across higher dimensions.* The first three principal components
 12 of region-averaged condition similarity. Dark lines highlight the encoding geometry (connecting target coherence
 13 circles and showing the average direction for distractor coherence diamonds). Gray lines reflect the projection of
 14 these trends on different planes of the representational space. See legend and Figure 4B for figure details. Note that
 15 in IPS, whereas targets and distractors are encoded orthogonally in the first two dimensions (floor), there appears to
 16 be some alignment in higher dimensions (right wall). In SPL, features appear to be aligned in all dimensions.
 17

18
 19
 20
 21
 22
 23
 24



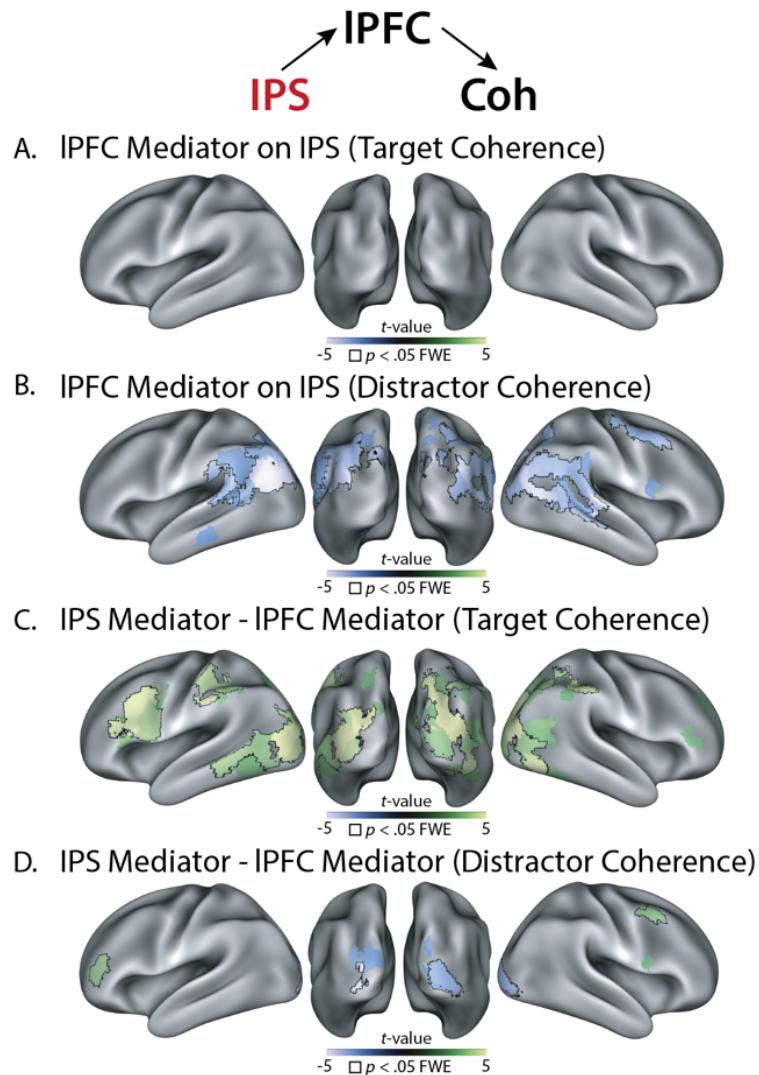
Supplementary Figure 9. Feature encoding in frontal networks. **A)** Similarity matrices for ‘Salience / Ventral Attention (SVA)’ and ‘Control’ networks in dACC and IPFC, correlating feature evidence (‘Evid’), feature coherence (‘Coh’), and feature congruence (‘Cong’). Encoding strength on diagonal (right-tailed p -value), encoding alignment on off-diagonal (two-tailed p -value). **B)** Classical MDS embedding of target (circle) and distractor (diamond) representations at different levels of evidence. Colors denote responses, hues denote coherence. GLMs: A: Feature MV, B: Evidence Levels, see Table 2.



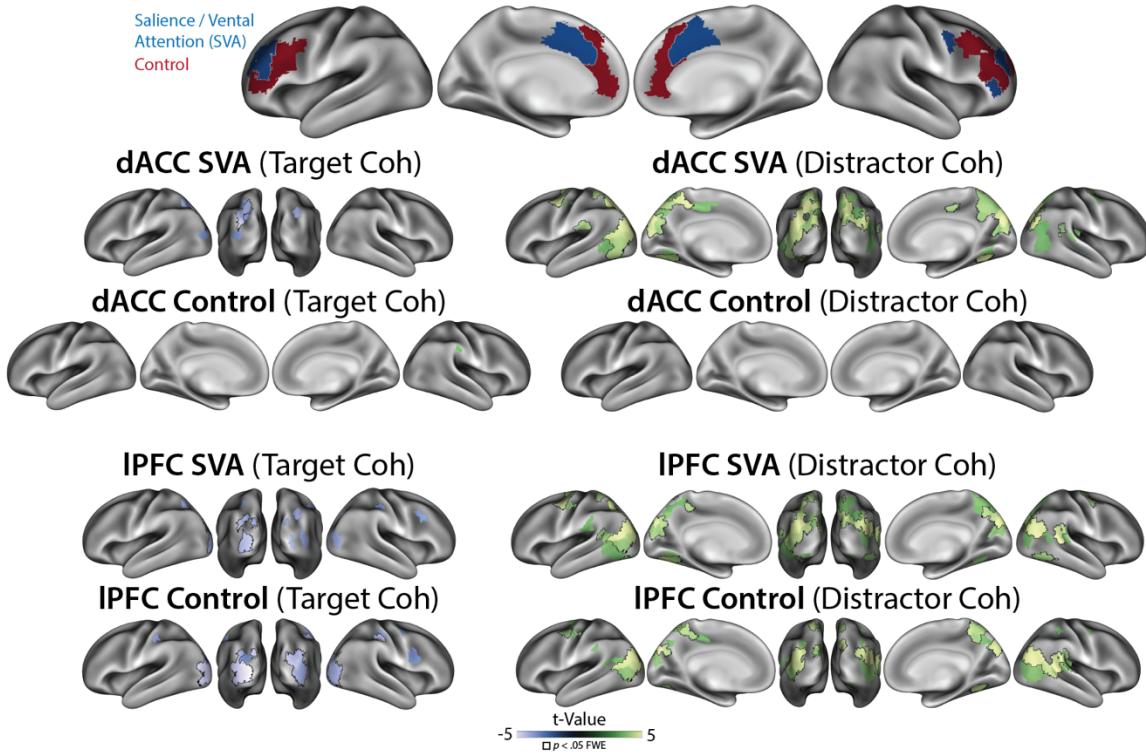
Supplementary Figure 10. Performance encoding. Encoding Strength (across-run reliability) for **A)** Accuracy and **B)** Reaction Time (**C**). **C**) Alignment between Accuracy and Reaction Time encoding. Outlined parcels are significant at $p < .05$ FWE (max-statistic randomization test). Parcels in **C** are thresholded based on the reliability in **A** and **B** (both $p < .001$). GLM: Performance.



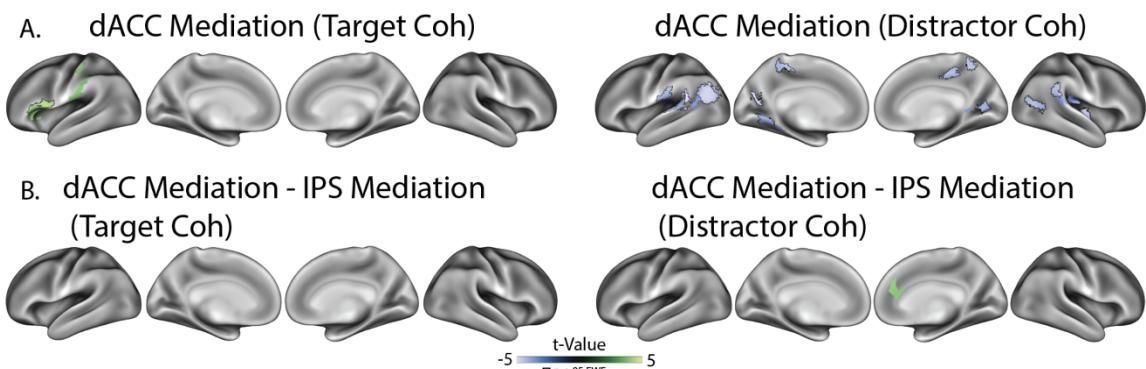
Supplementary Figure 11. Connectivity Alignment Schematic. We estimated connectivity encoding by getting the aggregated residual timeseries from our seed regions (eigenvariate; left), including these timeseries in our whole-brain GLM (middle), and then testing the alignment between connectivity encoding patterns and task encoding patterns (right).



Supplementary Figure 12. IPFC mediation. IPS → IPFC → Coherence mediation for target coherence (A) and distractor coherence (B; compare to Figure 7B). Contrast between IPS-mediation and IPFC-mediation for target coherence (C) and distractor coherence (D).

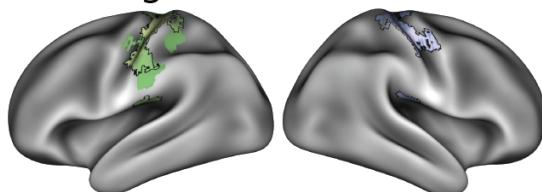


Supplementary Figure 13. Coherence alignment with frontal networks. Activity in ‘Salience / Ventral Attention (SVA)’ and ‘Control’ networks within dACC and IPFC (rows), aligned with target and distractor coherence (columns). Note the similarity between dACC SVA parcels and IPFC parcels.

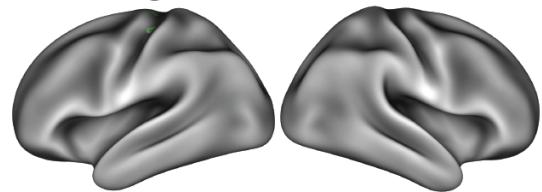


Supplementary Figure 14. IPS mediation of dACC connectivity. A) IPS mediation of dACC connectivity (difference in dACC-coherence alignment with and without including IPS predictors). B) Difference between ‘IPS mediation of dACC’ and ‘dACC mediation of IPS’. The lack of activation suggests that this relationship is bidirectional, or originates from a common cause. dACC seed is from the ‘Salience / Ventral Attention’ network (see Supplementary Figure 11).

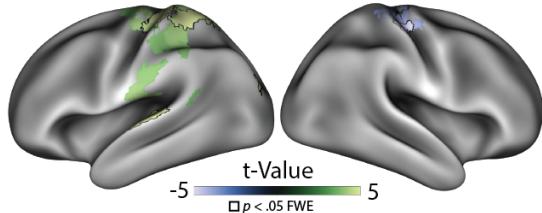
A. Target Evidence ~ SPL



B. Target Evidence ~ IPS

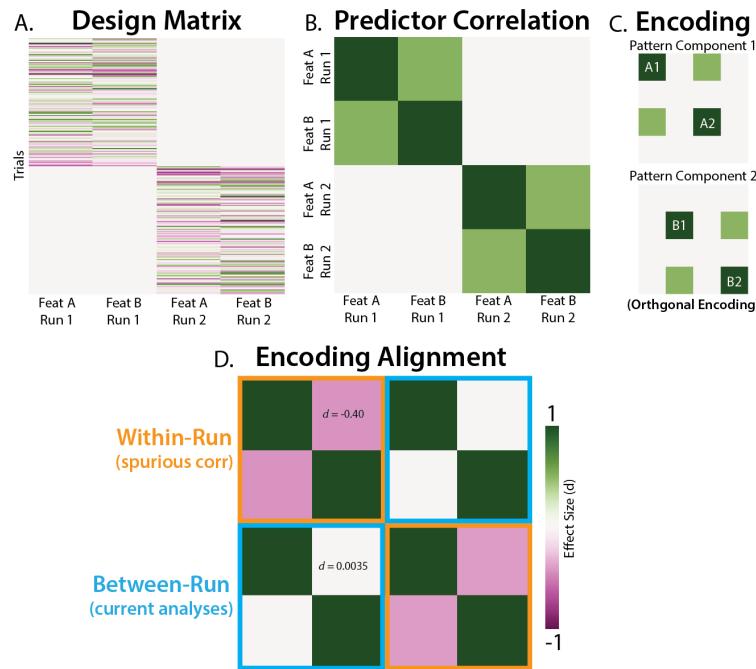


C. (Target Evidence ~ SPL) -
(Target Evidence ~ IPS)



t -Value
 $\square p < .05$ FWE

Supplementary Figure 15. SPL alignment with evidence encoding. **A)** Alignment between SPL activity and target evidence encoding. **B)** Alignment between IPS activity and target evidence encoding. **B)** Differences between SPL-evidence alignment and IPS-evidence alignment, showing stronger SPL connectivity. Note that target evidence encoding is signed according to the right-hand response (contralateral motor cortex should have a positive response).



1 **Supplementary Figure 16.** *Cross-validation avoids feature correlations biasing alignment.* We used pattern
 2 component modeling¹³¹ to simulate neural data, testing whether feature correlations could spuriously create
 3 encoding alignment. **A)** Our design matrix had two simulate runs of two feature timeseries. **B)** Our features were
 4 correlated by design (i.e., the columns of the design matrix were correlated). **C)** Despite correlation in the design
 5 matrix, these features were independently encoding in our simulated neural population (i.e., in two distinct pattern
 6 components, which were each reliable across runs). **D)** Correlating our estimated encoding profiles, we found that
 7 within-run alignment (orange) had a spurious negative correlation (the opposite direction of the feature correlations).
 8 Critically, our analyses used between-run alignment (cyan), which avoids this biasing effect of feature correlations.
 9 Intuitively, since features are not correlated across runs (i.e., they come from different trials), they do not produce
 10 spurious correlations. Effect sizes are computed across 10,000 simulations.
 11

Correlation	Covariates	dACC	IPFC	SPL	IPS
Target-Accuracy, Target-RT	Target, Accuracy, RT	$r_{(27)} = -0.32$ $p = .11$	$r_{(27)} = -0.36$ $p = .067$	$r_{(27)} = -0.11$ $p = .56$	$r_{(27)} = -0.47$ $p = .017$
Distractor-Accuracy, Distractor-RT	Distractor, Accuracy, RT	$r_{(27)} = -0.71$ $p = 0.50 \times 10^{-4}$	$r_{(27)} = -0.43$ $p = .027$	$r_{(27)} = -0.48$ $p = .012$	$r_{(27)} = -0.59$ $p = .0014$

12 **Supplementary Table 1.** *Partial correlations between coherence and performance.* Correlations
 13 between individual differences in coherence-performance alignment, controlling for coherence
 14 and performance encoding reliability. Since reliability determines alignment⁶⁶, similarity in
 15 alignment may be confounded with similarity in reliability. Overall, these results are
 16 qualitatively similar to the zero-order correlation (see Figure 6), albeit with weaker correlations
 17 for target coherence. These correlations are particularly robust in IPS.
 18

1 References

- 2 1. Musslick, S., Shenhav, A., Botvinick, M. & Cohen, J. A Computational Model of Control
3 Allocation based on the Expected Value of Control. in *2nd Multidisciplinary Conference on*
4 *Reinforcement Learning and Decision Making* (2015).
- 5 2. Badre, D., Bhandari, A., Keglovits, H. & Kikumoto, A. The dimensionality of neural
6 representations for control. *Curr Opin Behav Sci* **38**, 20–28 (2021).
- 7 3. Danielmeier, C., Eichele, T., Forstmann, B. U., Tittgemeyer, M. & Ullsperger, M. Posterior
8 medial frontal cortex activity predicts post-error adaptations in task-related visual and motor
9 areas. *J. Neurosci.* **31**, 1780–1789 (2011).
- 10 4. Egner, T. Multiple conflict-driven control mechanisms in the human brain. *Trends Cogn.
Sci.* **12**, 374–380 (2008).
- 11 5. Ritz, H., Leng, X. & Shenhav, A. Cognitive Control as a Multivariate Optimization
13 Problem. *J. Cogn. Neurosci.* **34**, 569–591 (2022).
- 14 6. Friedman, N. P. & Miyake, A. Unity and diversity of executive functions: Individual
15 differences as a window on cognitive structure. *Cortex* **86**, 186–204 (2017).
- 16 7. Danielmeier, C. & Ullsperger, M. Post-error adjustments. *Front. Psychol.* **2**, 233 (2011).
- 17 8. Fischer, A. G., Nigbur, R., Klein, T. A., Danielmeier, C. & Ullsperger, M. Cortical beta
18 power reflects decision dynamics and uncovers multiple facets of post-error adaptation. *Nat.
Commun.* **9**, 5038 (2018).
- 20 9. Leng, X., Yee, D., Ritz, H. & Shenhav, A. Dissociable influences of reward and punishment
21 on adaptive cognitive control. *PLoS Comput. Biol.* **17**, e1009737 (2021).
- 22 10. Ritz, H. & Shenhav, A. Humans reconfigure target and distractor processing to address
23 distinct task demands. *bioRxiv* 2021.09.08.459546 (2021) doi:10.1101/2021.09.08.459546.

- 1 11. Kalman, R. E. On the general theory of control systems. *IFAC Proceedings Volumes* **1**,
2 491–502 (1960).
- 3 12. Shenhav, A., Botvinick, M. M. & Cohen, J. D. The expected value of control: an integrative
4 theory of anterior cingulate cortex function. *Neuron* **79**, 217–240 (2013).
- 5 13. MacDonald, A. W., 3rd, Cohen, J. D., Stenger, V. A. & Carter, C. S. Dissociating the role
6 of the dorsolateral prefrontal and anterior cingulate cortex in cognitive control. *Science* **288**,
7 1835–1838 (2000).
- 8 14. Smith, E. H. *et al.* Widespread temporal coding of cognitive control in the human prefrontal
9 cortex. *Nat. Neurosci.* **66**, 83 (2019).
- 10 15. Menon, V. & D’Esposito, M. The role of PFC networks in cognitive control and executive
11 function. *Neuropsychopharmacology* 1–14 (2021) doi:10.1038/s41386-021-01152-w.
- 12 16. Kerns, J. G. *et al.* Anterior cingulate conflict monitoring and adjustments in control. *Science*
13 **303**, 1023–1026 (2004).
- 14 17. Gottlieb, J., Cohanpour, M., Li, Y., Singletary, N. & Zabeh, E. Curiosity, information
15 demand and attentional priority. *Current Opinion in Behavioral Sciences* **35**, 83–91 (2020).
- 16 18. Gordon, E. M. *et al.* Precision Functional Mapping of Individual Human Brains. *Neuron* **95**,
17 791-807.e7 (2017).
- 18 19. Gratton, C., Laumann, T. O., Gordon, E. M., Adeyemo, B. & Petersen, S. E. Evidence for
19 Two Independent Factors that Modify Brain Networks to Meet Task Goals. *Cell Rep.* **17**,
20 1276–1288 (2016).
- 21 20. Kragel, P. A. *et al.* Generalizable representations of pain, cognitive control, and negative
22 emotion in medial frontal cortex. *Nat. Neurosci.* **21**, 283–289 (2018).

- 1 21. Fu, Z. *et al.* The geometry of domain-general performance monitoring in the human medial
2 frontal cortex. *Science* **376**, eabm9922 (2022).
- 3 22. Vermeylen, L. *et al.* Shared Neural Representations of Cognitive Conflict and Negative
4 Affect in the Medial Frontal Cortex. *J. Neurosci.* **40**, 8715–8725 (2020).
- 5 23. Brown, J. W. & Braver, T. S. Learned predictions of error likelihood in the anterior
6 cingulate cortex. *Science* **307**, 1118–1121 (2005).
- 7 24. Rushworth, M. F. & Behrens, T. E. Choice, uncertainty and value in prefrontal and
8 cingulate cortex. *Nat. Neurosci.* **11**, 389–397 (2008).
- 9 25. Grinband, J. *et al.* The dorsal medial frontal cortex is sensitive to time on task, not response
10 conflict or error likelihood. *Neuroimage* **57**, 303–311 (2011).
- 11 26. Yarkoni, T., Barch, D. M., Gray, J. R., Conturo, T. E. & Braver, T. S. BOLD correlates of
12 trial-by-trial reaction time variability in gray and white matter: a multi-study fMRI analysis.
13 *PLoS One* **4**, e4257 (2009).
- 14 27. Mumford, J. A. *et al.* The response time paradox in functional magnetic resonance imaging
15 analyses. *bioRxiv* 2023.02.15.528677 (2023) doi:10.1101/2023.02.15.528677.
- 16 28. Pylyshyn, Z. W. & Storm, R. W. Tracking multiple independent targets: evidence for a
17 parallel tracking mechanism. *Spat. Vis.* **3**, 179–197 (1988).
- 18 29. Beldzik, E. & Ullsperger, M. A thin line between conflict and reaction time effects on EEG
19 and fMRI brain signals. *bioRxiv* 2023.02.14.528515 (2023)
20 doi:10.1101/2023.02.14.528515.
- 21 30. Shenhav, A., Straccia, M. A., Musslick, S., Cohen, J. D. & Botvinick, M. M. Dissociable
22 neural mechanisms track evidence accumulation for selection of attention versus action.
23 *Nat. Commun.* **9**, 2485 (2018).

- 1 31. Taren, A. A., Venkatraman, V. & Huettel, S. A. A parallel functional topography between
2 medial and lateral prefrontal cortex: evidence and implications for cognitive control. *J.*
3 *Neurosci.* **31**, 5026–5031 (2011).
- 4 32. Venkatraman, V., Rosati, A. G., Taren, A. A. & Huettel, S. A. Resolving response,
5 decision, and strategic control: evidence for a functional topography in dorsomedial
6 prefrontal cortex. *J. Neurosci.* **29**, 13158–13164 (2009).
- 7 33. Fu, Z. *et al.* Single-Neuron Correlates of Error Monitoring and Post-Error Adjustments in
8 Human Medial Frontal Cortex. *Neuron* **101**, 165-177.e5 (2019).
- 9 34. Zarr, N. & Brown, J. W. Hierarchical error representation in medial prefrontal cortex.
10 *Neuroimage* **124**, 238–247 (2016).
- 11 35. Ebitz, B. R. *et al.* Human dorsal anterior cingulate neurons signal conflict by amplifying
12 task-relevant information. *bioRxiv* 2020.03.14.991745 (2020)
13 doi:10.1101/2020.03.14.991745.
- 14 36. Flesch, T., Juechems, K., Dumbalska, T., Saxe, A. & Summerfield, C. Orthogonal
15 representations for robust context-dependent task performance in brains and neural
16 networks. *Neuron* **0**, (2022).
- 17 37. Minxha, J., Adolphs, R., Fusi, S., Mamelak, A. N. & Rutishauser, U. Flexible recruitment of
18 memory-based choice representations by the human medial frontal cortex. *Science* **368**,
19 (2020).
- 20 38. Ebitz, R. B. & Hayden, B. Y. The population doctrine in cognitive neuroscience. *Neuron*
21 (2021) doi:10.1016/j.neuron.2021.07.011.
- 22 39. Cunningham, J. P. & Yu, B. M. Dimensionality reduction for large-scale neural recordings.
23 *Nat. Neurosci.* **17**, 1500–1509 (2014).

- 1 40. Rigotti, M. *et al.* The importance of mixed selectivity in complex cognitive tasks. *Nature*
2 **497**, 585–590 (2013).
- 3 41. Mante, V., Sussillo, D., Shenoy, K. V. & Newsome, W. T. Context-dependent computation
4 by recurrent dynamics in prefrontal cortex. *Nature* **503**, 78–84 (2013).
- 5 42. Kayser, A. S., Erickson, D. T., Buchsbaum, B. R. & D’Esposito, M. Neural representations
6 of relevant and irrelevant features in perceptual decision making. *J. Neurosci.* **30**, 15778–
7 15789 (2010).
- 8 43. Bernardi, S. *et al.* The Geometry of Abstraction in the Hippocampus and Prefrontal Cortex.
9 *Cell* **0**, (2020).
- 10 44. Panichello, M. F. & Buschman, T. J. Shared mechanisms underlie the control of working
11 memory and attention. *Nature* 1–5 (2021) doi:10.1038/s41586-021-03390-w.
- 12 45. Takagi, Y., Hunt, L. T., Woolrich, M. W., Behrens, T. E. & Klein-Flügge, M. C. Adapting
13 non-invasive human recordings along multiple task-axes shows unfolding of spontaneous
14 and over-trained choice. *Elife* **10**, (2021).
- 15 46. Cohen, J. D., Dunbar, K. & McClelland, J. L. On the control of automatic processes: a
16 parallel distributed processing account of the Stroop effect. *Psychol. Rev.* **97**, 332–361
17 (1990).
- 18 47. Woolgar, A., Hampshire, A., Thompson, R. & Duncan, J. Adaptive coding of task-relevant
19 information in human frontoparietal cortex. *J. Neurosci.* **31**, 14592–14599 (2011).
- 20 48. Nee, D. E., Wager, T. D. & Jonides, J. Interference resolution: insights from a meta-analysis
21 of neuroimaging tasks. *Cogn. Affect. Behav. Neurosci.* **7**, 1–17 (2007).

- 1 49. Shenhav, A., Straccia, M. A., Botvinick, M. M. & Cohen, J. D. Dorsal anterior cingulate
2 and ventromedial prefrontal cortex have inverse roles in both foraging and economic
3 choice. *Cogn. Affect. Behav. Neurosci.* **16**, 1127–1139 (2016) doi:10.3758/s13415-016-0458-8.
- 4 50. Fleming, S. M., van der Putten, E. J. & Daw, N. D. Neural mediators of changes of mind
5 about perceptual decisions. *Nat. Neurosci.* **21**, 617–624 (2018).
- 6 51. Shenhav, A. & Karmarkar, U. R. Dissociable components of the reward circuit are involved
7 in appraisal versus choice. *Sci. Rep.* **9**, 1958 (2019).
- 8 52. Clairis, N. & Pessiglione, M. Value, confidence, deliberation: a functional partition of the
9 medial prefrontal cortex demonstrated across rating and choice tasks. *J. Neurosci.* **42**, 5580–
10 5592 (2022).
- 11 53. Yeo, B. T. T. *et al.* The organization of the human cerebral cortex estimated by intrinsic
12 functional connectivity. *J. Neurophysiol.* **106**, 1125–1165 (2011).
- 13 54. Kong, R. *et al.* Individual-Specific Areal-Level Parcellations Improve Functional
14 Connectivity Prediction of Behavior. *Cereb. Cortex* **31**, 4477–4500 (2021).
- 15 55. Schaefer, A. *et al.* Local-Global Parcellation of the Human Cerebral Cortex from Intrinsic
16 Functional Connectivity MRI. *Cereb. Cortex* **28**, 3095–3114 (2018).
- 17 56. Kong, R. *et al.* Spatial Topography of Individual-Specific Cortical Networks Predicts
18 Human Cognition, Personality, and Emotion. *Cereb. Cortex* **29**, 2533–2551 (2019).
- 19 57. Bisley, J. W. & Mirpour, K. The neural instantiation of a priority map. *Curr. Opin. Psychol.*
20 **29**, 108–112 (2019).
- 21 58. Yantis, S. & Serences, J. T. Cortical mechanisms of space-based and object-based
22 attentional control. *Curr. Opin. Neurobiol.* **13**, 187–193 (2003).

- 1 59. Kriegeskorte, N. & Diedrichsen, J. Peeling the Onion of Brain Representations. *Annu. Rev.*
2 *Neurosci.* **42**, 407–432 (2019).
- 3 60. Cohen, M. R. & Maunsell, J. H. R. A neuronal population measure of attention predicts
4 behavioral performance on individual trials. *J. Neurosci.* **30**, 15241–15253 (2010).
- 5 61. Libby, A. & Buschman, T. J. Rotational dynamics reduce interference between sensory and
6 memory representations. *Nat. Neurosci.* 1–12 (2021) doi:10.1038/s41593-021-00821-9.
- 7 62. Kimmel, D. L., Elsayed, G. F., Cunningham, J. P. & Newsome, W. T. Value and choice as
8 separable and stable representations in orbitofrontal cortex. *Nat. Commun.* **11**, 3466 (2020).
- 9 63. Walther, A. *et al.* Reliability of dissimilarity measures for multi-voxel pattern analysis.
10 *Neuroimage* **137**, 188–200 (2016).
- 11 64. Diedrichsen, J. & Kriegeskorte, N. Representational models: A common framework for
12 understanding encoding, pattern-component, and representational-similarity analysis. *PLoS*
13 *Comput. Biol.* **13**, e1005508 (2017).
- 14 65. Nili, H. *et al.* A toolbox for representational similarity analysis. *PLoS Comput. Biol.* **10**,
15 e1003553 (2014).
- 16 66. Spearman, C. The Proof and Measurement of Association between Two Things. *Am. J.*
17 *Psychol.* **100**, 441–471 (1987).
- 18 67. Thornton, M. A. & Mitchell, J. P. Consistent Neural Activity Patterns Represent Personally
19 Familiar People. *J. Cogn. Neurosci.* **29**, 1583–1594 (2017).
- 20 68. Hunt, L. T. *et al.* Mechanisms underlying cortical activity during value-guided choice. *Nat.*
21 *Neurosci.* **15**, 470–6, S1-3 (2012).

- 1 69. Kayser, A. S., Buchsbaum, B. R., Erickson, D. T. & D'Esposito, M. The functional
2 anatomy of a perceptual decision in the human brain. *J. Neurophysiol.* **103**, 1179–1194
3 (2010).
- 4 70. Kriegeskorte, N., Goebel, R. & Bandettini, P. Information-based functional brain mapping.
5 *Proc. Natl. Acad. Sci. U. S. A.* **103**, 3863–3868 (2006).
- 6 71. Woolgar, A., Williams, M. A. & Rich, A. N. Attention enhances multi-voxel representation
7 of novel objects in frontal, parietal and visual cortices. *Neuroimage* **109**, 429–437 (2015).
- 8 72. Woolgar, A., Afshar, S., Williams, M. A. & Rich, A. N. Flexible Coding of Task Rules in
9 Frontoparietal Cortex: An Adaptive System for Flexible Cognitive Control. *J. Cogn.*
10 *Neurosci.* **27**, 1895–1911 (2015).
- 11 73. Jackson, J., Rich, A. N., Williams, M. A. & Woolgar, A. Feature-selective Attention in
12 Frontoparietal Cortex: Multivoxel Codes Adjust to Prioritize Task-relevant Information. *J.*
13 *Cogn. Neurosci.* **29**, 310–321 (2017).
- 14 74. Woolgar, A., Thompson, R., Bor, D. & Duncan, J. Multi-voxel coding of stimuli, rules, and
15 responses in human frontoparietal cortex. *Neuroimage* **56**, 744–752 (2011).
- 16 75. Aoi, M. C., Mante, V. & Pillow, J. W. Prefrontal cortex exhibits multidimensional dynamic
17 encoding during decision-making. *Nat. Neurosci.* (2020) doi:10.1038/s41593-020-0696-5.
- 18 76. Pagan, M. *et al.* A new theoretical framework jointly explains behavioral and neural
19 variability across subjects performing flexible decision-making. *bioRxiv* 2022.11.28.518207
20 (2022) doi:10.1101/2022.11.28.518207.
- 21 77. Miller, E. K. & Cohen, J. D. An integrative theory of prefrontal cortex function. *Annu. Rev.*
22 *Neurosci.* **24**, 167–202 (2001).

- 1 78. Stringer, C. *et al.* Spontaneous behaviors drive multidimensional, brainwide activity.
- 2 *Science* **364**, 255 (2019).
- 3 79. Goldman-Rakic, P. S. Topography of cognition: parallel distributed networks in primate
- 4 association cortex. *Annu. Rev. Neurosci.* **11**, 137–156 (1988).
- 5 80. Corbetta, M. & Shulman, G. L. Control of goal-directed and stimulus-driven attention in the
- 6 brain. *Nat. Rev. Neurosci.* **3**, 201–215 (2002).
- 7 81. Suzuki, M. & Gottlieb, J. Distinct neural mechanisms of distractor suppression in the frontal
- 8 and parietal lobe. *Nat. Neurosci.* **16**, 98–104 (2013).
- 9 82. Kastner, S. & Ungerleider, L. G. Mechanisms of visual attention in the human cortex. *Annu.*
- 10 *Rev. Neurosci.* **23**, 315–341 (2000).
- 11 83. Kay, K. N. & Yeatman, J. D. Bottom-up and top-down computations in word- and face-
- 12 selective cortex. *Elife* **6**, (2017).
- 13 84. Saalmann, Y. B., Pigarev, I. N. & Vidyasagar, T. R. Neural mechanisms of visual attention:
- 14 how top-down feedback highlights relevant locations. *Science* **316**, 1612–1615 (2007).
- 15 85. Reid, A. T. *et al.* Advancing functional connectivity research from association to causation.
- 16 *Nat. Neurosci.* **22**, 1751–1760 (2019).
- 17 86. MacKinnon, D. P., Fairchild, A. J. & Fritz, M. S. Mediation analysis. *Annu. Rev. Psychol.*
- 18 **58**, 593–614 (2007).
- 19 87. Friedman, N. P. & Robbins, T. W. The role of prefrontal cortex in cognitive control and
- 20 executive function. *Neuropsychopharmacology* 1–18 (2021) doi:10.1038/s41386-021-
- 21 01132-0.
- 22 88. Holroyd, C. B. & McClure, S. M. Hierarchical control over effortful behavior by rodent
- 23 medial frontal cortex: A computational model. *Psychol. Rev.* **122**, 54–83 (2015).

- 1 89. Holroyd, C. B. & Yeung, N. An Integrative Theory of Anterior Cingulate Cortex Function:
2 Option Selection in Hierarchical Reinforcement Learning. *Neural Basis of Motivational and*
3 *Cognitive Control* 332–349 Preprint at
4 <https://doi.org/10.7551/mitpress/9780262016438.003.0018> (2011).
- 5 90. Vassena, E., Deraeve, J. & Alexander, W. H. Predicting Motivation: Computational Models
6 of PFC Can Explain Neural Coding of Motivation and Effort-based Decision-making in
7 Health and Disease. *J. Cogn. Neurosci.* **29**, 1633–1645 (2017).
- 8 91. Koechlin, E. & Summerfield, C. An information theoretical approach to prefrontal
9 executive function. *Trends Cogn. Sci.* **11**, 229–235 (2007).
- 10 92. Badre, D. & D'Esposito, M. Is the rostro-caudal axis of the frontal lobe hierarchical? *Nat.*
11 *Rev. Neurosci.* **10**, 659–669 (2009).
- 12 93. Badre, D. & Nee, D. E. Frontal Cortex and the Hierarchical Control of Behavior. *Trends*
13 *Cogn. Sci.* **22**, 170–188 (2018).
- 14 94. Serences, J. T. & Yantis, S. Spatially selective representations of voluntary and stimulus-
15 driven attentional priority in human occipital, parietal, and frontal cortex. *Cereb. Cortex* **17**,
16 284–293 (2007).
- 17 95. Yantis, S. *et al.* Transient neural activity in human parietal cortex during spatial attention
18 shifts. *Nat. Neurosci.* **5**, 995–1002 (2002).
- 19 96. Greenberg, A. S., Esterman, M., Wilson, D., Serences, J. T. & Yantis, S. Control of spatial
20 and feature-based attention in frontoparietal cortex. *J. Neurosci.* **30**, 14330–14339 (2010).
- 21 97. Esterman, M., Chiu, Y.-C., Tamber-Rosenau, B. J. & Yantis, S. Decoding cognitive control
22 in human parietal cortex. *Proc. Natl. Acad. Sci. U. S. A.* **106**, 17974–17979 (2009).

- 1 98. Serences, J. T., Schwarzbach, J., Courtney, S. M., Golay, X. & Yantis, S. Control of object-
2 based attention in human cortex. *Cereb. Cortex* **14**, 1346–1357 (2004).
- 3 99. Molenberghs, P., Mesulam, M. M., Peeters, R. & Vandenberghe, R. R. C. Remapping
4 attentional priorities: differential contribution of superior parietal lobule and intraparietal
5 sulcus. *Cereb. Cortex* **17**, 2703–2712 (2007).
- 6 100. Adam, K. C. S. & Serences, J. T. History modulates early sensory processing of salient
7 distractors. *J. Neurosci.* (2021) doi:10.1523/JNEUROSCI.3099-20.2021.
- 8 101. Soutschek, A., Stelzel, C., Paschke, L., Walter, H. & Schubert, T. Dissociable effects of
9 motivation and expectancy on conflict processing: an fMRI study. *J. Cogn. Neurosci.* **27**,
10 409–423 (2015).
- 11 102. Bisley, J. W. & Goldberg, M. E. Attention, intention, and priority in the parietal lobe. *Annu.
12 Rev. Neurosci.* **33**, 1–21 (2010).
- 13 103. Rust, N. C. & Cohen, M. R. Priority coding in the visual system. *Nat. Rev. Neurosci.* 1–13
14 (2022) doi:10.1038/s41583-022-00582-9.
- 15 104. Etzel, J. A., Cole, M. W., Zacks, J. M., Kay, K. N. & Braver, T. S. Reward Motivation
16 Enhances Task Coding in Frontoparietal Cortex. *Cereb. Cortex* **26**, 1647–1659 (2016).
- 17 105. Hall-McMaster, S., Muhle-Karbe, P. S., Myers, N. E. & Stokes, M. G. Reward Boosts
18 Neural Coding of Task Rules to Optimize Cognitive Flexibility. *J. Neurosci.* **39**, 8549–8561
19 (2019).
- 20 106. Desimone, R. & Duncan, J. Neural mechanisms of selective visual attention. *Annu. Rev.
21 Neurosci.* **18**, 193–222 (1995).
- 22 107. Lauritzen, T. Z., D'Esposito, M., Heeger, D. J. & Silver, M. A. Top-down flow of visual
23 spatial attention signals from parietal to occipital cortex. *J. Vis.* **9**, 18–18 (2009).

- 1 108. Semedo, J. D., Zandvakili, A., Machens, C. K., Yu, B. M. & Kohn, A. Cortical Areas
2 Interact through a Communication Subspace. *Neuron* **102**, 249-259.e4 (2019).
- 3 109. Srinath, R., Ruff, D. A. & Cohen, M. R. Attention improves information flow between
4 neuronal populations without changing the communication subspace. *bioRxiv*
5 2021.03.31.437940 (2021) doi:10.1101/2021.03.31.437940.
- 6 110. Petrides, M. & Pandya, D. N. Efferent association pathways originating in the caudal
7 prefrontal cortex in the macaque monkey. *J. Comp. Neurol.* **498**, 227–251 (2006).
- 8 111. Jackson, J. B., Feredoes, E., Rich, A. N., Lindner, M. & Woolgar, A. Concurrent
9 neuroimaging and neurostimulation reveals a causal role for dlPFC in coding of task-
10 relevant information. *Commun Biol* **4**, 588 (2021).
- 11 112. Vul, E., Alvarez, G., Tenenbaum, J. & Black, M. Explaining human multiple object
12 tracking as resource-constrained approximate inference in a dynamic probabilistic model. in
13 *Advances in Neural Information Processing Systems* (eds. Bengio, Y., Schuurmans, D.,
14 Lafferty, J., Williams, C. & Culotta, A.) vol. 22 (Curran Associates, Inc., 2009).
- 15 113. Ritz, H., Wild, C. J. & Johnsrude, I. S. Parametric Cognitive Load Reveals Hidden Costs in
16 the Neural Processing of Perfectly Intelligible Degraded Speech. *J. Neurosci.* **42**, 4619–
17 4628 (2022).
- 18 114. Culham, J. C., Cavanagh, P. & Kanwisher, N. G. Attention response functions:
19 characterizing brain areas using fMRI activation during parametric variations of attentional
20 load. *Neuron* **32**, 737–745 (2001).
- 21 115. Culham, J. C. *et al.* Cortical fMRI activation produced by attentive tracking of moving
22 targets. *J. Neurophysiol.* **80**, 2657–2670 (1998).

- 1 116. Howe, P. D., Horowitz, T. S., Morocz, I. A., Wolfe, J. & Livingstone, M. S. Using fMRI to
2 distinguish components of the multiple object tracking task. *J. Vis.* **9**, 10.1-11 (2009).
- 3 117. Jovicich, J. *et al.* Brain areas specific for attentional load in a motion-tracking task. *J. Cogn.*
4 *Neurosci.* **13**, 1048–1058 (2001).
- 5 118. Peck, C. J., Jangraw, D. C., Suzuki, M., Efem, R. & Gottlieb, J. Reward modulates attention
6 independently of action value in posterior parietal cortex. *J. Neurosci.* **29**, 11182–11191
7 (2009).
- 8 119. Wisniewski, D., Reverberi, C., Momennejad, I., Kahnt, T. & Haynes, J.-D. The Role of the
9 Parietal Cortex in the Representation of Task–Reward Associations. *J. Neurosci.* **35**,
10 12355–12365 (2015).
- 11 120. Parro, C., Dixon, M. L. & Christoff, K. The neural basis of motivational influences on
12 cognitive control. *Hum. Brain Mapp.* **39**, 5097–5111 (2018).
- 13 121. Weichert, E. R., Turner, B. M. & Sederberg, P. B. A model of dynamic, within-trial conflict
14 resolution for decision making. *Psychol. Rev.* (2020) doi:10.1037/rev0000191.
- 15 122. Esteban, O. *et al.* fMRIprep: a robust preprocessing pipeline for functional MRI. *Nat.*
16 *Methods* **16**, 111–116 (2019).
- 17 123. Gorgolewski, K. *et al.* Nipype: a flexible, lightweight and extensible neuroimaging data
18 processing framework in python. *Front. Neuroinform.* **5**, 13 (2011).
- 19 124. Diedrichsen, J. & Shadmehr, R. Detecting and adjusting for artifacts in fMRI time series
20 data. *Neuroimage* **27**, 624–634 (2005).
- 21 125. Jones, M. S., Zhu, Z., Bajracharya, A., Luor, A. & Peelle, J. E. A multi-dataset evaluation
22 of frame censoring for task-based fMRI. *bioRxiv* 2021.10.12.464075 (2021)
23 doi:10.1101/2021.10.12.464075.

- 1 126. Smith, S. M. & Nichols, T. E. Threshold-free cluster enhancement: addressing problems of
2 smoothing, threshold dependence and localisation in cluster inference. *Neuroimage* **44**, 83–
3 98 (2009).
- 4 127. Winkler, A. M., Ridgway, G. R., Webster, M. A., Smith, S. M. & Nichols, T. E.
5 Permutation inference for the general linear model. *Neuroimage* **92**, 381–397 (2014).
- 6 128. Vos de Wael, R. *et al.* BrainSpace: a toolbox for the analysis of macroscale gradients in
7 neuroimaging and connectomics datasets. *Commun Biol* **3**, 103 (2020).
- 8 129. Gale, D. J., Vos de Wael, R., Benkarim, O. & Bernhardt, B. *Surfplot: Publication-ready
9 brain surface figures*. (2021). doi:10.5281/zenodo.5567926.
- 10 130. Yarkoni, T., Poldrack, R. A., Nichols, T. E., Van Essen, D. C. & Wager, T. D. Large-scale
11 automated synthesis of human functional neuroimaging data. *Nat. Methods* **8**, 665–670
12 (2011).
- 13 131. Diedrichsen, J., Yokoi, A. & Arbuckle, S. A. Pattern component modeling: A flexible
14 approach for understanding the representational structure of brain activity patterns.
15 *Neuroimage* **180**, 119–133 (2018).
- 16 132. Rouder, J. N., Morey, R. D., Speckman, P. L. & Province, J. M. Default Bayes factors for
17 ANOVA designs. *J. Math. Psychol.* **56**, 356–374 (2012).
- 18 133. Nichols, T. E. & Holmes, A. P. Nonparametric permutation tests for functional
19 neuroimaging: a primer with examples. *Hum. Brain Mapp.* **15**, 1–25 (2002).
- 20 134. Belsley, D. A., Kuh, E. & Welsch, R. E. Wiley Series in Probability and Statistics.
21 *Regression Diagnostics: Identifying Influential Data and Sources of Collinearity* 293–300
22 (1980).
- 23