

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 (Badre et al., 2021;
4 Danielmeier et al., 2011; Egner, 2008; Friedman and Miyake, 2017; Musslick et al., 2015; Ritz et
5 al., 2022a). For example, someone may respond to a mistake by becoming more cautious,
6 enhancing task-relevant processing, or suppressing task-irrelevant processing (Danielmeier and
7 Ullsperger, 2011), and previous work has shown that people simultaneously deploy multiple
8 such strategies at the same time in response to different task demands (Danielmeier et al., 2011;
9 Fischer et al., 2018; Leng et al., 2021; Ritz and Shenhav, 2021). Flexibly coordinating multiple
10 cognitive processes requires a control system that can monitor multiple forms of task demands
11 and deploy multiple forms of control (also referred to as the necessity for *observability* and
12 *controllability*; (Kalman, 1960)). These monitoring and regulation processes are fundamental to
13 control, and are thought to be underpinned by distinct cingulo-opercular and frontoparietal neural
14 systems (Gordon et al., 2017; Gottlieb et al., 2020; Gratton et al., 2016; Kerns et al., 2004;
15 MacDonald et al., 2000; Menon and D’Esposito, 2021; Shenhav et al., 2013; Smith et al., 2019).
16 However, much is still unknown about how multiple forms of control are represented across
17 these domains.

18
19 Past research on the neural mechanisms of cognitive control has often sought to identify
20 representations that integrate over multiple different sources of task demands (i.e., represent
21 these different sources in *alignment*). For instance, previous studies has proposed that dorsal
22 anterior cingulate cortex (dACC) tracks integrative features like response conflict, effort, value,
23 error likelihood, and time-on-task (Brown and Braver, 2005; Fu et al., 2022; Grinband et al.,
24 2011; Kragel et al., 2018; Mumford et al., 2023; Rushworth and Behrens, 2008; Vermeulen et
25 al., 2020; Yarkoni et al., 2009). Because they integrate over different task features instead of
26 differentiating between them, these forms of ‘aligned encoding’ (Figure 1a) are ill-suited for
27 carrying out multidimensional control. Multidimensional cognitive control requires independent
28 representations that can track multiple sources of difficulty and regulate multiple cognitive
29 processes (e.g., prioritize multiple sources of information (Pylyshyn and Storm, 1988)).

30
31 An alternative to aligned encoding – one that would allow the brain to separately control
32 multiple processes – is *independent* encoding, which can come in at least two forms. One way
33 the brain can have independent representations is by encoding different task features in spatially
34 segregated neural populations (‘segregated encoding’; Figure 1b). For example, past work has
35 shown that different subregions within dACC encode distinct task demands, including various
36 forms of errors and processing conflict (Beldzik and Ullsperger, 2023; Fu et al., 2019; Shenhav
37 et al., 2018; Taren et al., 2011; Venkatraman et al., 2009; Zarr and Brown, 2016). The brain can
38 instead have independent representations that are distributed across units within the same
39 population, as has also been observed in dACC (Ebitz et al., 2020; Flesch et al., 2022; Minxha et

1 al., 2020). Within a shared population, independent encoding of information occurs along a set of
2 orthogonal dimensions or *subspaces* (Figure 1c, ‘subspace encoding’; (Cunningham and Yu,
3 2014; Ebitz and Hayden, 2021; Mante et al., 2013; Rigotti et al., 2013)). Despite this exciting
4 recent work, it remains unclear to what extent different components of the cognitive control
5 system leverage these aligned, segregated, or orthogonal encoding strategies for monitoring
6 multiple task demands and prioritizing multiple sources of information.

7
8 To gain new insight into the representations supporting cognitive control, we drew upon two key
9 innovations. First, we leveraged an experimental paradigm we developed to tag multiple control
10 processes (Ritz and Shenhav, 2021). Building on prior work (Danielmeier et al., 2011; Kayser et
11 al., 2010b; Mante et al., 2013; Shenhav et al., 2018), this task incorporates elements of
12 perceptual decision-making (discrimination of a target feature) and inhibitory control
13 (overcoming a salient and prepotent distractor). We have previously shown that we can
14 separately tag target and distractor processing from participants’ performance on this task, and
15 that target and distractor processing are independently controlled. For example, participants
16 adjust target and distractor sensitivity in response to distinct task demands (e.g., previous conflict
17 or incentives; (Ritz and Shenhav, 2021)). In conjunction with this process-tagging approach, our
18 second innovation was to develop a novel multivariate fMRI analysis for measuring relationships
19 between feature encoding (i.e., *encoding geometry*). Extending recent statistical approaches in
20 systems neuroscience (Bernardi et al., 2020; Ebitz et al., 2020; Panichello and Buschman, 2021),
21 we combined the strengths of multivariate encoding analyses and representation similarity
22 analyses into a method we call ‘Encoding Geometry Analysis’ (EGA). We used EGA to
23 characterize whether putative markers of monitoring and prioritization leverage independent
24 representations for targets and distractors.

25
26 In brief, we found that key nodes within the cognitive control network use orthogonal
27 representations of target and distractor information to support cognitive control. In the dorsal
28 anterior cingulate cortex (dACC), encoding of target and distractor difficulty was spatially
29 segregated and arranged along a rostrocaudal gradient. By contrast, in the intraparietal sulcus
30 (IPS), encoding of target and distractor coherence was encoded along orthogonal neural
31 subspaces. These regional distinctions are consistent with hypothesized roles in planning and
32 implementing (multivariate) attentional policies (Gottlieb et al., 2020; Shenhav et al., 2013).
33 Furthermore, we found that coherence encoding depended on control demands, and was aligned
34 with both task performance and frontoparietal activity, consistent with these coherence
35 representations playing a critical role in cognitive control (e.g., feature prioritization). Together,
36 these results suggest that cognitive control uses representational formats that allow the brain to
37 monitor and control multiple streams of information processing.

38

1 Results

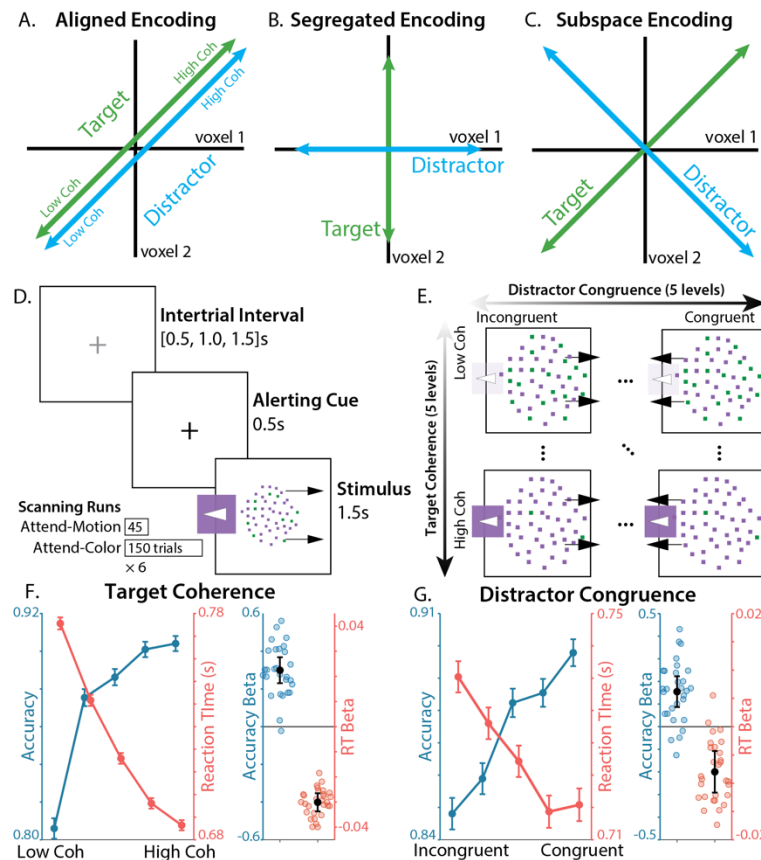
2 Task overview

3 Twenty-nine human participants performed the Parametric Attentional Control Task (PACT;
4 (Ritz and Shenhav, 2021) during fMRI. On each trial, participants responded to an array of
5 colored moving dots (colored random dot kinematogram; Figure 1d). In the critical condition
6 (Attend-Color), participants respond with a left/right keypress based on which of two colors were
7 in the majority. In alternating scanner runs, participants instead responded based on motion
8 (Attend-Motion), which was designed to be less control-demanding due to the (Simon-like)
9 congruence between motion direction and response hand (Danielmeier et al., 2011; Ritz and
10 Shenhav, 2021). Across trials, we independently and parametrically manipulated target and
11 distractor information across five levels of target coherence (e.g., percentage of dots in the
12 majority color, regardless of which color) and distractor congruence (e.g., percentage of dots
13 moving either in the congruent or incongruent direction relative to the correct color response;
14 Figure 1e). This task allowed us to ‘tag’ participants’ sensitivity to each dimension by measuring
15 behavioral and neural responses to independently manipulated target and distractor features.
16 Unlike a similar task used to study post-error adjustments (Danielmeier et al., 2011), our
17 parametric manipulation of target and distractor coherence allows us to better measure feature-
18 specific representations. Unlike similar tasks used to study contextual decision-making (Mante et
19 al., 2013; Shenhav et al., 2018; Takagi et al., 2021), this task pits more control-demanding
20 responses (towards color) against more automatic responses (towards motion), allowing
21 comparisons between Attend-Color and Attend-Motion tasks to isolate the contributions of
22 cognitive control (Cohen et al., 1990; Woolgar et al., 2011a).

23 Behavior

24 Participants had overall good performance on the task, with a high level of accuracy (median
25 Accuracy = 89%, IQR = [84% - 92%]), and a low rate of missed responses (median lapse rate =
26 2%, IQR = [0% - 5%]). We used mixed effects regressions to characterize how target coherence
27 and distractor congruence influenced participants’ accuracy and log-transformed correct reaction
28 times. Replicating previous behavioral findings using this task, participants were sensitive to
29 both target and distractor information (Ritz and Shenhav, 2021). When target coherence was
30 weaker, participants responded slower ($t_{(27.6)} = 16.1, p = 1.60 \times 10^{-15}$) and less accurately ($t_{(28)} = -$
31 $8.90, p = 1.19 \times 10^{-9}$; Figure 1f). When distractors were more incongruent, participants also
32 responded slower ($t_{(28.8)} = 5.09, p = 2.15 \times 10^{-5}$) and less accurately ($t_{(28)} = -4.66, p = 6.99 \times 10^{-5}$;
33 Figure 1g). Also replicating prior findings with this task, interactions between targets and
34 distractors were not significant for reaction time ($t_{(28.2)} = 0.143, p = .887$) and had a weak
35 influence on accuracy ($t_{(28)} = 2.36, p = .0257$), with model omitting target-distractor interactions
36 providing a better complexity-penalized fit (RT Δ AIC = 17.7, Accuracy Δ AIC = 1.38).

1



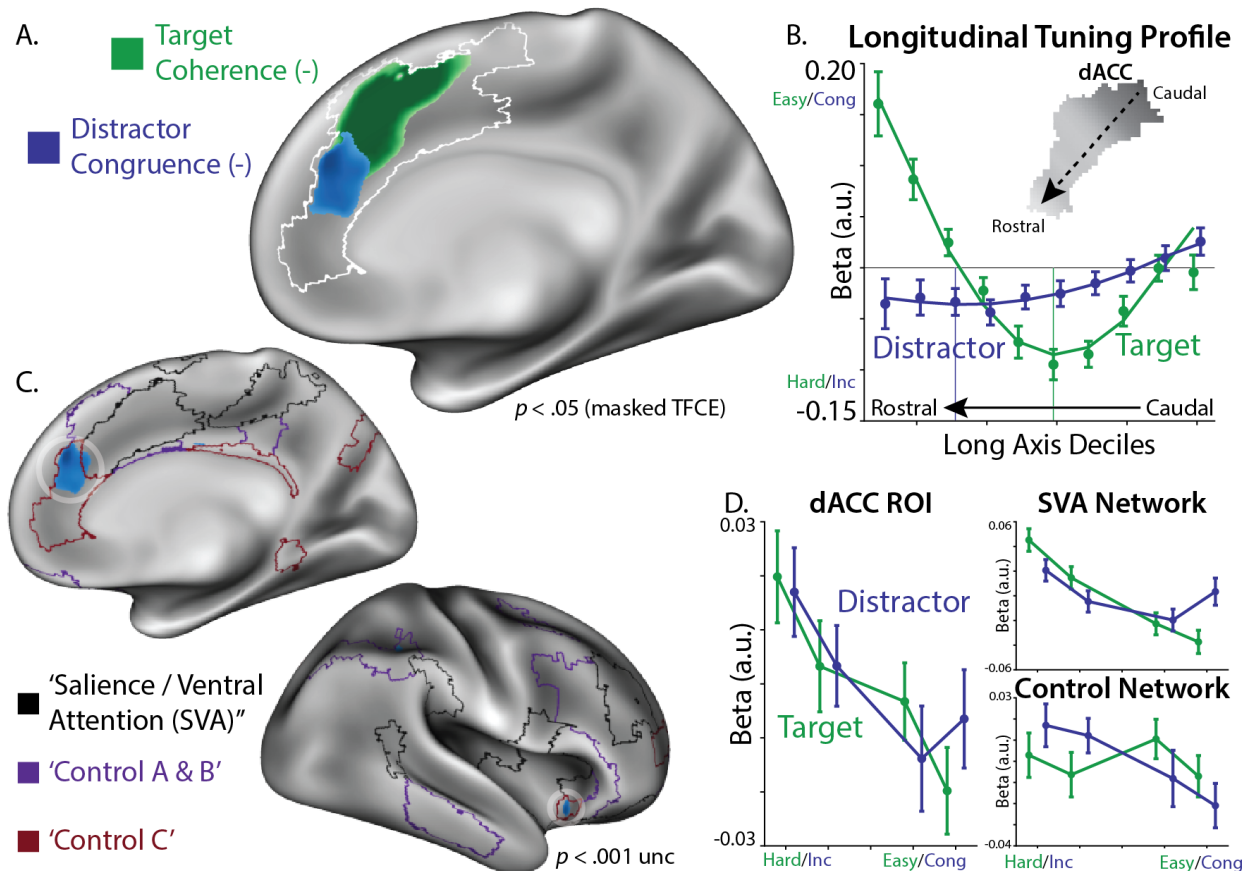
2

3 **Figure 1. Task and Behavior.** A-C) Three hypothesized encoding schemes. A) In *aligned encoding* features are
4 represented similarly, e.g., encode performance variables like error likelihood or time-on-task. B) In *segregated*
5 *encoding* features are encoded independently, in distinct voxel populations (i.e., voxel-level pure selectivity (Rigotti
6 et al., 2013)). C) In *subspace encoding*, features are encoding independently, in overlapping voxel populations (i.e.,
7 voxel-level mixed selectivity). D) Participants responded to a color-motion random dot kinematogram (RDK) with a
8 button press. Participants either responded to the left/right motion direction of the RDK (Attend-Motion runs) or
9 based on the majority color (Attend-Color runs; critical condition). E) We parametrically and independently
10 manipulated target coherence (% of dots in the majority color) and distractor congruence (motion coherence signed
11 relative to the target response). F) Participants were faster and more accurate when the target was more coherent. G)
12 Participants were faster and more accurate when the distractor was more congruent with the target. Error bars on line
13 plots reflect within-participant SEM, error bars for regression fixed-effect betas reflect 95% CI.

14 Segregated encoding of target and distractor difficulty in dACC

15 Past work has separately shown that the dACC tracks task demands related to perceptual
16 discrimination (induced in our task when target information is weaker) and related to the need to
17 suppress a salient distractor (induced in our task when distractor information is more strongly
18 incongruent with the target (Nee et al., 2007; Shenhav et al., 2018, 2013; Taren et al., 2011;
19 Venkatraman et al., 2009)). Our task allowed us to test whether these two sources of increasing
20 control demand are tracked within common regions of dACC (reflecting an aggregated
21 representation of multiple sources of task demands), or whether they are tracked by separate

1 regions (potentially reflecting a specialized representation according to the nature of the
 2 demands).
 3
 4 Targeting a large region of dACC – a conjunction of a cortical parcellation with a meta-analytic
 5 mask for ‘cognitive control’ (see ‘fMRI univariate analyses’ in Methods) – we found spatially
 6 distinct signatures of target difficulty and distractor congruence within dACC. In caudal dACC,
 7 we found significant clusters encoding the parametric effect of target difficulty (Figure 2a;
 8 negative effect of target coherence in green), and in more rostral dACC we found clusters
 9 encoding parametric distractor incongruence (negative effect of distractor congruence in blue).
 10 Supporting this dissociation, the spatial patterns of target and distractor regression weights were
 11 uncorrelated across dACC voxels ($t_{(28,0)} = 1.32, p = .197, \log\text{BF} = -0.363$). These analyses
 12 control for omission errors, and additionally controlling for commission errors produced the
 13 same whole-brain pattern at a reduced threshold (see Supplementary Figure 1). While distractor
 14 congruence was marginally significant when correcting for multiple comparisons across the
 15 entire brain (one-sided $p = .08$, whole-brain TFCE), extensive previous research predicts
 16 congruence effects in this ROI and in this direction (Kragel et al., 2018; Nee et al., 2007;
 17 Shenhav et al., 2013), suggesting that a whole-brain corrected estimate is overly conservative.
 18



19
 20 **Figure 2. Distinct coding of target and distractor difficulty in dACC.** **A)** We looked for linear target coherence and
 21 distractor congruence signals within an a priori dACC mask (white outline; overlapping Kong22 parcels and medial

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

13
14 To further quantify how feature encoding changed along the longitudinal axis of dACC, we used
15 principal component analysis to extract the axis position of dACC voxels (see ‘dACC
16 longitudinal axis analyses’ in Methods), and then regressed target and distractor beta weights
17 onto these axis scores. We found that targets had stronger difficulty coding in more caudal
18 voxels ($t_{(27.9)} = 3.74, p = .000840$), with a quadratic trend ($t_{(26.5)} = 4.48, p = .000129$; Figure 2b).
19 In line with previous work on both perceptual and value-based decision-making (Clairis and
20 Pessiglione, 2022; Fleming et al., 2018; Shenhav et al., 2018, 2016; Shenhav and Karmarkar,
21 2019), we found that signatures of target discrimination difficulty (negative correlation with
22 target coherence) in caudal dACC were paralleled by signals of target discrimination *ease*
23 (positive correlation with target coherence) within the rostral-most extent of our dACC ROI
24 (Supplementary Figure 2). In contrast to targets, distractors had stronger incongruence coding in
25 more rostral voxel ($t_{(28.0)} = -3.26, p = .00294$), without a significant quadratic trend. We used
26 participants’ random effects terms to estimate the gradient location where target and distractor
27 coding were at their most negative, finding that the target minimum was significantly more
28 caudal than the distractor minimum (signed-rank test, $z_{(28)} = 2.41, p = .0159$). Target and
29 distractor minima were uncorrelated across subjects ($r_{(27)} = .0282, p = .880, \log\text{BF} = -0.839$),
30 again consistent with independent encoding of targets and distractors.

31
32 As additional evidence that target-related and distractor-related demands have a dissociable
33 encoding profile, we found that the crossover between target and distractor encoding in dACC
34 occurred at the boundary between two well-characterized functional networks (Kong et al., 2021;
35 Schaefer et al., 2018; Yeo et al., 2011). Whereas distractor-related demands were more strongly
36 encoded rostrally in the Control Network (particularly within regions of dACC and insula
37 corresponding to the ‘Control C’ Sub-Network; (Kong et al., 2021, 2019)), target-related
38 demands were more strongly encoded caudally within the ‘Salience / Ventral Attention (SVA)’
39 Network (Figure 2C-D). Including network membership alongside long axis location predicted
40 target and distractor encoding better than models with either network membership or axis
41 location alone ($\Delta\text{BIC} > 1675$).

1 Subspace encoding of target and distractor coherence in intraparietal 2 sulcus

3 We found that dACC appeared to dissociably encode target and distractor difficulty through
4 spatially segregated encoding, consistent with a role in monitoring different task demands and/or
5 specifying different control signals (Shenhav et al., 2013). To identify neural mechanisms for the
6 implementation of this control through the prioritization targets versus distractors, we next tested
7 for regions that encode target and distractor coherence (the amount of information in a feature,
8 regardless of which response it supports). Based on previous research, we might expect to find
9 this form of selective attention in posterior parietal cortex (Bisley and Mirpour, 2019; Gottlieb et
10 al., 2020; Yantis and Serences, 2003). We explored whether target and distractor coherence share
11 a common neural code (e.g., as a global index of spatial salience), compared to where these
12 features are encoded distinctly (e.g., as separate targets of control).

13
14 An initial whole-brain univariate analysis showed that overlapping regions throughout occipital,
15 parietal, and prefrontal cortices track the feature coherence (proportion of dots in the majority
16 category) for both targets and distractors (Figure 3a; conjunction in orange). These regions
17 showed elevated responses to lower target coherence and higher distractor coherence, potentially
18 reflecting the relevance of each feature for task performance. Note that in contrast to distractor
19 congruence, distractor *coherence* had an inconsistent relationship with task performance (RT:
20 $t_{(27.0)} = 2.08, p = .048$; Accuracy: $t_{(28)} = -0.845, p = .406$), suggesting that these neural responses
21 are unlikely to reflect task difficulty per se.

22
23 While these univariate activations point towards widespread and coarsely overlapping encoding
24 of the feature coherence (potentially consistent with aligned encoding; Figure 1a), they lack
25 information about how these features are encoded at finer spatial scales. To interrogate the
26 relationship between target and distractor encoding, we developed a multivariate analysis that
27 combines multivariate encoding analyses with pattern similarity analyses, which we term
28 Encoding Geometry Analysis (EGA). Whereas pattern similarity analyses typically quantify
29 relationships between representations of specific stimuli or responses (e.g., whether they could
30 be classified, (Kriegeskorte and Diedrichsen, 2019)), EGA characterizes relationships between
31 encoding subspaces (patterns of contrast weights) across different task features, consistent with
32 recent analyses trends in systems neuroscience (Bernardi et al., 2020; Cohen and Maunsell,
33 2010; Ebitz et al., 2020; Flesch et al., 2022; Kimmel et al., 2020; Libby and Buschman, 2021). A
34 stronger correlation between encoding subspaces (either positive or negative) indicates that
35 features are similarly encoded (i.e., that their representations are aligned and thus confusable by
36 a decoder; Figure 1a), whereas weak correlation indicate that these representations are
37 orthogonal (and thus distinguishable by a decoder; (Kriegeskorte and Diedrichsen, 2019)). In
38 contrast to standard pattern similarity, the sign of these relationships is interpretable in EGA,
39 reflecting how features are coded relative to one another. Compared to standard encoding

1 analyses, EGA is less sensitive to noise (Supplementary Figure 3). We estimated this encoding
2 alignment within each parcel, correlating unsmoothed and spatially pre-whitened patterns of
3 parametric regression betas across scanner runs to minimize spatiotemporal autocorrelation
4 (Diedrichsen and Kriegeskorte, 2017; Nili et al., 2014; Walther et al., 2016). This cross-validated
5 similarity further allowed us to anchor our analysis on the measurement reliability of encoding
6 profiles (i.e., the self-correlation of encoding patterns across cross-validation folds (Spearman,
7 1987; Thornton and Mitchell, 2017)).

8
9 Focusing on regions that encoded both target and distractor information (parcels where both
10 group-level $p < .001$), EGA revealed clear dissociations between regions that represent these
11 features in alignment versus orthogonally. Within visual cortex and the superior parietal lobule
12 (SPL), target and distractor representations demonstrated significant negative correlations
13 (Figure 3b, red), reflecting (negatively) aligned encoding. In contrast, early visual cortex and
14 intraparietal sulcus (IPS; see Figure 3c for anatomical boundaries) demonstrated target-distractor
15 correlations near zero (Figure 3b, black), suggesting encoding along orthogonal subspaces.

16
17 To bolster our interpretation of the latter findings as reflecting orthogonal (i.e., uncorrelated)
18 representations rather than merely small but non-significant correlations, we employed Bayesian
19 t-tests at the group level to estimate the relative (log-10) likelihood that these encoding
20 dimensions were orthogonal or correlated. Consistent with our previous analyses, we found
21 strong evidence for correlation (positive log bayes factors) in more medial regions of occipital
22 and posterior parietal cortex (e.g., SPL), and strong evidence for orthogonality (negative log
23 bayes factors) in more lateral regions of occipital and posterior parietal cortex (e.g., IPS; Figure
24 3D). Control analyses confirmed that coherence orthogonality was not due to encoding
25 reliability, as a similar topography was observed with disattenuated correlations (normalizing
26 correlations by their reliability; see Supplementary Figure 4). Further supporting these results,
27 our Bayes factor analyses were robust to the choice of priors (see Supplementary Figure 5).

28
29

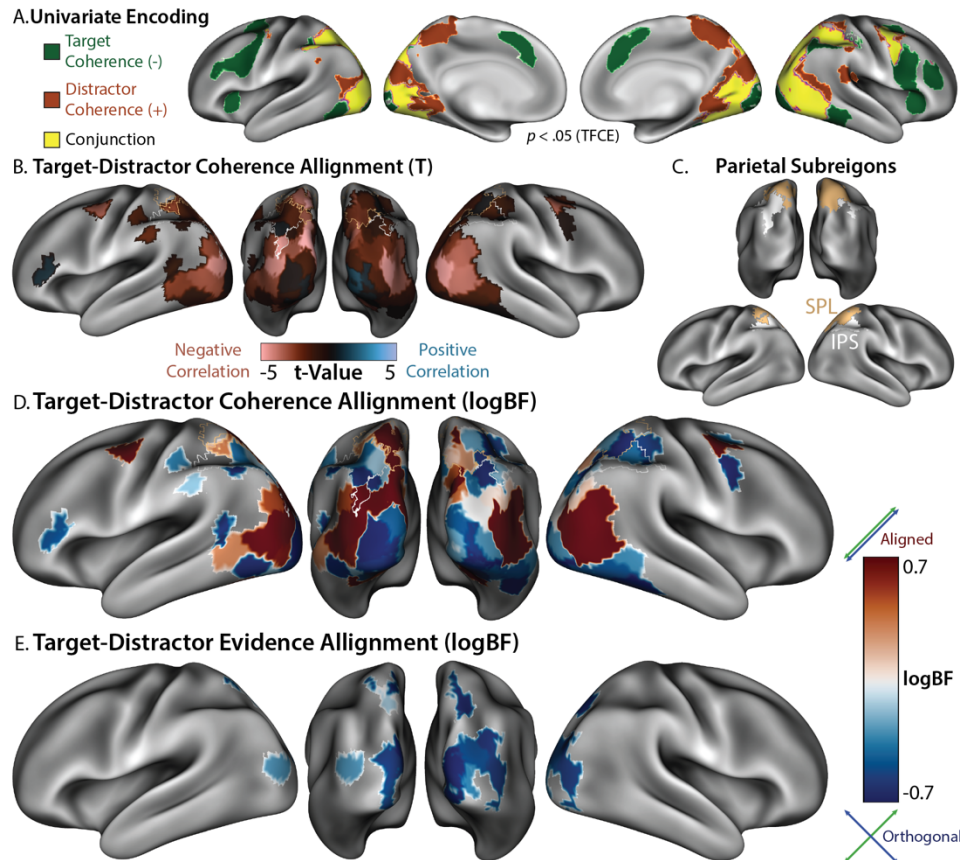
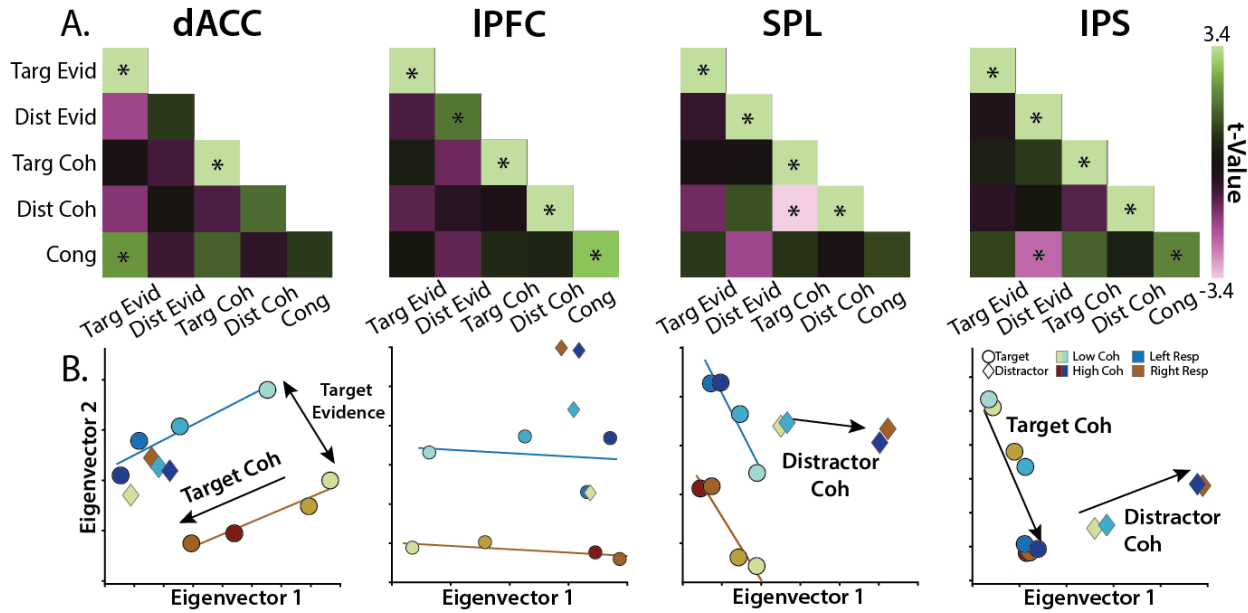


Figure 3. Encoding Geometry Analysis (EGA) dissociates target and distractor encoding. **A)** Parametric univariate responses to weak target coherence (green; percentage of dots in majority color), strong distractor coherence (orange; percentage of dots with coherent motion), and their conjunction (yellow). 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 selectivity (Rigotti et al., 2013)). 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

1 amount of evidence supporting specific response (Hunt et al., 2012; Kayser et al., 2010a, 2010b),
 2 which we found was also true for our task. In addition to encoding target and distractor
 3 coherence, SPL and visual cortex also tracked target and distractor ‘evidence’ (proportion of dots
 4 supporting a rightward vs leftward response; Figure 3e). EGA revealed orthogonal evidence
 5 representations between targets and distractors, in the same areas with aligned coherence
 6 representations (compare Figure 3d and 3e), consistent with previous observations of multiple
 7 decision-related signals in SPL (Hunt et al., 2012).
 8



9
 10 **Figure 4. Region-specific feature encoding.** **A)** Similarity matrices for dACC, IPFC, SPL, and IPS, correlating
 11 feature evidence (‘Evid’), feature coherence (‘Coh’), and feature congruence (‘Cong’). Encoding strength on
 12 diagonal (right-tailed p -value), encoding alignment on off-diagonal (two-tailed p -value). **B)** Classical MDS
 13 embedding of target (circle) and distractor (diamond) representations at different levels of evidence. Colors denote
 14 responses, hues denote coherence. GLMs: A: Feature MV, B: Evidence Levels, see Table 2.

15
 16 We complemented our whole-brain analyses with ROI analyses in areas exhibiting reliable
 17 encoding of key variables, focusing on core frontal regions linked with cognitive control (dACC
 18 and lateral PFC [IPFC]), and parietal regions linked with decision-making and attention (SPL
 19 and IPS; (Menon and D’Esposito, 2021; Shenhav et al., 2013)). Consistent with our analyses
 20 above, we found that target and distractor coherence encoding was aligned in SPL, but not in IPS
 21 (Figure 4a, compare to Figure 3d), whereas SPL encoded target and distractor evidence. Directly
 22 comparing these regions (see Table 1), we found stronger encoding of target evidence in SPL,
 23 stronger encoding of target coherence in IPS, and stronger alignment between target-distractor
 24 coherence alignment in SPL. Unlike our univariate results, we did not find distractor congruence
 25 encoding in dACC (though this was found in IPFC and IPS). Instead, dACC showed multivariate
 26 encoding of target coherence and evidence.
 27

Task Feature	SPL	IPS	SPL – IPS
Target Evidence	$t_{(28)} = 10.39, p = 4.07 \times 10^{-11}, \log\text{BF} = 8.37$	$t_{(28)} = 5.82, p = 3.00 \times 10^{-6}, \log\text{BF} = 3.81$	$t_{(28)} = 3.89, p = .000562, \log\text{BF} = 1.75$
Distractor Evidence	$t_{(28)} = 4.42, p = 1.34 \times 10^{-4}, \log\text{BF} = 2.30$	$t_{(28)} = 3.62, p = 0.0012, \log\text{BF} = 1.47$	$t_{(28)} = 0.896, p = .378, \log\text{BF} = -0.545$
Target-Distractor Evidence Alignment	$t_{(28)} = -0.703, p = .488, \log\text{BF} = -0.606$	$t_{(28)} = -0.436, p = 0.666, \log\text{BF} = -0.667$	$t_{(28)} = -0.145, p = .886, \log\text{BF} = -0.701$
Target Coherence	$t_{(28)} = 5.82, p = 2.94 \times 10^{-6}, \log\text{BF} = 3.82$	$t_{(28)} = 7.73, p = 2.00 \times 10^{-8}, \log\text{BF} = 5.83$	$t_{(28)} = -3.89, p = 9.36 \times 10^{-9}, \log\text{BF} = 6.14$
Distractor Coherence	$t_{(28)} = 8.88, p = 1.25 \times 10^{-9}, \log\text{BF} = 6.97$	$t_{(28)} = 8.53, p = 2.80 \times 10^{-9}, \log\text{BF} = 6.63$	$t_{(28)} = 1.40, p = .170, \log\text{BF} = -0.320$
Target-Distractor Coherence Alignment	$t_{(28)} = -4.75, p = 5.50 \times 10^{-5}, \log\text{BF} = 2.65$	$t_{(28)} = -1.06, p = 0.294, \log\text{BF} = -0.479$	$t_{(28)} = -2.99, p = .00580, \log\text{BF} = 0.861$

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; (Diedrichsen and Kriegeskorte, 2017; Kriegeskorte et al., 2006)). 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). In addition to visualizing cross-region dissociations in these representations, these analyses also help to validate that task features are encoded in a monotonic fashion.

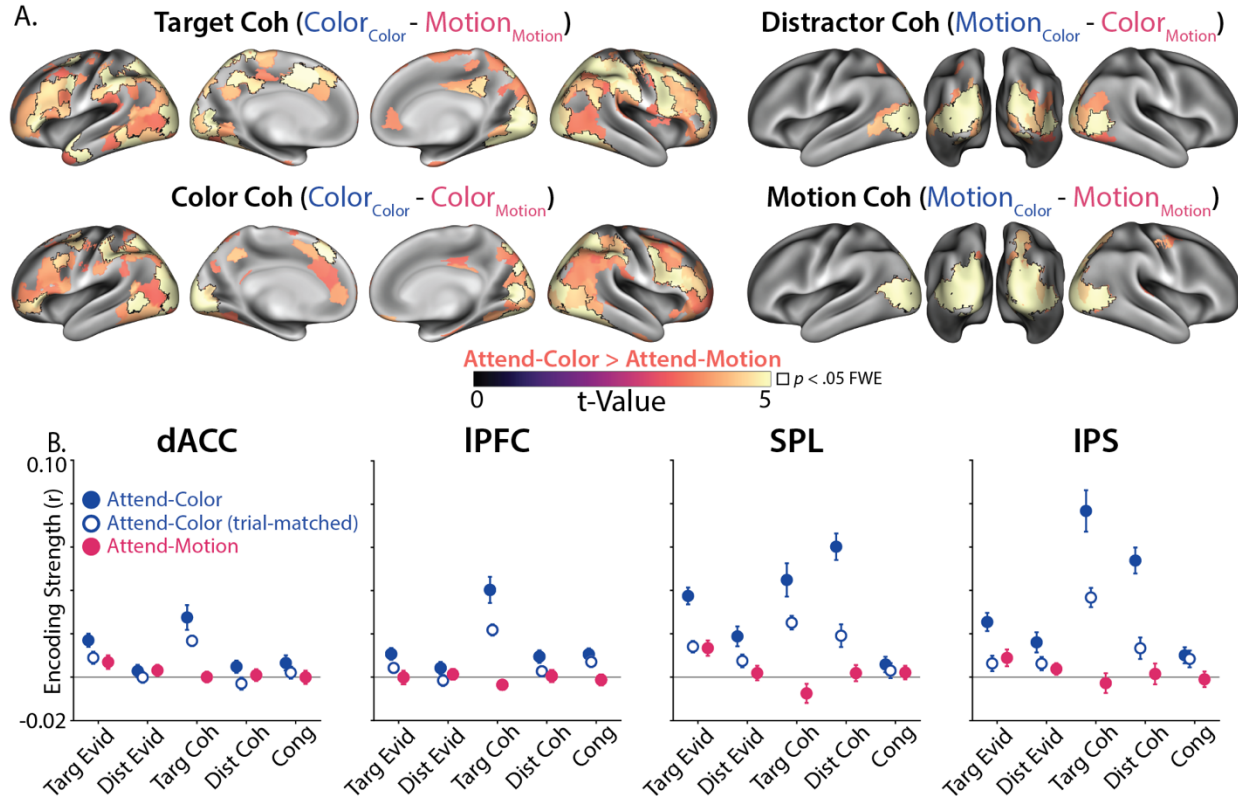
Finally, to explore the divisions between SVA and Control networks evident in the univariate analyses, we split up our two prefrontal ROIs by their network membership (Supplementary Figure 8). In dACC, we found that SVA parcels tended to have stronger feature encoding than Control parcels. Interestingly, in these SVA parcels several features were aligned with the target evidence dimension, consistent with recent human electrophysiology findings (Ebitz et al., 2020). In IPFC, we found that Control parcels, but not SVA parcels, encoded distractor congruence (Control: $t_{(28)} = 3.60$, two-tailed $p = .0012$, $\log\text{BF} = 1.45$; SVA: $t_{(28)} = 0.57$, $p = .57$,

1 logBF = -0.64; Control – SVA: $t_{(28)} = 3.27, p = .0029$, logBF = 1.12). This distractor congruence
2 encoding was present in IPFC both in ‘Control A/B’ parcels ($t_{(28)} = 3.66, p = .001$) and
3 marginally in ‘Control C’ parcels ($t_{(28)} = 1.86, p = .073$). This network-selective encoding of
4 congruence is consistent with the univariate results in dACC (see Figure 2).
5

6 Control demands dissociate coherence and evidence encoding

7 Our findings thus far demonstrate two sets of dissociations within and across brain regions. In
8 dACC, we find that distinct regions encode the control demands related to discriminating targets
9 (caudal dACC) versus overcoming distractor incongruence (rostral dACC). In posterior parietal
10 cortex, we find that overlapping regions track the coherence of these two stimulus features, but
11 that distinct regions represent these features in alignment (SPL) versus orthogonally (IPS). While
12 these findings suggest that this set of regions was involved in translating between feature
13 information and goal-directed responding, they only focus on the information that was presented
14 to the participant on a given trial. To provide a more direct link between feature-specific
15 encoding and control, we examined how the encoding of feature coherence differed between
16 matched task that placed stronger or weaker demands on cognitive control. So far, our analyses
17 have focused on conditions in which participants needed to respond to the color feature while
18 ignoring the motion feature (Attend-Color task), but on alternating scanner runs participants
19 instead responded to the motion dimension and ignored the color dimension (Attend-Motion
20 task). These tasks were matched in their visual properties (identical stimuli) and motor outputs
21 (left/right responses), but critically differed in their control demands. Attend-Motion was
22 designed to be much easier than Attend-Color, as the left/right motion directions are compatible
23 with the left/right response directions (i.e., Simon facilitation; (Danielmeier et al., 2011; Ritz and
24 Shenhav, 2021)). Comparing these tasks allows us to disambiguate bottom-up attentional
25 salience from the top-down contributions to attentional priority (Jackson et al., 2017; Woolgar et
26 al., 2011a, 2015b, 2015a).

27
28 Consistent with previous work (Ritz and Shenhav, 2021), performance on the Attend-Motion
29 task was better overall (mean RT: 565ms vs 725ms, sign-rank $p = 2.56 \times 10^{-6}$; mean Accuracy:
30 93.7% vs 87.5%, sign-rank $p = .000318$). Unlike the Attend-Color task, performance was not
31 impaired by distractor incongruence (i.e., color distractors; RT: $t_{(28)} = -1.39, p = .176$; Accuracy:
32 $t_{(28)} = 0.674, p = .506$). To investigate these task-dependent feature representations, we fit a GLM
33 that included both tasks. To control for performance differences across tasks, we only analyzed
34 accurate trials and included trial-wise RT as a nuisance covariate, concatenating RT across tasks.
35



1
 2 **Figure 5. Task-dependent encoding strength.** **A)** Across cortex, feature coherence encoding was stronger during
 3 Attend-Color than Attend-Motion, matched for the same number of trials. Attend-Color had stronger encoding when
 4 comparing target coherence (top left), distractor coherence (top right), color coherence (bottom left) and motion
 5 coherence (bottom right). Parcels are thresholded at $p < .001$ (uncorrected); outlined parcels are significant at $p < .05$
 6 I (max-statistic randomization test across all parcels). Titles are coded 'Feature_{Task}'. **B)** Target and distractor
 7 coherence information was encoded more strongly during Attend-Color than Attend-Motion in dACC, IPFC, SPL
 8 and IPS. Attend-Color encoding plotted from the whole sample (blue fill) and a trial-matched sample (first 45 trials
 9 of each run; white fill) In Attend-Motion runs, only target evidence was significantly encoded (magenta). **C)** Target
 10 and distractor coherence was not reliably encoded during the Attend-Motion task (liberally thresholded at $p < .01$
 11 uncorrected). GLM: Between-Task, see Table 2.

12
 13 Whereas the encoding of both color and motion coherence was widespread during the Attend-
 14 Color task (Figure 3), coherence encoding was consistently weaker during the less demanding
 15 Attend-Motion task (Figure 5A). Coherence encoding was weaker during Attend-Motion
 16 whether classifying according to goal-relevance (comparing targets or distractors) or the features
 17 themselves (comparing motion or color). Task-relevant ROIs revealed that coherence encoding
 18 was effectively absent during the easy Attend-Motion task (Figure 5B), suggesting that they
 19 depend on the control demands of the Attend-Color task (Woolgar et al., 2011a, 2011b).

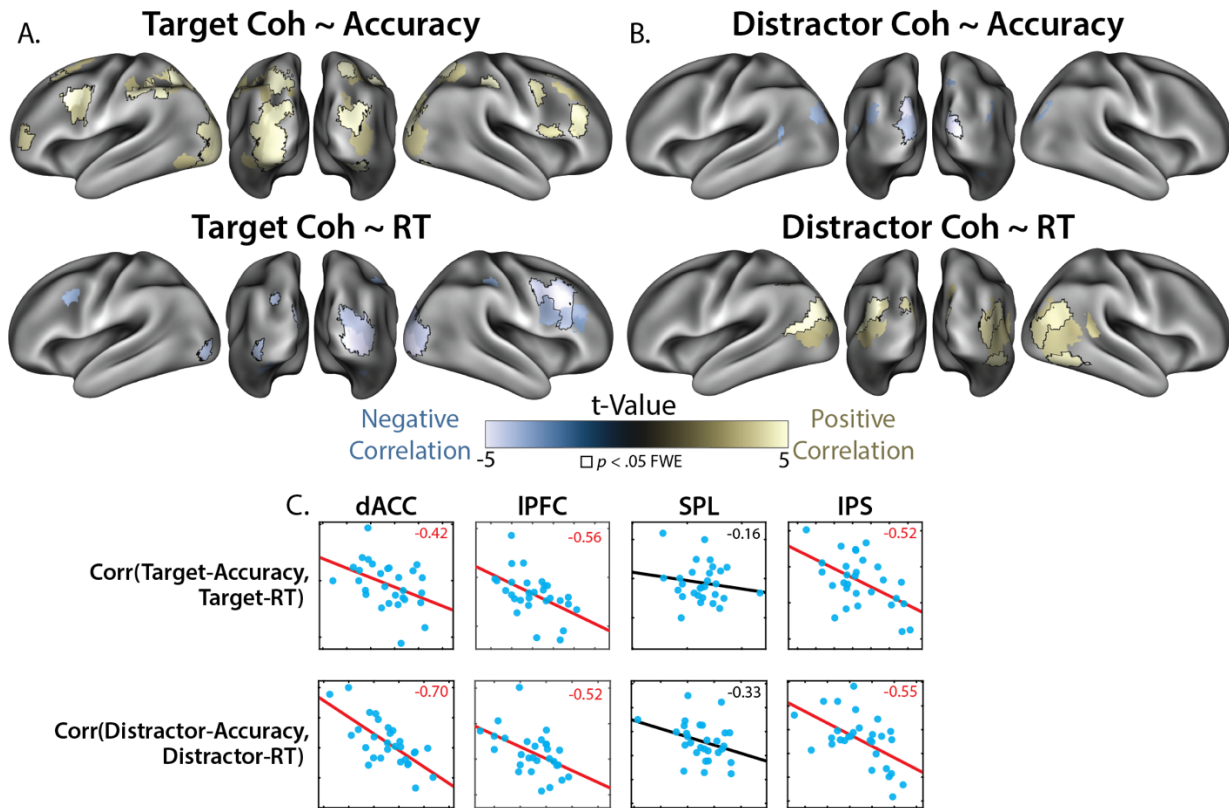
20
 21 In contrast to these stark task-related differences in coherence encoding, we found that neural
 22 encoding of the target evidence (color evidence in the Attend-Color task and motion evidence in
 23 the Attend-Motion task) was preserved across tasks, including within dACC, IPFC, SPL, and IPS
 24 (Figure 5B). Consistent with previous experiments examining context-dependent decision-

1 making (Aoi et al., 2020; Flesch et al., 2022; Jackson et al., 2017; Kayser et al., 2010b; Mante et
2 al., 2013; Pagan et al., 2022; Takagi et al., 2021), we found stronger target evidence encoding
3 relative to distractor evidence encoding, in our case in the evidence-encoding SPL (Attend-
4 Color: $t_{(28)} = 4.26$, one-tailed $p = 0.0001$; Attend-Motion: $t_{(28)} = 2.37$, one-tailed $p = 0.0124$). We
5 also found that target evidence encoding during Attend-Motion was aligned with Attend-Color,
6 both for *motion* evidence encoding ('stimulus axis'; SPL: one-tailed $p = .0236$, IPS: one-tailed p
7 = .0166) and *target* evidence encoding ('decision axis'; SPL: one-tailed $p = 1.29 \times 10^{-6}$; IPS:
8 one-tailed $p = .0005$), again in agreement with these previous experiments. Whereas our
9 experiment replicates previous observations of the neural representations supporting contextual
10 decision-making, we now extended these findings to understand how putative attention signals
11 (i.e., feature coherence) are encoded in response to the asymmetric inference that is characteristic
12 of cognitive control (Miller and Cohen, 2001).

13 Aligned encoding dimensions for feature coherence and task 14 performance

15 We next explored whether the encoding of feature coherence, seemingly in the service of
16 cognitive control, was related to how well participants performed the task. We tested this
17 question by determining whether feature coherence representations were aligned with
18 representations of behavior (i.e., alignment between stimulus and behavioral subspaces (Stringer
19 et al., 2019)). Specifically, we included trial-level reaction time and accuracy in our first-level
20 GLMs. Encoding of performance was itself highly robust: most parcels encoded reaction time
21 and accuracy, with the strongest encoding in cognitive control regions (Supplementary Figure 9).
22 Across cortex, reaction time and accuracy were negatively correlated, again most prominently
23 across the cognitive control network. To explore the behavioral relevance of coherence
24 representations, we tested whether coherence encoding was aligned with the voxel patterns
25 encoding task performance.

26
27
28



1
2 **Figure 6.** Alignment between feature and performance encoding. **A)** Alignment between encoding of target
3 coherence and performance (top row: Accuracy, bottom row: RT). **B)** Alignment between encoding of distractor
4 coherence and performance (top row: Accuracy, bottom row: RT). Across A and B, parcels are thresholded at $p <$
5 $.001$ (uncorrected, in jointly reliable parcels), and outlined parcels are significant at $p < .05$ I (max-statistic
6 randomization test across jointly reliable parcels). **C)** Individual differences in feature-RT alignment correlated with
7 feature-accuracy alignment across regions (correlation values in top right; $p < .05$ in red). See Supplementary Table
8 1 for partial correlations controlling for reliability. GLM: Performance, see Table 2.
9

10 We found that the encoding of target and distractor coherence was aligned with performance
11 across frontoparietal and visual regions (Figure 6a-b). If a regions' encoding of target coherence
12 reflects how sensitive the participant was to target information on that trial (e.g., due to top-down
13 priority), we would expect target encoding to be positively aligned with performance on a given
14 trial, such that stronger target coherence encoding is associated with better performance and
15 weaker target coherence encoding is associated with poorer performance. We would also expect
16 distractor encoding to demonstrate the opposite pattern – stronger encoding associated with
17 poorer performance and weaker encoding associated with better performance. We found
18 evidence for both patterns of feature-performance alignment across visual and frontoparietal
19 cortex: target encoding was aligned with better performance (faster RTs and higher accuracy;
20 Figure 6a), whereas distractor encoding was aligned with worse performance (slower RTs and
21 lower accuracy; Figure 6b).
22

1 Next, we examined whether performance-coherence alignment reflected individual differences in
2 participants' task performance in our main task-related ROIs (see Figures 3-4). In particular, we
3 tested whether the alignment between features and behavior reflects specific relationships with
4 speed or accuracy, or whether they reflected overall increases in evidence accumulation (e.g.,
5 faster responding and higher accuracy). Within each ROI, we correlated feature-RT alignment
6 with feature-accuracy alignment across subjects. We found that in dACC and IPS, participants
7 showed the negative correlation between performance alignment measures predicted by an
8 increase in processing speed (Figure 6c). People with stronger alignment between target
9 coherence and shorter RTs tended to have stronger alignment between target coherence and
10 higher accuracy, with the opposite found for distractors. While these between-participant
11 correlations were present within targets and distractors, we did not find any significant
12 correlations across features (between-feature: all $ps > .10$), again consistent with feature-specific
13 processing. These analyses were qualitatively similar after partialing out the reliability of
14 coherence and performance encoding, albeit with dACC and IPFC now showing marginal
15 correlations for target coherence (see Supplementary Table 1). While between-participant
16 analyses using small sample sizes warrant a note of caution, these findings are consistent across
17 features and regions. In conjunction with our within-participant evidence that feature coherence
18 representations are aligned with performance efficiency, these findings support a role for
19 coherence encoding in adaptive control.

20 Coherence encoding aligns with frontoparietal activity

21 Across frontal, parietal, and visual cortex, encoding of target and distractor coherence depended
22 on task demands and was aligned with performance. Since this widespread encoding of task
23 information likely reflects distributed network involvement in cognitive control (Corbetta and
24 Shulman, 2002; Goldman-Rakic, 1988; Miller and Cohen, 2001), we sought to understand how
25 frontal and parietal systems interact. We focused our analyses on IPS and lateral PFC (IPFC),
26 linking the core parietal site of orthogonal coherence encoding (IPS) to an prefrontal site
27 previous work suggests provides top-down feedback during cognitive control (Goldman-Rakic,
28 1988; Kastner and Ungerleider, 2000; Suzuki and Gottlieb, 2013; Yantis and Serences, 2003).
29 Previous work has found that IPS attentional biases lower-level stimulus encoding in visual
30 cortices (Kay and Yeatman, 2017; Saalmann et al., 2007), and that IPS mediates directed
31 connectivity between IPFC and visual cortex during perceptual decision-making (Kayser et al.,
32 2010b). Here, we extended these experiments to test how IPS mediates the relationship between
33 prefrontal feedback and stimulus encoding.

34
35 To investigate these putative cortical interactions, we developed a novel multivariate
36 connectivity analysis to test whether coherence encoding was aligned with prefrontal activity,
37 and whether this IPFC-coherence alignment was mediated by IPS. We first estimated the voxel-
38 averaged residual timeseries in IPFC (SPM12's eigenvariate), and then included this residual

1 timeseries alongside task predictors in a whole-brain regression analysis (Figure 7A). This
2 analysis can be schematized as:

3

$$\begin{aligned} 4 \quad \beta_{seed} &= GLM(X, Y_{seed}) \\ 5 \quad e_{seed} &= PCA(Y_{seed} - X\beta_{seed}) \\ 6 \quad \beta_{all} &= GLM([X, e_{seed}], Y_{all}) \end{aligned}$$

7

8 The GLM function performs regression using design matrix X and multivariate voxel timeseries
9 Y, and the PCA function extracts the first principal component of the residuals. Finally, we used
10 EGA to test whether there was alignment between patterns encoding IPFC functional
11 connectivity (i.e., betas from the residual timeseries predictor e_{seed}) and patterns encoding
12 target and distractor coherence. Note that while these results reflect functional connectivity, all
13 correlational measures are subject to potential confounding (Reid et al., 2019).

14

15 We found that IPFC connectivity patterns were aligned with coherence-encoding patterns in
16 visual cortex (Figure 7B). Stronger prefrontal functional connectivity was aligned with weaker
17 target coherence and stronger distractor coherence, consistent with prefrontal recruitment during
18 difficult trials. Notably, IPS connectivity was also aligned with target and distractor coherence in
19 overlapping parcels, even when controlling for IPFC connectivity. These effects were liberally
20 thresholded for visualization, as significant direct and indirect effects are not necessary for
21 significant mediation (MacKinnon et al., 2007).

22

23 Our critical test was whether IPS mediated the relationship between IPFC activity and coherence
24 encoding. We compared regression estimates between a model that only included IPFC residuals
25 ('solo' model) to a model that included both IPFC and IPS residuals ('both' model). Comparing
26 the strength of IPFC-coherence alignment with and without IPS is a test of whether parietal
27 cortex mediates IPFC-coherence alignment (MacKinnon et al., 2007). These models can be
28 schematized as:

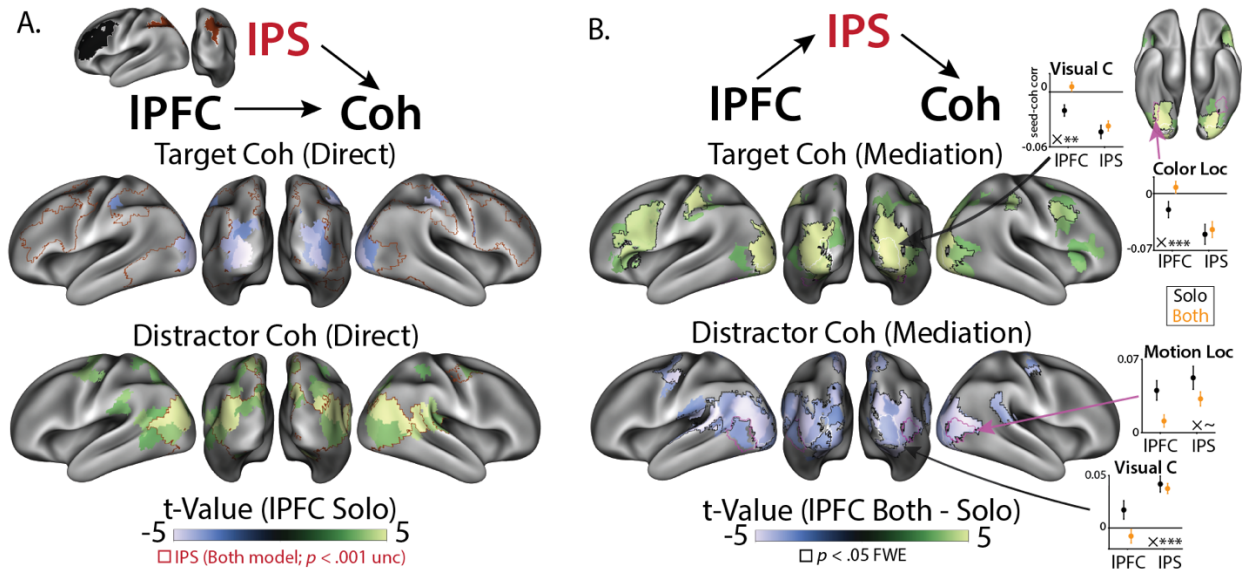
29

$$\begin{aligned} 30 \quad \beta_{solo} &= GLM([X, e_{IPFC}], Y_{all}) \\ 31 \quad \beta_{both} &= GLM([X, e_{IPFC}, e_{IPS}], Y_{all}) \end{aligned}$$

32

33 We found that this mediation was strongest in early visual cortex, where the alignment between
34 IPFC and feature coherence was reduced in a model that included IPS relative to a model without
35 IPS (Figure 7C). The negatively correlated target-IPFC relationship became more positive when
36 IPS was included (top), and the positively correlated distractor-PFC relationship became more
37 negative when IPS was included (bottom). Critically, we found that IPS reduced prefrontal-
38 coherence alignment in early visual cortex more than IPFC reduced parietal-coherence alignment
39 (Figure 7C inset; Supplementary Figure 10a-b), consistent with frontal-to-parietal directed
40 connectivity in previous research (Kayser et al., 2010b; Suzuki and Gottlieb, 2013). Looking

1 within color- and motion-sensitive parcels, determined using task-free localizer runs (see
 2 Methods), we found this mediation was robustly significant in color-sensitive cortex and
 3 marginally significant in motion-sensitive cortex. The opposite relationship, IPFC mediation of
 4 IPS connectivity, appeared in higher-level visual cortex for distractor coherence (Supplementary
 5 Figure 10c-d), though these effects were not reliable in explicit contrasts and may reflect
 6 projections from both regions.
 7



8
 9 **Figure 7. IPS mediates alignment between IPFC and feature encoding. A)** Connectivity patterns from IPFC (color)
 10 and IPS (red outline) were aligned with target and distractor coherence patterns ($p < .001$ uncorrected, in jointly
 11 reliable parcels). IPS effects are outlined to show overlap, with all effects in a consistent direction to IPFC. **B)** IPFC-
 12 feature alignment contrasted between IPFC-only model ('Solo') and IPFC + IPS model ('Both'). Including IPS
 13 reduced the alignment between IPFC and feature encoding (compare the sign of the main effect in B to the contrast
 14 in C). Parcels are thresholded at $p < .001$ (uncorrected, jointly reliable parcels), and outlined parcels are significant
 15 at $p < .05$ I (max-statistic randomization test across jointly reliable parcels). Inset graphs: seed-coherence alignment
 16 in Solo models (black) and Both model (orange) across visual regions. 'Visual C' is defined by our parcellation
 17 (Kong et al., 2021); Color and Motion localizers are parcels near the peak response identified during feature
 18 localizer runs (see Methods). In general, IPFC alignment was more affected by IPS than IPS alignment was affected
 19 by IPFC (inset 'X': difference of differences; $\sim p < .10$, $* p < .01$, $*** p < .001$). GLM: Performance CX, see Table
 20 2.
 21

22 While we were primarily interested in alignment with IPFC due to previous work implicating
 23 these regions in top-down control (for reviews, see (Friedman and Robbins, 2021; Shenhav et al.,
 24 2013)), for completeness we also examined how different subnetworks in both IPFC and dACC
 25 aligned with coherence encoding. In IPFC, we found that SVA and Control subnetworks had
 26 similar patterns of alignment (Supplementary Figure 11). In dACC we found that the SVA
 27 subnetwork had a qualitatively similar profile of coherence alignment as IPFC, but this alignment
 28 was absent in the Control subnetwork. Whereas this seed-coherence alignment was similar

1 across IPFC and SVA dACC, unlike IPFC we found that SVA dACC failed to demonstrate
2 strong evidence for mediation by IPS (Supplementary Figure 12).

3
4 A final set of analyses examined whether SPL and IPS demonstrated different patterns of task-
5 related functional connectivity with other regions, given that we found that these regions
6 differentially encoded evidence and coherence. When seeding our connectivity analyses with
7 SPL activity, we found that SPL activity aligned with evidence encoding in bilateral motor
8 cortex (Supplementary Figure 13). In contrast, IPS activity did not significantly align with
9 evidence encoding, and this seed-evidence alignment in motor cortex was stronger for SPL than
10 IPS, consistent with a putative role for SPL in response selection (Hunt et al., 2012).

11

12 Discussion

13 In this experiment, we explored whether neural control systems use representations with the
14 same dimensionality as the processes they regulate (Badre et al., 2021; Kalman, 1960; Ritz et al.,
15 2022a). Inspired by behavioral evidence that participants can independently control their
16 sensitivity to targets and distractors (Ritz and Shenhav, 2021), we set out to understand whether
17 the neural correlates of monitoring and prioritization leverage independent encoding for feature-
18 selective control (Figures 1a-c). We found that key nodes of canonical cognitive control
19 networks had orthogonal neural representations of targets and distractors. Within dACC,
20 orthogonal representations of target and distractor difficulty arose from segregated encoding
21 along a rostrocaudal axis. Within IPS, orthogonal representations of target and distractor
22 coherence arose from orthogonal subspaces in overlapping voxels. Consistent with a role in
23 attentional priority, coherence representations depended on control demands, task performance,
24 and frontoparietal activity. Together, these results reveal a neural mechanism for how cognitive
25 control prioritizes multiple streams of information during decision-making.

26

27 Neurocomputational theories have proposed that dACC is involved in planning control across
28 multiple levels of abstraction (Holroyd and McClure, 2015; Holroyd and Yeung, 2011; Shenhav
29 et al., 2013; Vassena et al., 2017). Past work has found that control abstraction is hierarchically
30 organized along dACC's rostrocaudal axis, with more caudal dACC involved in lower-level
31 action control, and more rostral dACC involved in higher-level strategy control (Shenhav et al.,
32 2018; Taren et al., 2011; Venkatraman et al., 2009; Zarr and Brown, 2016), an organization that
33 may reflect a more general hierarchy of abstraction within PFC (Badre and D'Esposito, 2009;
34 Badre and Nee, 2018; Koechlin and Summerfield, 2007; Taren et al., 2011). Consistent with this
35 account, we found that caudal dACC tracked the coherence of the target and distractor
36 dimensions, especially within the SVA network. In contrast, more rostral dACC tracked
37 incongruence between targets and distractors, especially within the Control network.
38 Speculatively, our results are consistent with caudal dACC tracking the first-order difficulty

1 arising from the relative salience of feature-specific information, and more rostral dACC
2 tracking the second-order difficulty arising from cross-feature (in)compatibility (Badre and
3 D’Esposito, 2009), the latter of which may require additional disengagement from distractor-
4 dependent attentional capture.
5

6 Whereas dACC encoded feature difficulty (e.g., distractor incongruence), in parietal cortex we
7 found overlapping representations of feature coherence (e.g., distractor coherence). In SPL,
8 features had correlated coherence encoding (similarly representing low target coherence and high
9 distractor coherence), consistent with this region’s transient and non-selective role in attentional
10 control (Esterman et al., 2009; Greenberg et al., 2010; Molenberghs et al., 2007; Serences et al.,
11 2004; Serences and Yantis, 2007; Yantis et al., 2002). In contrast, IPS had orthogonal
12 representations of feature coherence, consistent with selective prioritization of task-relevant
13 information (Adam and Serences, 2021; Greenberg et al., 2010; Jackson et al., 2017; Kay and
14 Yeatman, 2017; Molenberghs et al., 2007; Serences and Yantis, 2007; Suzuki and Gottlieb,
15 2013; Woolgar et al., 2015b, 2015a, 2011a; Yantis et al., 2002). Our previous work has
16 demonstrated behavioral evidence for independent control over target and distractor attentional
17 priority in this task (Ritz and Shenhav, 2021), with different task variables selectively enhancing
18 target or distractor sensitivity (see also (Egner, 2008; Soutschek et al., 2015)). Orthogonal
19 feature representation in IPS may offer a mechanism for this feature-selective control, consistent
20 with theoretical accounts of IPS implementing a priority map that combines stimulus- or value-
21 dependent salience with goal-dependent feedback from PFC (Bisley and Goldberg, 2010; Bisley
22 and Mirpour, 2019; Corbetta and Shulman, 2002; Gottlieb et al., 2020; Yantis and Serences,
23 2003).
24

25 In dACC, we found that target and distractor difficulty encoding was consistent with the
26 segregated encoding hypothesis, with features evoking univariate responses in distinct but
27 adjacent regions. Interestingly, we did not find corresponding encoding of distractor congruence
28 in our multivariate analyses within dACC, potentially reflecting the spatial smoothness of this
29 response. However, we did find multivariate encoding of distractor congruence in IPFC, and
30 multivariate encoding of target and distractor coherence in IPS. These multivariate profiles were
31 consistent with our subspace encoding hypothesis. The reasons for a mix of segregated and
32 subspace encoding across cortex is unclear, but this may speculatively reflect the segregation
33 across functional networks. Like in dACC, distractor congruence had stronger encoding within
34 the IPFC Control network, albeit without the feature segregation (IPFC Control parcels also
35 encoded target coherence in an orthogonal subspace). It is possible that these network
36 segregations help bind related control processes (Corbetta and Shulman, 2002; Gordon et al.,
37 2017; Menon and D’Esposito, 2021), a hypothesis that future experiments should test with
38 targeted paradigms (e.g., with subject-specific functional networks).
39

1 By comparing two different task goals (Attend-Color vs. Attend-Motion), our study was able to
2 test whether coherence representations reflect control-dependent prioritization of information
3 processing. Previous research has shown that these tasks differ dramatically in their control
4 demands (Ritz and Shenhav, 2021). As in previous work, task performance was much better in
5 Attend-Motion runs than Attend-Color runs, and participants were not sensitive to color
6 distractors. Consistent with previous work on context-dependent decision-making, target
7 evidence had similarly strong encoding across tasks, with generalizable encoding dimensions for
8 choice and motion directions (Flesch et al., 2022; Mante et al., 2013; Takagi et al., 2021). In
9 contrast to these putative decision representations, we found that coherence representations
10 disappeared in the easier Attend-Motion task. On its own, weaker encoding of color distractors in
11 Attend-Motion could be explained by the weaker bottom-up salience of the color dimension.
12 However, the stark drop in the encoding of target (motion) coherence in these blocks cannot be
13 similarly accounted for – these differences in target coherence encoding showed the opposite
14 relationship expected from salience: better encoding of low-salience color targets (hard Attend-
15 Color task) and weaker encoding of high-salience motion targets (easy Attend-Motion task).
16 Instead, this encoding profile is consistent with previous research finding that feature decoding is
17 stronger for more difficult tasks (Rust and Cohen, 2022; Woolgar et al., 2015b, 2015a, 2011a) or
18 when people are incentivized to use cognitive control (Etzel et al., 2016; Hall-McMaster et al.,
19 2019).

20
21 Critically, stimuli and responses were matched across tasks, helping to rule out alternative
22 accounts of coherence encoding based on ‘bottom-up’ stimulus salience, decision-making, or eye
23 movements. Difficulty-dependent coherence encoding may instead reflect the involvement of an
24 attention control system that can separately regulate target and distractor processing,
25 speculatively indexing the top-down ‘gain’ or ‘priority’ on these features (Bisley and Goldberg,
26 2010; Gottlieb et al., 2020; Yantis and Serences, 2003). Supporting this account, coherence
27 representations in cognitive control regions like IPS were aligned with performance
28 representations, with target encoding strength aligned with better performance and distractor
29 encoding strength aligned with poorer performance. Individual difference in feature-performance
30 alignment was correlated across features, consistent with these representations reflecting the
31 underlying processes (e.g., priority) that give rise to behavior, rather than performance
32 monitoring or surprise (which would likely have the opposite relationship, e.g., high target
33 coherence aligned with poorer performance).

34
35 Classic models of prefrontal involvement in cognitive control (Desimone and Duncan, 1995;
36 Kastner and Ungerleider, 2000; Miller and Cohen, 2001) propose that prefrontal cortex biases
37 information processing in sensory regions. In line with this macro-scale organization, we found
38 that coherence encoding in visual cortex was related to functional connectivity with the
39 frontoparietal network. In particular, coherence encoding in visual cortex was aligned with
40 patterns of functional connectivity to lateral prefrontal cortex, and this feature-seed relationship

1 was mediated by IPS. The results of this novel multivariate connectivity analysis are consistent
2 with previous research supporting a role for IPS in top-down control of visual encoding (Kay and
3 Yeatman, 2017; Lauritzen et al., 2009; Saalman et al., 2007), as well as a granger-causal PFC-
4 IPS-visual pathway during a similar decision-making task (Kayser et al., 2010b). Here, we
5 demonstrate stable ‘communication subspaces’ between visual cortex and PFC (Semedo et al.,
6 2019; Srinath et al., 2021), which can plausibly communicate feedback adjustments to feature
7 gain. With that said, while our interpretation of the direction of communication is therefore
8 supported by prior work, these connectivity methods are correlational (Reid et al., 2019), and
9 cannot rule out the possibility that our mediation findings reflect a bottom-up pattern of
10 communication (e.g., visual-IPS-PFC). The asymmetric mediation between regions (i.e., IPS
11 mediates IPFC more than IPFC mediates IPS; Supplementary Figure 10) rules out a range of
12 potential confounders, and these regions were selected based on the anatomical connectivity
13 within the frontoparietal network, notably through the superior longitudinal fasciculus (Petrides
14 and Pandya, 2006). Future research should use temporally precise neuroimaging to account for
15 directionality, causal manipulations to account for causality (e.g., (Jackson et al., 2021)), and
16 should explore the higher dimensional connectivity subspaces that link different regions (Rust
17 and Cohen, 2022; Srinath et al., 2021). These considerations notwithstanding, our findings are
18 consistent with IPS, a critical site for orthogonal feature representations, playing a key role in
19 linking prefrontal cortex with early perceptual processing.

20
21 Collectively, our findings provide new insights into how the brain may control multiple streams
22 of information processing. While evidence for multivariate control has a long history in
23 attentional tracking (Pylyshyn and Storm, 1988; Vul et al., 2009), including parametric
24 relationships between attentional load and IPS activity (Culham et al., 2001, 1998; Howe et al.,
25 2009; Jovicich et al., 2001; Ritz et al., 2022b), little is known about how the brain coordinates
26 multiple control signals (Badre et al., 2021; Ritz et al., 2022a). Future experiments should further
27 elaborate on this frontoparietal control circuit, for instance by interrogating how incentives
28 influence different task representations (Etzel et al., 2016; Hall-McMaster et al., 2019; Parro et
29 al., 2018; Peck et al., 2009; Wisniewski et al., 2015), or how neural and behavioral indices of
30 control causally depend on perturbations of neural activity (Jackson et al., 2021). 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 (Ritz and Shenhav,
33 2021; Weichart et al., 2020). In sum, this experiment provides new insights into the large-scale
34 neural networks involved in multivariate cognitive control, and points towards new avenues for
35 developing a richer understanding of goal-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 (Ritz and
11 Shenhav, 2021). On each trial, participants saw a random dot kinematogram (RDK) against a
12 black background. This RDK consisted of colored dots that moved left or right, and participants
13 responded to the stimulus with button presses using their left or right thumbs.
14

15 In Attend-Color blocks (six blocks of 150 trials), participants responded depending on which
16 color was in the majority. Two colors were mapped to each response (four colors total), and dots
17 were a mixture of one color from each possible response. Dots colors were approximately
18 isolument (uncalibrated; RGB: [239, 143, 143], [191, 239, 143], [143, 239, 239], [191, 143,
19 239]), and we counterbalanced their assignment to responses across participants.
20

21 In Attend-Motion blocks (six blocks of 45 trials), participants responded based on the dot motion
22 instead of the dot color. Dot motion consisted of a mixture between dots moving coherently
23 (either left or right) and dots moving in a random direction. Attend-Motion blocks were shorter
24 because they acted to reinforce motion sensitivity and provide a test of stimulus-dependent
25 effects.
26

27 Critically, dots always had color and motion, and we varied the strength of color coherence
28 (percentage of dots in the majority) and motion coherence (percentage of dots moving
29 coherently) across trials. Our previous experiments have found that in Attend-Color blocks,
30 participants are still influenced by motion information, introducing a response conflict when
31 color and motion are associated with different responses (Ritz and Shenhav, 2021). Target
32 coherence (e.g., color coherence during Attend-Color) was linearly spaced between 65% and
33 95% with 5 levels, and distractor congruence (signed coherence relative to the target response)
34 was linearly spaced between -95% and 95% with 5 levels. In order to increase the salience of the
35 motion dimension relative to the color dimension, the display was large (~10 degrees of visual
36 angle) and dots moved quickly (~10 degrees of visual angle per second).

1
2 Participants had 1.5 seconds from the onset of the stimulus to make their response, and the RDK
3 stayed on the screen for this full duration to avoid confusing reaction time and visual stimulation
4 (the fixation cross changed from white to gray to register the response). The inter-trial interval
5 was uniformly sampled from 1.0, 1.5, or 2.0 seconds. This ITI was relatively short in order to
6 maximize the behavioral effect, and because efficiency simulations showed that it increased
7 power to detect parametric effects of target and distractor coherence (e.g., relative to a more
8 standard 5 second ITI). The fixation cross changed from gray to white for the last 0.5 seconds
9 before the stimulus to provide an alerting cue.

10 Procedure

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

22
23 Following the twelve blocks of the scanner task, participants underwent localizers for color and
24 motion, based on the tasks used in our previous experiments (Shenhav et al., 2018). Both
25 localizers were block designs, alternating between 16 seconds of feature present and 16 seconds
26 of feature absent for seven cycles. For the color localizer, participants saw an aperture the same
27 size as the task, either filled with colored squares that were resampled every second during
28 stimulus-on ('Mondrian stimulus'), or luminance-matched gray squares that were similarly
29 resampled during stimulus-off. For the motion localizer, participants saw white dots that were
30 moving with full coherence in a different direction every second during stimulus-on, or still dots
31 for stimulus-off. No responses were required during the localizers.

32 MRI sequence

33 We scanned participants with a Siemens Prisma 3T MR system. We used the following sequence
34 parameters for our functional runs: field of view (FOV) = 211 mm × 211 mm (60 slices), voxel
35 size = 2.4 mm³, repetition time (TR) = 1.2 sec with interleaved multiband acquisitions
36 (acceleration factor 4), echo time (TE) = 33 ms, and flip angle (FA) = 62°. Slices were acquired

1 anterior to posterior, with an auto-aligned slice orientation tilted 15° relative to the AC/PC plane.
2 At the start of the imaging session, we collected a high-resolution structural MPRAGE with the
3 following sequence parameters: FOV = 205 mm × 205 mm (192 slices), voxel size = 0.8 mm³,
4 TR = 2.4 sec, TE1 = 1.86 ms, TE2 = 3.78 ms, TE3 = 5.7 ms, TE4 = 7.62, and FA = 7°. At the
5 end of the scan, we collected a field map for susceptibility distortion correction (TR = 588ms,
6 TE1 = 4.92 ms, TE2 = 7.38 ms, FA = 60°).

7 fMRI preprocessing

8 We preprocessed our structural and functional data using fMRIPrep (v20.2.6; (Esteban et al.,
9 2019) based on the Nipype platform (Gorgolewski et al., 2011). We used FreeSurfer and ANTs
10 to nonlinearly register structural T1w images to the MNI152NLin6Asym template (resampling to
11 2mm). To preprocess functional T2w images, we applied susceptibility distortion correction
12 using fMRIPrep, co-registered our functional images to our T1w images using FreeSurfer, and
13 slice-time corrected to the midpoint of the acquisition using AFNI. We then registered our
14 images into MNI152NLin6Asym space using the transformation that ANTs computed for the
15 T1w images, resampling our functional images in a single step. For univariate analyses, we
16 smoothed our functional images using a Gaussian kernel (8mm FWHM, as dACC responses
17 often have a large spatial extent). For multivariate analyses, we worked in the unsmoothed
18 template space (see below).

19 fMRI univariate analyses

20 We used SPM12 (v7771) for our univariate general linear model (GLM) analyses. Due to high
21 trial-to-trial collinearity from to our short ITIs, we performed all analyses across trials, rather
22 than extracting single-trial estimates. Our regression models used whole trials as events (i.e., a
23 1.5 second boxcar aligned to the stimulus onset). We parametrically modulated these events with
24 standardized trial-level predictors (e.g., linear-coded target coherence, or contrast-coded errors),
25 and then convolved these predictors with SPM's canonical HRF, concatenating our voxel
26 timeseries across runs. We included nuisance regressors to capture 1) run intercepts and 2) the
27 average timeseries across white matter and CSF (as segmented by fMRIPrep). To reduce the
28 influence of motion artifacts, we used robust weighted least-squares (Diedrichsen and Shadmehr,
29 2005; Jones et al., 2021), a procedure for optimally down-weighting noisy TRs.

30
31 We estimated contrast maps at the subject-level, which we then used for one-sample t-tests at the
32 group-level. We controlled for family-wise error rate using threshold-free cluster enhancement
33 (Smith and Nichols, 2009), testing whether voxels have an unlikely degree of clustering under a
34 randomized null distribution (Implemented in PALM (Winkler et al., 2014); 10,000
35 randomizations). To improve the specificity of our coverage (e.g., reducing white-matter
36 contributions) and to facilitate our inference about functional networks (see below), we limited

1 these analyses to voxels within the Kong2022 whole-brain parcellation (Kong et al., 2021;
2 Schaefer et al., 2018). This parcellation assigns regional labels to parcels (e.g., whether parcels
3 are in ‘SPL’ or ‘IPS’), which was used through-out to generate ROIs. Surface renders were
4 generated using surfplot (Gale et al., 2021; Vos de Wael et al., 2020), projecting from MNI space
5 to the Human Connectome Project’s fsLR space (164,000 vertices).

6 dACC longitudinal axis analyses

7 To characterize the spatial organization of target difficulty and distractor congruence signals in
8 dACC, we constructed an analysis mask that provided broad coverage across cingulate cortex
9 and preSMA. This mask was constructed by 1) getting a meta-analytic mask of cingulate
10 responses co-occurring with ‘cognitive control’ (Neurosynth uniformity test; (Yarkoni et al.,
11 2011), and taking any parcels from the whole-brain Schaefer parcellation (400 parcels; (Kong et
12 al., 2021; Schaefer et al., 2018) that had a 50 voxel overlap with this meta-analytic mask. We
13 used this parcellation because it provided more selective gray matter coverage than the
14 Neurosynth mask alone and it allowed us to categorize voxels membership in putative functional
15 networks.

16
17 To characterize the spatial organization within dACC, we first performed PCA on the masked
18 voxel coordinates (y and z), getting a score for eac’ voxel’s position on the longitudinal axis of
19 this ROI. We then regressed voxel’s gradient scores against their regression weights from a
20 model including linear target coherence and distractor congruence (both coded -1 to 1 across
21 difficulty levels). We used linear mixed effects analysis to partially pool across subjects and
22 accommodate within-subject correlations between voxels. Our model predicted gradient score
23 from the linear and quadratic expansions of the target and distractor betas ($\text{gradientScore} \sim 1 +$
24 $\text{target} + \text{target}^2 + \text{distractor} + \text{distractor}^2 + (1 + \text{target} + \text{target}^2 + \text{distractor} + \text{distractor}^2 |$
25 $\text{subject})$). To characterize the network-dependent organization of target and distractor encoding,
26 we complexity-penalized fits between models that either 1) predicted target or distractor betas
27 from linear and quadratic expansions of gradient scores, or 2) predicted target/distractor betas
28 from dummy-coded network assignment from the Schaefer parcellation, comparing these models
29 against a model that used both network and gradient information.

30 Encoding Geometry Analysis (EGA)

31 We adapted functions from the pcm-toolbox and rsatoolbox packages for our multivariate
32 analyses (Diedrichsen et al., 2018; Nili et al., 2014). We first fit whole-brain GLMs without
33 spatial smoothing, separately for each scanner run. These GLMs estimated the parametric
34 relationship between task variables and BOLD response (e.g., linearly coded target coherence),
35 with a pattern of these parametric betas across voxels reflecting linear encoding subspace
36 (Kriegeskorte and Diedrichsen, 2019). Within each Schaefer parcel (n=400), we spatially pre-

1 whitened these beta maps, reducing noise correlations between voxels that can inflate pattern
2 similarity and reduce reliability (Walther et al., 2016). We then computed the cross-validated
3 Pearson correlation, estimating the similarity of whitened patterns across scanner runs. We used
4 a correlation metric to estimate the alignment between encoding subspaces, rather than distances
5 between condition patterns, to normalize biases and scaling across stimuli (e.g., greater
6 sensitivity to targets vs distractors) and across time (e.g., representational drift). We found
7 convergent results when using (un-centered) cosine similarity, suggesting that our results were
8 not trivially due to parcels' univariate response, but a correlation metric had the best reliability
9 across runs. Note that this analysis approach is related to 'Parallelism Scores' (Bernardi et al.,
10 2020), but here we use parametric encoding models and emphasize not only deviations from
11 parallel/orthogonal, but also the direction of alignment between features (e.g., Figures 5 and 7).

12
13 We computed subspace alignment between contrasts of interest within each participant, and then
14 tested these against zero at the group level. Since our correlations were less than $r = |0.5|$, we did
15 not transform correlations before analysis. We used a Bayesian t-test to test for orthogonality
16 (bayesFactor toolbox in MATLAB, based on (Rouder et al., 2012)). The Bayes factor from this t-
17 test gives evidence for either non-orthogonality (BF_{10} further from zero) or orthogonality (BF_{10}
18 closer to zero, often defined as the reciprocal BF_{01}). Using a standard prior (Cauchy, width =
19 0.707), our strongest possible evidence for the orthogonality is $BF_{01} = 5.07$ or equivalently
20 $\log BF = -0.705$ (i.e., the Bayes factor when $t_{(28)} = 0$).

21
22 Our measure of encoding strength was whether encoding subspaces were reliable across blocks
23 (i.e., whether within-feature encoding pattern correlations across runs were significantly above
24 zero at the group level). We used pattern reliability as a geometric proxy for how well a linear
25 encoder would predict held-out brain data, as reliability provides the similarity between the
26 cross-validated model and the best linear unbiased estimator of the within-sample data. We
27 confirmed through simulations that pattern reliability is a good proxy for the traditional encoding
28 metric of predicting held-out timeseries (Kriegeskorte and Diedrichsen, 2019). However, we
29 found that pattern reliability is more powerful, due to it being much less sensitive to the
30 magnitude of residual variance (these two methods are identical in the noise-free case; see
31 Supplementary Figure 3).

32
33 When looking at alignment between two subspaces across parcels, we first selected parcels that
34 significantly encoded both factors ('jointly reliable parcels', both $p < .001$ uncorrected). This
35 selection process acts as a thresholded version of classical correlation disattenuation (Spearman,
36 1987; Thornton and Mitchell, 2017), and we confirmed through simulations this selection
37 procedure does not increase type 1 error rate. We corrected for multiple comparisons using non-
38 parametric max-statistic randomization tests across parcels (Nichols and Holmes, 2002). These
39 randomization tests determine the likelihood of an observed effects under a null distribution
40 generated by randomizing the sign of alignment correlations across participants and parcels

1 10,000 times. Within each randomization, we saved the max and min group-level effect sizes
2 across all parcels, estimating the strongest parcel-wise effect we'd expect if there wasn't a
3 systematic group-level effect.

4
5 Some of our first-level models had non-zero levels of multicollinearity, due to conditioning on
6 trials without omission errors or when including feature coherence and performance in the same
7 model. Multicollinearity was far below standard thresholds for concern, generally (much) less
8 than 5 for a standard threshold is 30 (ratio between largest and smallest singular values in the
9 design matrix, using MATLAB `colintest`; (Belsley et al., 1980)). However, we wanted to confirm
10 that predictor correlations wouldn't bias our estimates of encoding alignment. We simulated data
11 from a pattern component model (Diedrichsen et al., 2018) in which two variables were
12 orthogonal (generated by separate variance components with no covariance), but were generated
13 from a design matrix with correlated predictors. These simulations confirmed that cross-
14 validated similarity measures were not biased by predictor correlations (Supplementary Figure
15 14).

16
17 To provide further validation for our parametric analyses, we estimated encoding profiles using
18 an analysis with fewer parametric assumptions. First, we fit a GLM with separate predictors for
19 levels of target and distractor evidence ('Evidence Levels' GLM in Table 2). Next, we estimated
20 a traditional (cross-validated) representational dissimilarity matrix across all feature levels.
21 Finally, we visualized these encoding profiles using classical multidimensional scaling
22 (eigenvalue decomposition; see Figure 4B).

23 Multivariate Connectivity Analysis

24 To estimate what information is plausibly communicated between cortical areas, we measured
25 the alignment between multivariate connectivity patterns (i.e., the 'communication subspace'
26 with a seed region, (Semedo et al., 2019)) and local feature encoding patterns. First, we
27 residualized our Performance GLM (see Table 2) from a seed region's timeseries, and then
28 extracted the variance-weighted average timecourse (i.e., the leading eigenvariate from SPM12's
29 volume of interest function). We then re-estimated our Performance GLM, including the block-
30 specific seed timeseries as a covariate, and performed EGA between seed and coherence
31 patterns. We found convergent results when we residualized a quadratic expansion of our
32 Performance GLM from our seed region, helping to confirm that connectivity alignment wasn't
33 due to underfitting. Note that our cross-validated EGA helps avoid false positives due to any
34 correlations in the design matrix (see above). We localized this connectivity analysis to color-
35 and motion-sensitive cortex by finding the bilateral Kong22 parcels that roughly covered the area
36 of strongest block-level contrast during our localizer runs. Note that these analyses reflect
37 'functional connectivity', which is susceptible to unmodelled confounders (Reid et al., 2019).

38

1 To estimate the mediation of IPFC connectivity by IPS, we compared models in which just IPFC
 2 or just IPS were used for EGA against a model where both seeds were included as covariates in
 3 the same model (MacKinnon et al., 2007). Our test of mediation was the group-level difference
 4 in IPFC seed-coherence alignment before and after including IPS. While these analyses are
 5 inherently cross-sectional (i.e., IPFC and IPS are measured at the same time), we supplemented
 6 these analyses by showing that the mediating effect of IPS on IPFC was much larger than the
 7 mediating effect of IPFC on IPS (see Figure 7c; Supplementary Figure 10).

8
 9
 10
 11

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).

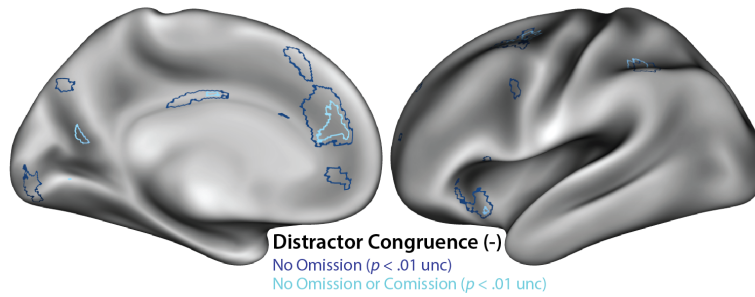
9
10 **Acknowledgements:** This work was supported by NIH grant R01MH124849 (A.S.), NSF
11 CAREER Award 2046111 (A.S.), NIH grant S10OD025181, and the C.V. Starr Postdoctoral
12 Fellowship (H.R.). We are grateful to Joonhwa Kim for her assistance in data collection, and to
13 Michael J. Frank, Matthew N. Nassar, Jonathan Cohen, Michael Esterman, Romy Frömer, Jörn
14 Diedrichsen, Apoorva Bhandari, Debbie Yee, Sam Nastase, and the rest of the Shenhav Lab for
15 helpful discussions.

16 **Conflicts of Interest:** None

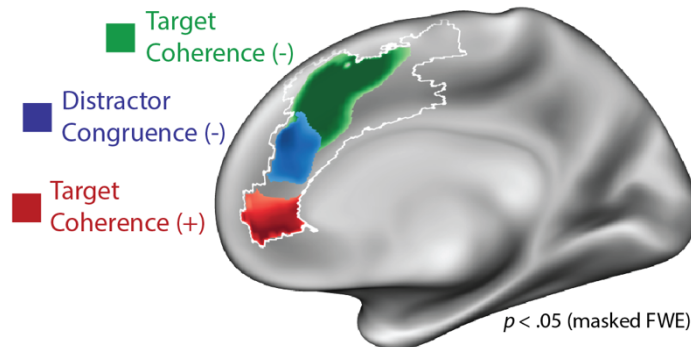
17 **Data Availability:** Data and analysis scripts will be made available upon publication.

18
19
20

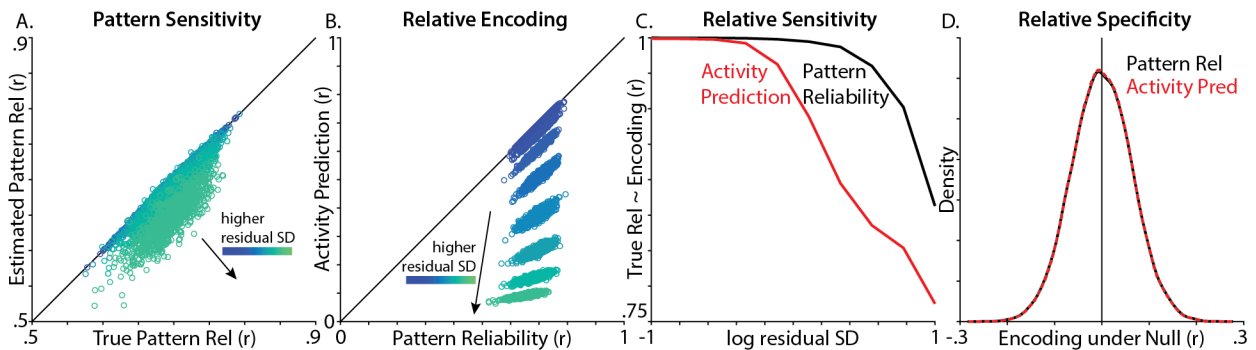
1 Supplementary Figures



2
3 **Figure S1. Error control analysis.** Distractor congruence effect when controlling for different types of errors. Our
4 primary analysis only analyzed trials without omission errors (navy), here plotted at a liberal uncorrected threshold.
5 When we analyze trials without omission errors and commission errors (cyan), we see a consistent whole-brain
6 topography, albeit at a lower statistical threshold. In both cases, relevant errors trials were included as nuisance
7 events.

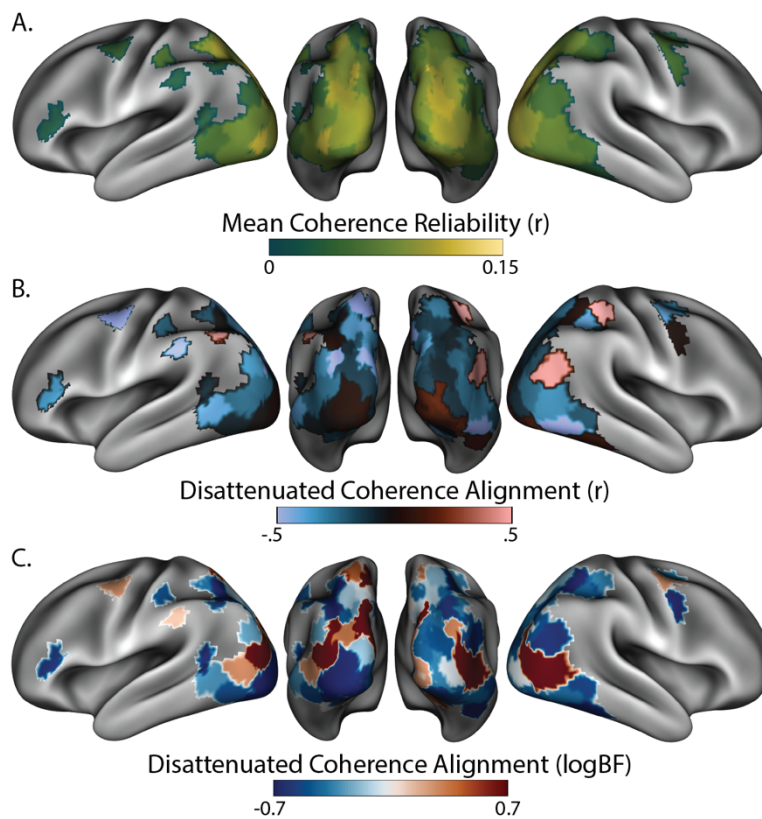


9
10 **Figure S2. Target ease.** Parametric effects of target coherence and distractor congruence, showing the rostral effect
11 of target ease (positive relationship with target coherence) in red.

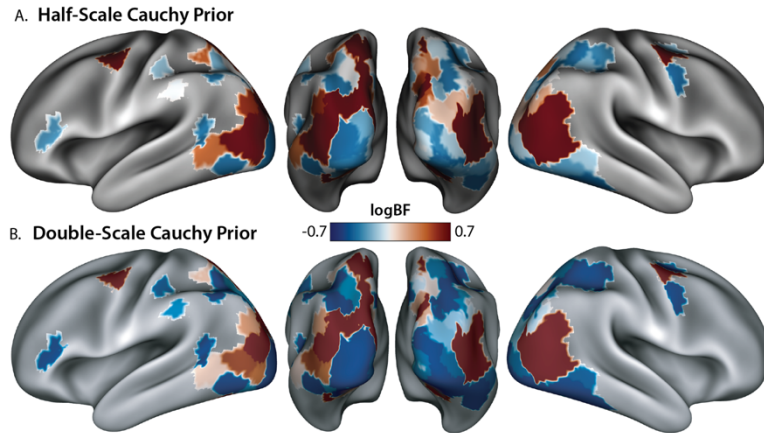


13
14 **Figure S3. Encoding Geometry Analysis (EGA) validation.** We validated how well we could recover the similarity
15 between linear Gaussian models (training: $Y = XB + \Sigma$, test: $Y' = X'B' + \Sigma$). Y is the $[1000 \times 250]$ activity
16 timeseries, X is the $[1000 \times 1]$ design matrix, B is the $[1 \times 250]$ encoding profile, and Σ reflects IID Gaussian
17 noise. In each of our 1000 simulations, we used two different methods to recover the similarity between the true
18 training encoding profile (B) and the true test encoding profile ($B' = B + \mathcal{N}(0,1)$), based on noisy activity

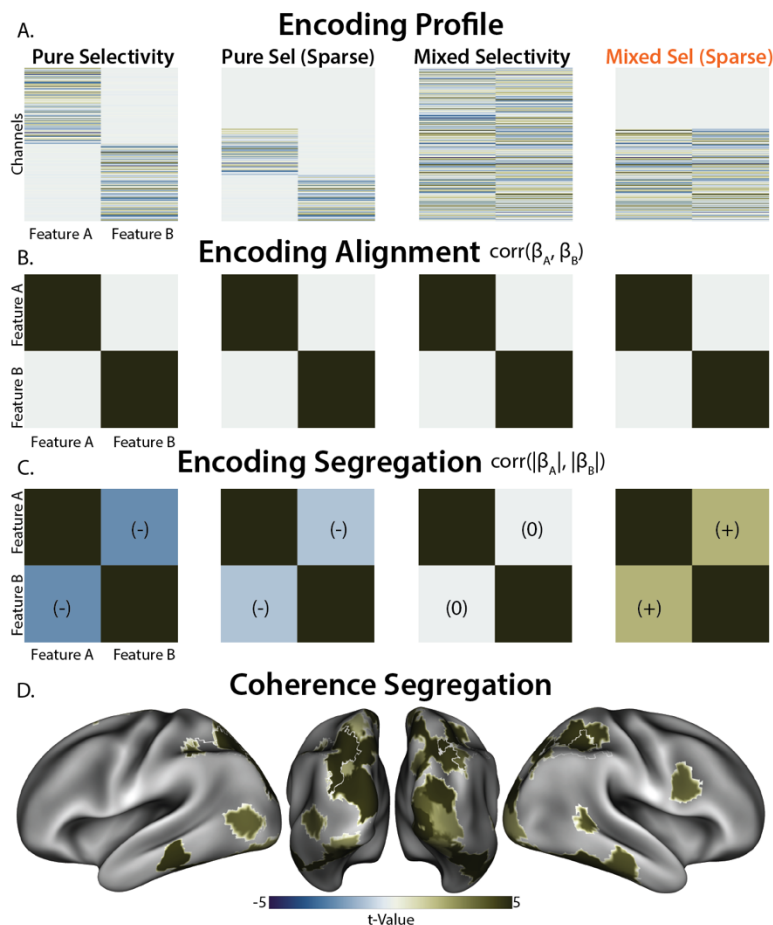
1 timeseries ($Y = XB + \mathcal{N}(0, \sigma_Y)$; $Y' = X'B' + \mathcal{N}(0, \sigma_Y)$). The first method was *pattern reliability* (i.e., our EGA
 2 method), correlating the encoding profile estimated during training ($\hat{B} = X^\dagger Y$, \dagger indicates pseudoinverse) with the
 3 encoding profile estimated during test ($\hat{B}' = X'^\dagger Y'$). The second method was *activity prediction* (i.e., the traditional
 4 encoding approach), correlating the ground-truth test activity (Y') with the predicted test activity ($\hat{Y}' = X'\hat{B}$). To
 5 simulate the high measurement noise inherent to fMRI, we compared these methods under different levels of
 6 residual SD (σ_Y). **A)** Estimated pattern reliability tracked the true pattern reliability across the full range of residual
 7 SD, with some attenuation at high levels of noise **B)** Unlike pattern reliability, activity prediction became much
 8 poorer as residual SD increased. **C)** Correlating the true pattern reliability (correlation between B and B') and
 9 estimated encoding strength (i.e., pattern reliability or activity prediction), we found pattern reliability was better
 10 correlated with the true reliability, particularly at higher levels of noise. **D)** Both methods had similar performance in
 11 the absence of a signal ($B'_{null} = \mathcal{N}(0,1)$).
 12
 13



14 **Figure S4. Reliability control analysis.** **A)** Geometric mean of target and distractor coherence reliability
 15 ($\sqrt{r_{targ} \times r_{dist}}$), plotted in the reliability-thresholded parcels used in Figure 4. Reliability provides the theoretical
 16 upper bound on correlation strength. Median across participants, excluding participants with non-positive reliability.
 17 **B)** Target-distractor correlations, normalized by target-distractor reliability (i.e., disattenuated correlations) **C)** Log
 18 bayes factors for disattenuated target-distractor correlations. Compare to Figure 4C.
 19
 20

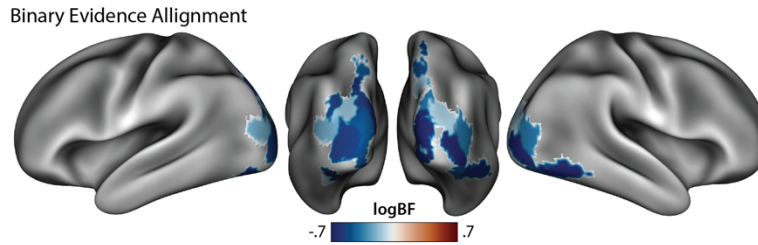


1
2 **Figure S5.** *Bayes factor prior control analysis.* **A)** Log bayes factors for target-distractor coherence alignment using
3 a narrower prior (one-half the default Cauchy scale = 0.35). Minimum logBF is -0.46 at $t_{(28)} = 0$. **B)** Same log bayes
4 factor using a wider prior (double the default Cauchy scale = 1.41). Minimum logBF = -0.99 at $t_{(28)} = 0$. Across
5 different prior parameterizations, note the similarity to Figure 4C.

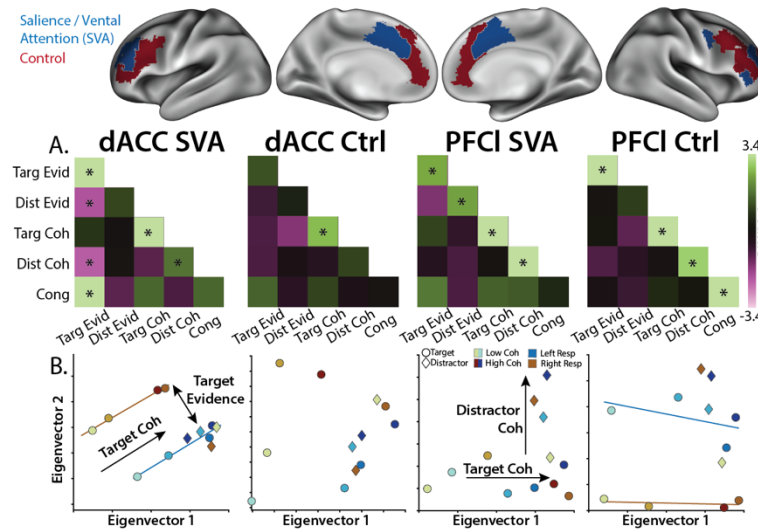


7
8 **Figure S6.** *Segregation Analysis.* **A)** We used pattern component modelling (Diedrichsen et al., 2018) to simulate
9 different candidate encoding profiles. ‘Pure Selectivity’ reflects the segregated encoding hypothesis, with different
10 voxels (rows) encoding different features (columns). ‘Mixed Selectivity’ reflects the orthogonal subspace
11 hypothesis, with the same voxels encoding both features. ‘Sparse’ models include non-selective voxels. **B)** By

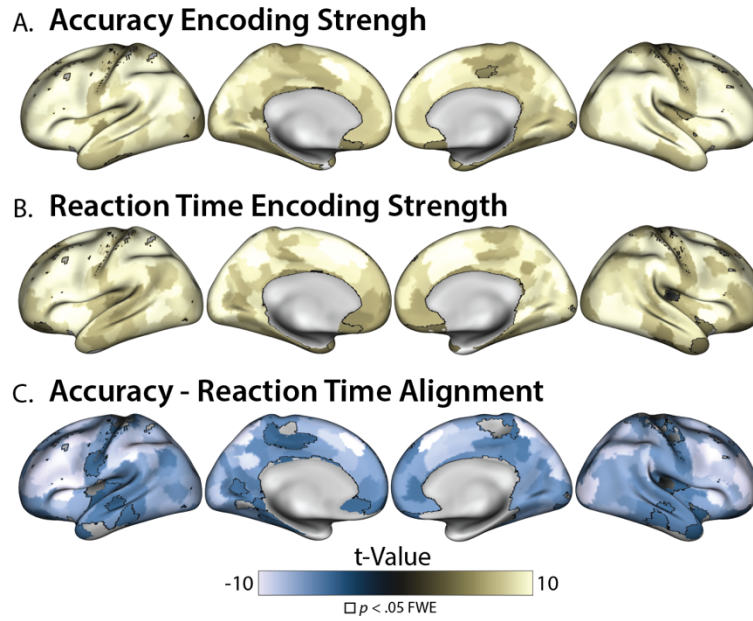
1 design, all of these encoding profiles had the same orthogonal encoding alignment (uncorrelated encoding weights),
 2 highlighting that this measure is unable to adjudicate between candidate encoding profiles. **C)** These models can be
 3 differentiated by correlating their absolute encoding weights, testing whether the sensitivity of a voxel to one feature
 4 is related to its sensitivity to the other feature, ignoring the direction of encoding. Pure selective encoding predicts a
 5 negative relationship, mixed selective encoding predicts no relationship, and sparse mixed selective encoding
 6 predicts a positive relationship. Similarity matrices averaged over 10,000 simulations. **D)** Correlating the absolute
 7 encoding weights, we found that the IPS profile was consistent with sparse mixed selective encoding.
 8



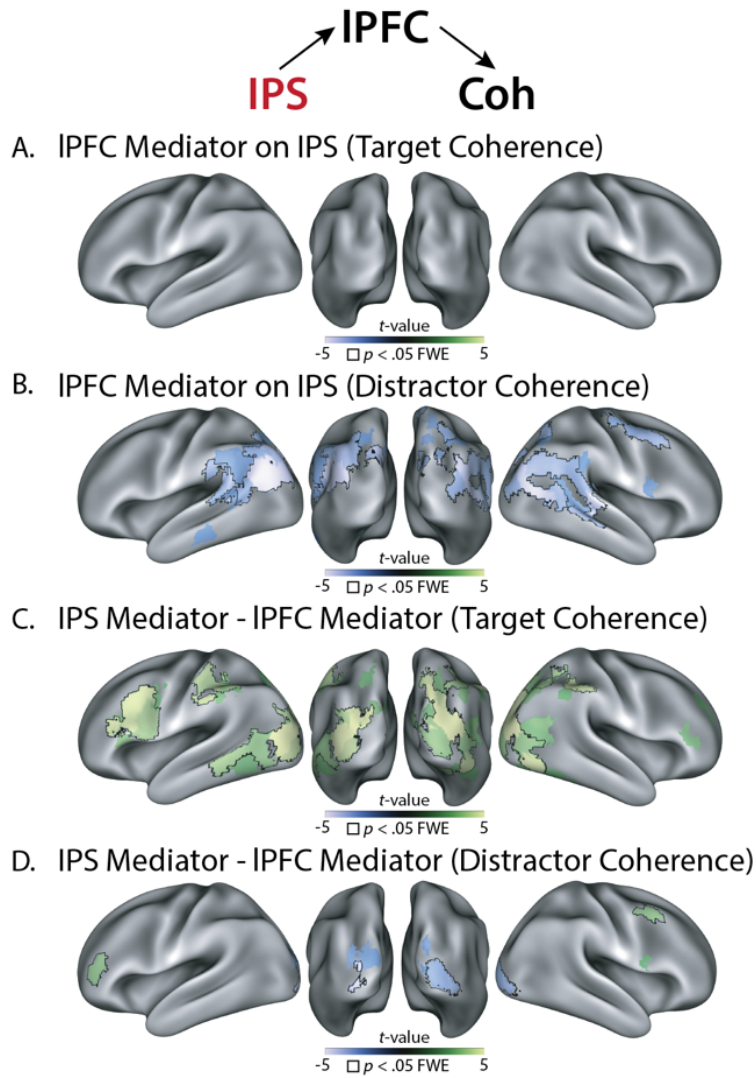
9
 10 **Figure S7. Binary evidence encoding control analysis.** Target-distractor response encoding alignment using binary
 11 evidence rather than coherence-modulated evidence. Note the similarity to Figure 4D.
 12



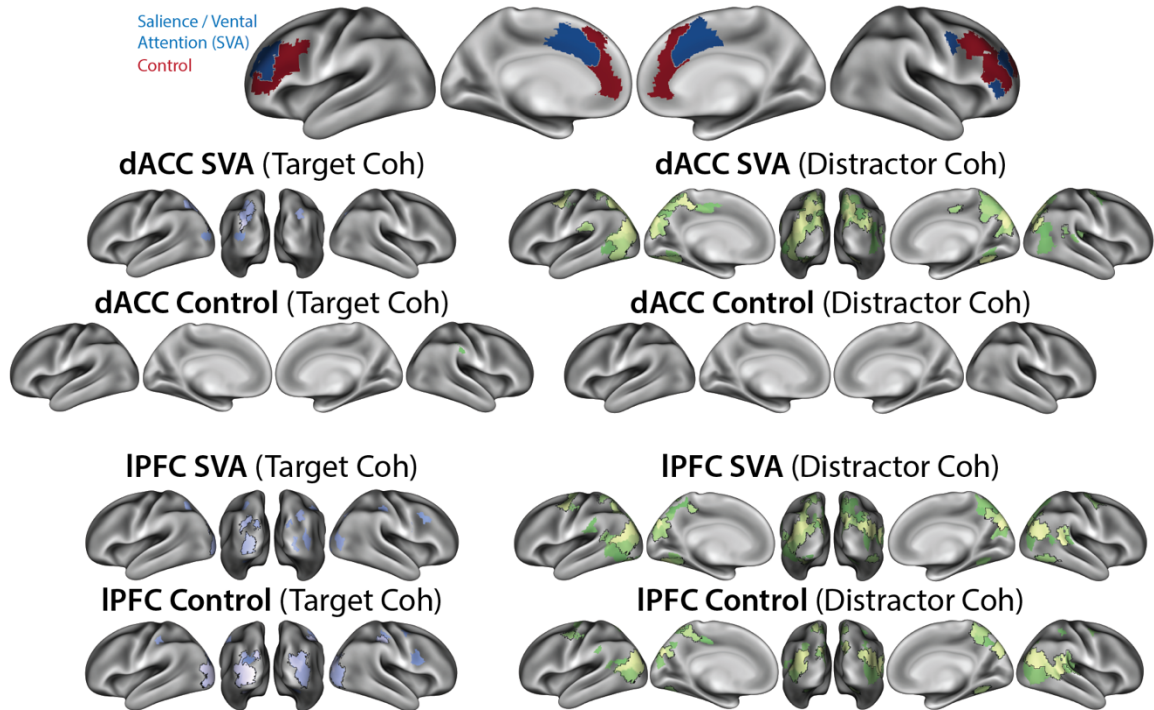
13
 14 **Figure S8. Feature encoding in frontal networks.** **A)** Similarity matrices for ‘Saliency / Ventral Attention (SVA)’
 15 and ‘Control’ networks in dACC and IPFC, correlating feature evidence (‘Evid’), feature coherence (‘Coh’), and
 16 feature congruence (‘Cong’). Encoding strength on diagonal (right-tailed p -value), encoding alignment on off-
 17 diagonal (two-tailed p -value). **B)** Classical MDS embedding of target (circle) and distractor (diamond)
 18 representations at different levels of evidence. Colors denote responses, hues denote coherence. GLMs: A: Feature
 19 MV, B: Evidence Levels, see Table 2.
 20



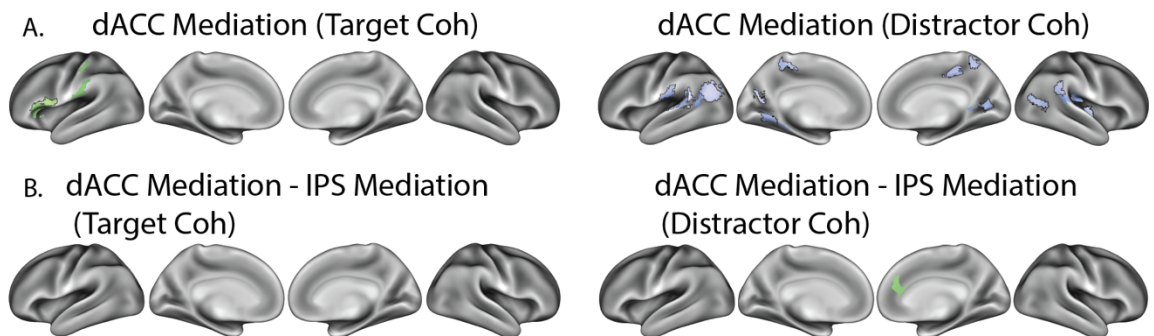
1
2 **Figure S9.** *Performance encoding.* Encoding Strength (across-run reliability) for **A)** Accuracy and **B)** Reaction
3 Time (**B**). **C)** Alignment between Accuracy and Reaction Time encoding. Outlined parcels are significant at $p < .05$
4 FWE (max-statistic randomization test). Parcels in **C** are thresholded based on the reliability in **A** and **B** (both $p <$
5 $.001$). GLM: Performance.
6



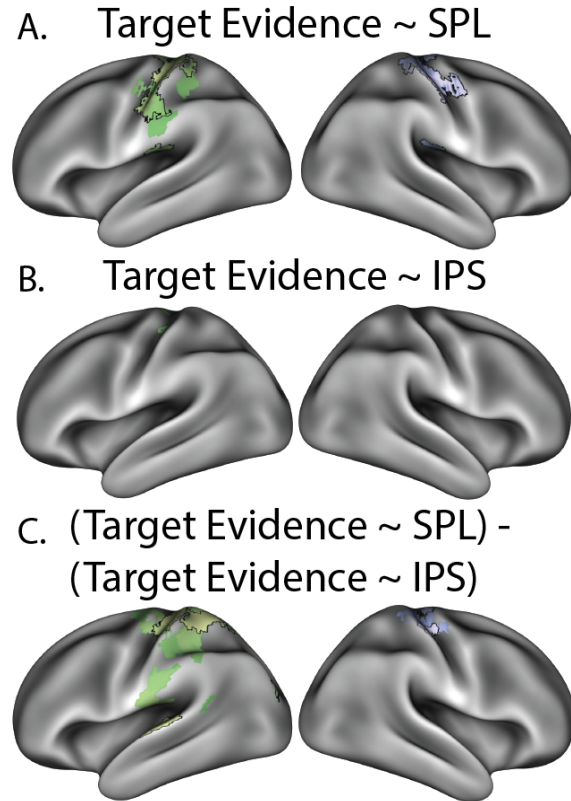
1
2 **Figure S10.** *IPFC mediation.* IPS→IPFC→Coherence mediation for target coherence (A) and distractor coherence
3 (B; compare to Figure 7c). Contrast between IPS-mediation and IPFC-mediation for target coherence (C) and
4 distractor coherence (D).
5



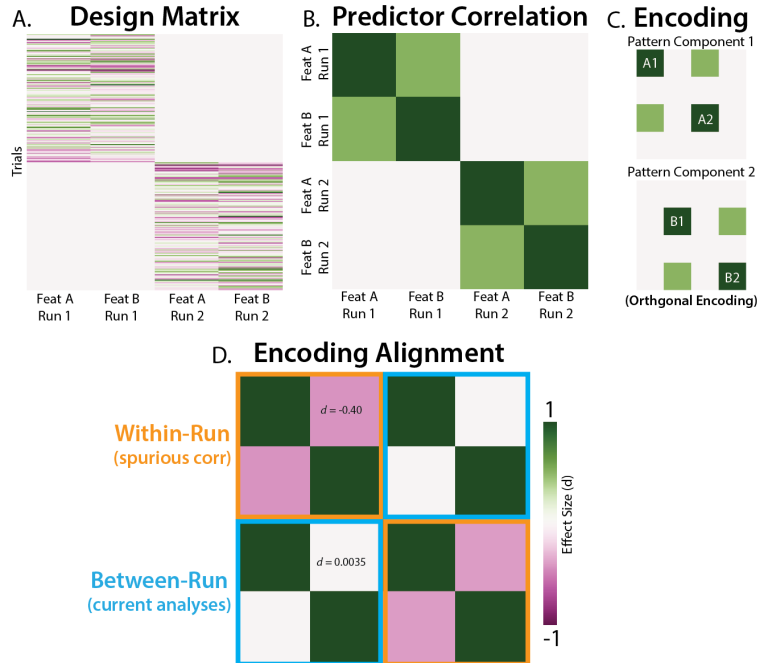
1
2 **Figure S11.** *Coherence alignment with frontal networks.* Activity in ‘Salience / Ventral Attention (SVA)’ and
3 ‘Control’ networks within dACC and IPFC (rows), aligned with target and distractor coherence (columns). Note the
4 similarity between dACC SVA parcels and IPFC parcels.
5



6
7 **Figure S12.** *IPS mediation of dACC connectivity.* **A)** IPS mediation of dACC connectivity (difference in dACC-
8 coherence alignment with and without including IPS predictors). **B)** Difference between ‘IPS mediation of dACC’
9 and ‘dACC mediation of IPS’. The lack of activation suggests that this relationship is bidirectional, or originates
10 from a common cause. dACC seed is from the ‘Salience / Ventral Attention’ network (see Supplementary Figure
11 11).
12



1
2 **Figure S13.** *SPL alignment with evidence encoding.* **A)** Alignment between SPL activity and target evidence
3 encoding. **B)** Alignment between IPS activity and target evidence encoding. **C)** Differences between SPL-evidence
4 alignment and IPS-evidence alignment, showing stronger SPL connectivity. Note that target evidence encoding is
5 signed according to the right-hand response (contralateral motor cortex should have a positive response).
6
7



1
2 **Figure S14.** *Cross-validation avoids feature correlations biasing alignment.* We used pattern component modeling
3 (Diedrichsen et al., 2018) to simulate neural data, testing whether feature correlations could spuriously create
4 encoding alignment. **A)** Our design matrix had two simulate runs of two feature timeseries. **B)** Our features were
5 correlated by design (i.e., the columns of the design matrix were correlated). **C)** Despite correlation in the design
6 matrix, these features were independently encoding in our simulated neural population (i.e., in two distinct pattern
7 components, which were each reliable across runs). **D)** Correlating our estimated encoding profiles, we found that
8 within-run alignment (orange) had a spurious negative correlation (the opposite direction of the feature correlations).
9 Critically, our analyses used between-run alignment (cyan), which avoids this biasing effect of feature correlations.
10 Intuitively, since features are not correlated across runs (i.e., they come from different trials), they do not produce
11 spurious correlations. Effect sizes are computed across 10,000 simulations.

12

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	$r_{(27)} = -0.71$ $p = 0.50 \times 10^{-4}$	$r_{(27)} = -0.43$ $p = .027$	$r_{(27)} = -0.48$ $p = .012$	$r_{(27)} = -0.59$ $p = .0014$

13 **Table S1.** *Partial correlations between coherence and performance.* Correlations between
14 individual differences in coherence-performance alignment, controlling for coherence and
15 performance encoding reliability. Since reliability determines alignment (Spearman, 1987),
16 similarity in alignment may be confounded with similarity in reliability. Overall, these results are
17 qualitatively similar to the zero-order correlation (see Figure 6), albeit with weaker correlations
18 for target coherence. These correlations are particularly robust in IPS.

19

1 References

- 2 Adam KCS, Serences JT. 2021. History modulates early sensory processing of salient distractors.
3 *J Neurosci*. doi:10.1523/JNEUROSCI.3099-20.2021
- 4 Aoi MC, Mante V, Pillow JW. 2020. Prefrontal cortex exhibits multidimensional dynamic
5 encoding during decision-making. *Nat Neurosci*. doi:10.1038/s41593-020-0696-5
- 6 Badre D, Bhandari A, Keglovits H, Kikumoto A. 2021. The dimensionality of neural
7 representations for control. *Curr Opin Behav Sci* **38**:20–28.
8 doi:10.1016/j.cobeha.2020.07.002
- 9 Badre D, D’Esposito M. 2009. Is the rostro-caudal axis of the frontal lobe hierarchical? *Nat Rev*
10 *Neurosci* **10**:659–669. doi:10.1038/nrn2667
- 11 Badre D, Nee DE. 2018. Frontal Cortex and the Hierarchical Control of Behavior. *Trends Cogn*
12 *Sci* **22**:170–188. doi:10.1016/j.tics.2017.11.005
- 13 Beldzik E, Ullsperger M. 2023. A thin line between conflict and reaction time effects on EEG
14 and fMRI brain signals. *bioRxiv*. doi:10.1101/2023.02.14.528515
- 15 Belsley DA, Kuh E, Welsch RE. 1980. Wiley Series in Probability and Statistics. *Regression*
16 *Diagnostics: Identifying Influential Data and Sources of Collinearity* 293–300.
- 17 Bernardi S, Benna MK, Rigotti M, Munuera J, Fusi S, Daniel Salzman C. 2020. The Geometry
18 of Abstraction in the Hippocampus and Prefrontal Cortex. *Cell* **0**.
19 doi:10.1016/j.cell.2020.09.031
- 20 Bisley JW, Goldberg ME. 2010. Attention, intention, and priority in the parietal lobe. *Annu Rev*
21 *Neurosci* **33**:1–21. doi:10.1146/annurev-neuro-060909-152823
- 22 Bisley JW, Mirpour K. 2019. The neural instantiation of a priority map. *Curr Opin Psychol*
23 **29**:108–112. doi:10.1016/j.copsyc.2019.01.002
- 24 Brown JW, Braver TS. 2005. Learned predictions of error likelihood in the anterior cingulate
25 cortex. *Science* **307**:1118–1121. doi:10.1126/science.1105783
- 26 Clairis N, Pessiglione M. 2022. Value, confidence, deliberation: a functional partition of the
27 medial prefrontal cortex demonstrated across rating and choice tasks. *J Neurosci* **42**:5580–
28 5592. doi:10.1523/JNEUROSCI.1795-21.2022
- 29 Cohen JD, Dunbar K, McClelland JL. 1990. On the control of automatic processes: a parallel
30 distributed processing account of the Stroop effect. *Psychol Rev* **97**:332–361.
31 doi:10.1037/0033-295x.97.3.332
- 32 Cohen MR, Maunsell JHR. 2010. A neuronal population measure of attention predicts behavioral
33 performance on individual trials. *J Neurosci* **30**:15241–15253.
34 doi:10.1523/JNEUROSCI.2171-10.2010
- 35 Corbetta M, Shulman GL. 2002. Control of goal-directed and stimulus-driven attention in the
36 brain. *Nat Rev Neurosci* **3**:201–215. doi:10.1038/nrn755

- 1 Culham JC, Brandt SA, Cavanagh P, Kanwisher NG, Dale AM, Tootell RB. 1998. Cortical fMRI
2 activation produced by attentive tracking of moving targets. *J Neurophysiol* **80**:2657–2670.
3 doi:10.1152/jn.1998.80.5.2657
- 4 Culham JC, Cavanagh P, Kanwisher NG. 2001. Attention response functions: characterizing
5 brain areas using fMRI activation during parametric variations of attentional load. *Neuron*
6 **32**:737–745. doi:10.1016/s0896-6273(01)00499-8
- 7 Cunningham JP, Yu BM. 2014. Dimensionality reduction for large-scale neural recordings. *Nat*
8 *Neurosci* **17**:1500–1509. doi:10.1038/nn.3776
- 9 Danielmeier C, Eichele T, Forstmann BU, Tittgemeyer M, Ullsperger M. 2011. Posterior medial
10 frontal cortex activity predicts post-error adaptations in task-related visual and motor areas. *J*
11 *Neurosci* **31**:1780–1789. doi:10.1523/JNEUROSCI.4299-10.2011
- 12 Danielmeier C, Ullsperger M. 2011. Post-error adjustments. *Front Psychol* **2**:233.
13 doi:10.3389/fpsyg.2011.00233
- 14 Desimone R, Duncan J. 1995. Neural mechanisms of selective visual attention. *Annu Rev*
15 *Neurosci* **18**:193–222. doi:10.1146/annurev.ne.18.030195.001205
- 16 Diedrichsen J, Kriegeskorte N. 2017. Representational models: A common framework for
17 understanding encoding, pattern-component, and representational-similarity analysis. *PLoS*
18 *Comput Biol* **13**:e1005508. doi:10.1371/journal.pcbi.1005508
- 19 Diedrichsen J, Shadmehr R. 2005. Detecting and adjusting for artifacts in fMRI time series data.
20 *Neuroimage* **27**:624–634. doi:10.1016/j.neuroimage.2005.04.039
- 21 Diedrichsen J, Yokoi A, Arbuuckle SA. 2018. Pattern component modeling: A flexible approach
22 for understanding the representational structure of brain activity patterns. *Neuroimage*
23 **180**:119–133. doi:10.1016/j.neuroimage.2017.08.051
- 24 Ebitz BR, Smith EH, Horga G, Schevon CA, Yates MJ, McKhann GM, Botvinick MM, Sheth
25 SA, Hayden BY. 2020. Human dorsal anterior cingulate neurons signal conflict by
26 amplifying task-relevant information. *bioRxiv*. doi:10.1101/2020.03.14.991745
- 27 Ebitz RB, Hayden BY. 2021. The population doctrine in cognitive neuroscience. *Neuron*.
28 doi:10.1016/j.neuron.2021.07.011
- 29 Egner T. 2008. Multiple conflict-driven control mechanisms in the human brain. *Trends Cogn*
30 *Sci* **12**:374–380. doi:10.1016/j.tics.2008.07.001
- 31 Esteban O, Markiewicz CJ, Blair RW, Moodie CA, Isik AI, Erramuzpe A, Kent JD, Goncalves
32 M, DuPre E, Snyder M, Oya H, Ghosh SS, Wright J, Durnez J, Poldrack RA, Gorgolewski
33 KJ. 2019. fMRIPrep: a robust preprocessing pipeline for functional MRI. *Nat Methods*
34 **16**:111–116. doi:10.1038/s41592-018-0235-4
- 35 Esterman M, Chiu Y-C, Tamber-Rosenau BJ, Yantis S. 2009. Decoding cognitive control in
36 human parietal cortex. *Proc Natl Acad Sci U S A* **106**:17974–17979.
37 doi:10.1073/pnas.0903593106
- 38 Etzel JA, Cole MW, Zacks JM, Kay KN, Braver TS. 2016. Reward Motivation Enhances Task
39 Coding in Frontoparietal Cortex. *Cereb Cortex* **26**:1647–1659. doi:10.1093/cercor/bhu327

- 1 Fischer AG, Nigbur R, Klein TA, Danielmeier C, Ullsperger M. 2018. Cortical beta power
2 reflects decision dynamics and uncovers multiple facets of post-error adaptation. *Nat*
3 *Commun* **9**:5038. doi:10.1038/s41467-018-07456-8
- 4 Fleming SM, van der Putten EJ, Daw ND. 2018. Neural mediators of changes of mind about
5 perceptual decisions. *Nat Neurosci* **21**:617–624. doi:10.1038/s41593-018-0104-6
- 6 Flesch T, Juechems K, Dumbalska T, Saxe A, Summerfield C. 2022. Orthogonal representations
7 for robust context-dependent task performance in brains and neural networks. *Neuron* **0**.
8 doi:10.1016/j.neuron.2022.01.005
- 9 Friedman NP, Miyake A. 2017. Unity and diversity of executive functions: Individual
10 differences as a window on cognitive structure. *Cortex* **86**:186–204.
11 doi:10.1016/j.cortex.2016.04.023
- 12 Friedman NP, Robbins TW. 2021. The role of prefrontal cortex in cognitive control and
13 executive function. *Neuropsychopharmacology* 1–18. doi:10.1038/s41386-021-01132-0
- 14 Fu Z, Beam D, Chung JM, Reed CM, Mamelak AN, Adolphs R, Rutishauser U. 2022. The
15 geometry of domain-general performance monitoring in the human medial frontal cortex.
16 *Science* **376**:eabm9922. doi:10.1126/science.abm9922
- 17 Fu Z, Wu D-AJ, Ross I, Chung JM, Mamelak AN, Adolphs R, Rutishauser U. 2019. Single-
18 Neuron Correlates of Error Monitoring and Post-Error Adjustments in Human Medial Frontal
19 Cortex. *Neuron* **101**:165-177.e5. doi:10.1016/j.neuron.2018.11.016
- 20 Gale DJ, Vos de Wael R, Benkarim O, Bernhardt B. 2021. Surfplot: Publication-ready brain
21 surface figures. doi:10.5281/zenodo.5567926
- 22 Goldman-Rakic PS. 1988. Topography of cognition: parallel distributed networks in primate
23 association cortex. *Annu Rev Neurosci* **11**:137–156.
24 doi:10.1146/annurev.ne.11.030188.001033
- 25 Gordon EM, Laumann TO, Gilmore AW, Newbold DJ, Greene DJ, Berg JJ, Ortega M, Hoyt-
26 Drazen C, Gratton C, Sun H, Hampton JM, Coalson RS, Nguyen AL, McDermott KB,
27 Shimony JS, Snyder AZ, Schlaggar BL, Petersen SE, Nelson SM, Dosenbach NUF. 2017.
28 Precision Functional Mapping of Individual Human Brains. *Neuron* **95**:791-807.e7.
29 doi:10.1016/j.neuron.2017.07.011
- 30 Gorgolewski K, Burns CD, Madison C, Clark D, Halchenko YO, Waskom ML, Ghosh SS. 2011.
31 Nipype: a flexible, lightweight and extensible neuroimaging data processing framework in
32 python. *Front Neuroinform* **5**:13. doi:10.3389/fninf.2011.00013
- 33 Gottlieb J, Cohanpour M, Li Y, Singletary N, Zabeh E. 2020. Curiosity, information demand and
34 attentional priority. *Current Opinion in Behavioral Sciences* **35**:83–91.
35 doi:10.1016/j.cobeha.2020.07.016
- 36 Gratton C, Laumann TO, Gordon EM, Adeyemo B, Petersen SE. 2016. Evidence for Two
37 Independent Factors that Modify Brain Networks to Meet Task Goals. *Cell Rep* **17**:1276–
38 1288. doi:10.1016/j.celrep.2016.10.002
- 39 Greenberg AS, Esterman M, Wilson D, Serences JT, Yantis S. 2010. Control of spatial and
40 feature-based attention in frontoparietal cortex. *J Neurosci* **30**:14330–14339.
41 doi:10.1523/JNEUROSCI.4248-09.2010

- 1 Grinband J, Savitskaya J, Wager TD, Teichert T, Ferrera VP, Hirsch J. 2011. The dorsal medial
2 frontal cortex is sensitive to time on task, not response conflict or error likelihood.
3 *Neuroimage* **57**:303–311. doi:10.1016/j.neuroimage.2010.12.027
- 4 Hall-McMaster S, Muhle-Karbe PS, Myers NE, Stokes MG. 2019. Reward Boosts Neural
5 Coding of Task Rules to Optimize Cognitive Flexibility. *J Neurosci* **39**:8549–8561.
6 doi:10.1523/JNEUROSCI.0631-19.2019
- 7 Holroyd CB, McClure SM. 2015. Hierarchical control over effortful behavior by rodent medial
8 frontal cortex: A computational model. *Psychol Rev* **122**:54–83. doi:10.1037/a0038339
- 9 Holroyd CB, Yeung N. 2011. An Integrative Theory of Anterior Cingulate Cortex Function:
10 Option Selection in Hierarchical Reinforcement Learning. *Neural Basis of Motivational and*
11 *Cognitive Control*. doi:10.7551/mitpress/9780262016438.003.0018
- 12 Howe PD, Horowitz TS, Morocz IA, Wolfe J, Livingstone MS. 2009. Using fMRI to distinguish
13 components of the multiple object tracking task. *J Vis* **9**:10.1-11. doi:10.1167/9.4.10
- 14 Hunt LT, Kolling N, Soltani A, Woolrich MW, Rushworth MFS, Behrens TEJ. 2012.
15 Mechanisms underlying cortical activity during value-guided choice. *Nat Neurosci* **15**:470–6,
16 S1-3. doi:10.1038/nn.3017
- 17 Jackson J, Rich AN, Williams MA, Woolgar A. 2017. Feature-selective Attention in
18 Frontoparietal Cortex: Multivoxel Codes Adjust to Prioritize Task-relevant Information. *J*
19 *Cogn Neurosci* **29**:310–321. doi:10.1162/jocn_a_01039
- 20 Jackson JB, Feredoes E, Rich AN, Lindner M, Woolgar A. 2021. Concurrent neuroimaging and
21 neurostimulation reveals a causal role for dlPFC in coding of task-relevant information.
22 *Commun Biol* **4**:588. doi:10.1038/s42003-021-02109-x
- 23 Jones MS, Zhu Z, Bajracharya A, Luor A, Peelle JE. 2021. A multi-dataset evaluation of frame
24 censoring for task-based fMRI. *bioRxiv*. doi:10.1101/2021.10.12.464075
- 25 Jovicich J, Peters RJ, Koch C, Braun J, Chang L, Ernst T. 2001. Brain areas specific for
26 attentional load in a motion-tracking task. *J Cogn Neurosci* **13**:1048–1058.
- 27 Kalman RE. 1960. On the general theory of control systems. *IFAC Proceedings Volumes* **1**:491–
28 502. doi:10.1016/S1474-6670(17)70094-8
- 29 Kastner S, Ungerleider LG. 2000. Mechanisms of visual attention in the human cortex. *Annu Rev*
30 *Neurosci* **23**:315–341. doi:10.1146/annurev.neuro.23.1.315
- 31 Kay KN, Yeatman JD. 2017. Bottom-up and top-down computations in word- and face-selective
32 cortex. *Elife* **6**. doi:10.7554/eLife.22341
- 33 Kayser AS, Buchsbaum BR, Erickson DT, D’Esposito M. 2010a. The functional anatomy of a
34 perceptual decision in the human brain. *J Neurophysiol* **103**:1179–1194.
35 doi:10.1152/jn.00364.2009
- 36 Kayser AS, Erickson DT, Buchsbaum BR, D’Esposito M. 2010b. Neural representations of
37 relevant and irrelevant features in perceptual decision making. *J Neurosci* **30**:15778–15789.
38 doi:10.1523/JNEUROSCI.3163-10.2010

- 1 Kerns JG, Cohen JD, MacDonald AW 3rd, Cho RY, Stenger VA, Carter CS. 2004. Anterior
2 cingulate conflict monitoring and adjustments in control. *Science* **303**:1023–1026.
3 doi:10.1126/science.1089910
- 4 Kimmel DL, Elsayed GF, Cunningham JP, Newsome WT. 2020. Value and choice as separable
5 and stable representations in orbitofrontal cortex. *Nat Commun* **11**:3466.
6 doi:10.1038/s41467-020-17058-y
- 7 Koechlin E, Summerfield C. 2007. An information theoretical approach to prefrontal executive
8 function. *Trends Cogn Sci* **11**:229–235. doi:10.1016/j.tics.2007.04.005
- 9 Kong R, Li J, Orban C, Sabuncu MR, Liu H, Schaefer A, Sun N, Zuo X-N, Holmes AJ, Eickhoff
10 SB, Yeo BTT. 2019. Spatial Topography of Individual-Specific Cortical Networks Predicts
11 Human Cognition, Personality, and Emotion. *Cereb Cortex* **29**:2533–2551.
12 doi:10.1093/cercor/bhy123
- 13 Kong R, Yang Q, Gordon E, Xue A, Yan X, Orban C, Zuo X-N, Spreng N, Ge T, Holmes A,
14 Eickhoff S, Yeo BTT. 2021. Individual-Specific Areal-Level Parcellations Improve
15 Functional Connectivity Prediction of Behavior. *Cereb Cortex* **31**:4477–4500.
16 doi:10.1093/cercor/bhab101
- 17 Kragel PA, Kano M, Van Oudenhove L, Ly HG, Dupont P, Rubio A, Delon-Martin C, Bonaz
18 BL, Manuck SB, Gianaros PJ, Ceko M, Reynolds Losin EA, Woo C-W, Nichols TE, Wager
19 TD. 2018. Generalizable representations of pain, cognitive control, and negative emotion in
20 medial frontal cortex. *Nat Neurosci* **21**:283–289. doi:10.1038/s41593-017-0051-7
- 21 Kriegeskorte N, Diedrichsen J. 2019. Peeling the Onion of Brain Representations. *Annu Rev*
22 *Neurosci* **42**:407–432. doi:10.1146/annurev-neuro-080317-061906
- 23 Kriegeskorte N, Goebel R, Bandettini P. 2006. Information-based functional brain mapping.
24 *Proc Natl Acad Sci U S A* **103**:3863–3868. doi:10.1073/pnas.0600244103
- 25 Lauritzen TZ, D’Esposito M, Heeger DJ, Silver MA. 2009. Top–down flow of visual spatial
26 attention signals from parietal to occipital cortex. *J Vis* **9**:18–18. doi:10.1167/9.13.18
- 27 Leng X, Yee D, Ritz H, Shenhav A. 2021. Dissociable influences of reward and punishment on
28 adaptive cognitive control. *PLoS Comput Biol* **17**:e1009737.
29 doi:10.1371/journal.pcbi.1009737
- 30 Libby A, Buschman TJ. 2021. Rotational dynamics reduce interference between sensory and
31 memory representations. *Nat Neurosci* 1–12. doi:10.1038/s41593-021-00821-9
- 32 MacDonald AW 3rd, Cohen JD, Stenger VA, Carter CS. 2000. Dissociating the role of the
33 dorsolateral prefrontal and anterior cingulate cortex in cognitive control. *Science* **288**:1835–
34 1838. doi:10.1126/science.288.5472.1835
- 35 MacKinnon DP, Fairchild AJ, Fritz MS. 2007. Mediation analysis. *Annu Rev Psychol* **58**:593–
36 614. doi:10.1146/annurev.psych.58.110405.085542
- 37 Mante V, Sussillo D, Shenoy KV, Newsome WT. 2013. Context-dependent computation by
38 recurrent dynamics in prefrontal cortex. *Nature* **503**:78–84. doi:10.1038/nature12742
- 39 Menon V, D’Esposito M. 2021. The role of PFC networks in cognitive control and executive
40 function. *Neuropsychopharmacology* 1–14. doi:10.1038/s41386-021-01152-w

- 1 Miller EK, Cohen JD. 2001. An integrative theory of prefrontal cortex function. *Annu Rev*
2 *Neurosci* **24**:167–202. doi:10.1146/annurev.neuro.24.1.167
- 3 Minxha J, Adolphs R, Fusi S, Mamelak AN, Rutishauser U. 2020. Flexible recruitment of
4 memory-based choice representations by the human medial frontal cortex. *Science* **368**.
5 doi:10.1126/science.aba3313
- 6 Molenberghs P, Mesulam MM, Peeters R, Vandenberghe RRC. 2007. Remapping attentional
7 priorities: differential contribution of superior parietal lobule and intraparietal sulcus. *Cereb*
8 *Cortex* **17**:2703–2712. doi:10.1093/cercor/bhl179
- 9 Mumford JA, Bissett PG, Jones HM, Shim S, Rios JAH, Poldrack RA. 2023. The response time
10 paradox in functional magnetic resonance imaging analyses. *bioRxiv*.
11 doi:10.1101/2023.02.15.528677
- 12 Musslick S, Shenhav A, Botvinick M, Cohen J. 2015. A Computational Model of Control
13 Allocation based on the Expected Value of Control2nd Multidisciplinary Conference on
14 Reinforcement Learning and Decision Making. Presented at the Multidisciplinary
15 Conference on Reinforcement Learning and Decision Making.
- 16 Nee DE, Wager TD, Jonides J. 2007. Interference resolution: insights from a meta-analysis of
17 neuroimaging tasks. *Cogn Affect Behav Neurosci* **7**:1–17. doi:10.3758/cabn.7.1.1
- 18 Nichols TE, Holmes AP. 2002. Nonparametric permutation tests for functional neuroimaging: a
19 primer with examples. *Hum Brain Mapp* **15**:1–25. doi:10.1002/hbm.1058
- 20 Nili H, Wingfield C, Walther A, Su L, Marslen-Wilson W, Kriegeskorte N. 2014. A toolbox for
21 representational similarity analysis. *PLoS Comput Biol* **10**:e1003553.
22 doi:10.1371/journal.pcbi.1003553
- 23 Pagan M, Tang VD, Aoi MC, Pillow JW, Mante V, Sussillo D, Brody CD. 2022. A new
24 theoretical framework jointly explains behavioral and neural variability across subjects
25 performing flexible decision-making. *bioRxiv*. doi:10.1101/2022.11.28.518207
- 26 Panichello MF, Buschman TJ. 2021. Shared mechanisms underlie the control of working
27 memory and attention. *Nature* 1–5. doi:10.1038/s41586-021-03390-w
- 28 Parro C, Dixon ML, Christoff K. 2018. The neural basis of motivational influences on cognitive
29 control. *Hum Brain Mapp* **39**:5097–5111. doi:10.1002/hbm.24348
- 30 Peck CJ, Jangraw DC, Suzuki M, Efem R, Gottlieb J. 2009. Reward modulates attention
31 independently of action value in posterior parietal cortex. *J Neurosci* **29**:11182–11191.
32 doi:10.1523/JNEUROSCI.1929-09.2009
- 33 Petrides M, Pandya DN. 2006. Efferent association pathways originating in the caudal prefrontal
34 cortex in the macaque monkey. *J Comp Neurol* **498**:227–251. doi:10.1002/cne.21048
- 35 Pylyshyn ZW, Storm RW. 1988. Tracking multiple independent targets: evidence for a parallel
36 tracking mechanism. *Spat Vis* **3**:179–197. doi:10.1163/156856888x00122
- 37 Reid AT, Headley DB, Mill RD, Sanchez-Romero R, Uddin LQ, Marinazzo D, Lurie DJ,
38 Valdés-Sosa PA, Hanson SJ, Biswal BB, Calhoun V, Poldrack RA, Cole MW. 2019.
39 Advancing functional connectivity research from association to causation. *Nat Neurosci*
40 **22**:1751–1760. doi:10.1038/s41593-019-0510-4

- 1 Rigotti M, Barak O, Warden MR, Wang X-J, Daw ND, Miller EK, Fusi S. 2013. The importance
2 of mixed selectivity in complex cognitive tasks. *Nature* **497**:585–590.
3 doi:10.1038/nature12160
- 4 Ritz H, Leng X, Shenhav A. 2022a. Cognitive Control as a Multivariate Optimization Problem. *J*
5 *Cogn Neurosci* **34**:569–591. doi:10.1162/jocn_a_01822
- 6 Ritz H, Shenhav A. 2021. Humans reconfigure target and distractor processing to address distinct
7 task demands. *bioRxiv* 2021.09.08.459546. doi:10.1101/2021.09.08.459546
- 8 Ritz H, Wild CJ, Johnsrude IS. 2022b. Parametric Cognitive Load Reveals Hidden Costs in the
9 Neural Processing of Perfectly Intelligible Degraded Speech. *J Neurosci* **42**:4619–4628.
10 doi:10.1523/JNEUROSCI.1777-21.2022
- 11 Rouder JN, Morey RD, Speckman PL, Province JM. 2012. Default Bayes factors for ANOVA
12 designs. *J Math Psychol* **56**:356–374. doi:10.1016/j.jmp.2012.08.001
- 13 Rushworth MF, Behrens TE. 2008. Choice, uncertainty and value in prefrontal and cingulate
14 cortex. *Nat Neurosci* **11**:389–397. doi:10.1038/nn2066
- 15 Rust NC, Cohen MR. 2022. Priority coding in the visual system. *Nat Rev Neurosci* 1–13.
16 doi:10.1038/s41583-022-00582-9
- 17 Saalmann YB, Pigarev IN, Vidyasagar TR. 2007. Neural mechanisms of visual attention: how
18 top-down feedback highlights relevant locations. *Science* **316**:1612–1615.
19 doi:10.1126/science.1139140
- 20 Schaefer A, Kong R, Gordon EM, Laumann TO, Zuo X-N, Holmes AJ, Eickhoff SB, Yeo BTT.
21 2018. Local-Global Parcellation of the Human Cerebral Cortex from Intrinsic Functional
22 Connectivity MRI. *Cereb Cortex* **28**:3095–3114. doi:10.1093/cercor/bhx179
- 23 Semedo JD, Zandvakili A, Machens CK, Yu BM, Kohn A. 2019. Cortical Areas Interact through
24 a Communication Subspace. *Neuron* **102**:249-259.e4. doi:10.1016/j.neuron.2019.01.026
- 25 Serences JT, Schwarzbach J, Courtney SM, Golay X, Yantis S. 2004. Control of object-based
26 attention in human cortex. *Cereb Cortex* **14**:1346–1357. doi:10.1093/cercor/bhh095
- 27 Serences JT, Yantis S. 2007. Spatially selective representations of voluntary and stimulus-driven
28 attentional priority in human occipital, parietal, and frontal cortex. *Cereb Cortex* **17**:284–293.
29 doi:10.1093/cercor/bhj146
- 30 Shenhav A, Botvinick MM, Cohen JD. 2013. The expected value of control: an integrative
31 theory of anterior cingulate cortex function. *Neuron* **79**:217–240.
32 doi:10.1016/j.neuron.2013.07.007
- 33 Shenhav A, Karmarkar UR. 2019. Dissociable components of the reward circuit are involved in
34 appraisal versus choice. *Sci Rep* **9**:1958. doi:10.1038/s41598-019-38927-7
- 35 Shenhav A, Straccia MA, Botvinick MM, Cohen JD. 2016. Dorsal anterior cingulate and
36 ventromedial prefrontal cortex have inverse roles in both foraging and economic choice.
37 *Cogn Affect Behav Neurosci*. doi:10.3758/s13415-016-0458-8
- 38 Shenhav A, Straccia MA, Musslick S, Cohen JD, Botvinick MM. 2018. Dissociable neural
39 mechanisms track evidence accumulation for selection of attention versus action. *Nat*
40 *Commun* **9**:2485. doi:10.1038/s41467-018-04841-1

- 1 Smith EH, Horga G, Yates MJ, Mikell CB, Banks GP, Pathak YJ, Schevon CA, McKhann GM,
2 Hayden BY, Botvinick MM, Sheth SA. 2019. Widespread temporal coding of cognitive
3 control in the human prefrontal cortex. *Nat Neurosci* **66**:83. doi:10.1038/s41593-019-0494-0
- 4 Smith SM, Nichols TE. 2009. Threshold-free cluster enhancement: addressing problems of
5 smoothing, threshold dependence and localisation in cluster inference. *Neuroimage* **44**:83–
6 98. doi:10.1016/j.neuroimage.2008.03.061
- 7 Soutschek A, Stelzel C, Paschke L, Walter H, Schubert T. 2015. Dissociable effects of
8 motivation and expectancy on conflict processing: an fMRI study. *J Cogn Neurosci* **27**:409–
9 423. doi:10.1162/jocn_a_00712
- 10 Spearman C. 1987. The Proof and Measurement of Association between Two Things. *Am J*
11 *Psychol* **100**:441–471. doi:10.2307/1422689
- 12 Srinath R, Ruff DA, Cohen MR. 2021. Attention improves information flow between neuronal
13 populations without changing the communication subspace. *bioRxiv*.
14 doi:10.1101/2021.03.31.437940
- 15 Stringer C, Pachitariu M, Steinmetz N, Reddy CB, Carandini M, Harris KD. 2019. Spontaneous
16 behaviors drive multidimensional, brainwide activity. *Science* **364**:255.
17 doi:10.1126/science.aav7893
- 18 Suzuki M, Gottlieb J. 2013. Distinct neural mechanisms of distractor suppression in the frontal
19 and parietal lobe. *Nat Neurosci* **16**:98–104. doi:10.1038/nn.3282
- 20 Takagi Y, Hunt LT, Woolrich MW, Behrens TE, Klein-Flügge MC. 2021. Adapting non-
21 invasive human recordings along multiple task-axes shows unfolding of spontaneous and
22 over-trained choice. *Elife* **10**. doi:10.7554/eLife.60988
- 23 Taren AA, Venkatraman V, Huettel SA. 2011. A parallel functional topography between medial
24 and lateral prefrontal cortex: evidence and implications for cognitive control. *J Neurosci*
25 **31**:5026–5031. doi:10.1523/JNEUROSCI.5762-10.2011
- 26 Thornton MA, Mitchell JP. 2017. Consistent Neural Activity Patterns Represent Personally
27 Familiar People. *J Cogn Neurosci* **29**:1583–1594. doi:10.1162/jocn_a_01151
- 28 Vassena E, Deraeve J, Alexander WH. 2017. Predicting Motivation: Computational Models of
29 PFC Can Explain Neural Coding of Motivation and Effort-based Decision-making in Health
30 and Disease. *J Cogn Neurosci* **29**:1633–1645. doi:10.1162/jocn_a_01160
- 31 Venkatraman V, Rosati AG, Taren AA, Huettel SA. 2009. Resolving response, decision, and
32 strategic control: evidence for a functional topography in dorsomedial prefrontal cortex. *J*
33 *Neurosci* **29**:13158–13164. doi:10.1523/JNEUROSCI.2708-09.2009
- 34 Vermeylen L, Wisniewski D, González-García C, Hoofs V, Notebaert W, Braem S. 2020. Shared
35 Neural Representations of Cognitive Conflict and Negative Affect in the Medial Frontal
36 Cortex. *J Neurosci* **40**:8715–8725. doi:10.1523/JNEUROSCI.1744-20.2020
- 37 Vos de Wael R, Benkarim O, Paquola C, Lariviere S, Royer J, Tavakol S, Xu T, Hong S-J,
38 Langs G, Valk S, Misic B, Milham M, Margulies D, Smallwood J, Bernhardt BC. 2020.
39 BrainSpace: a toolbox for the analysis of macroscale gradients in neuroimaging and
40 connectomics datasets. *Commun Biol* **3**:103. doi:10.1038/s42003-020-0794-7

- 1 Vul E, Alvarez G, Tenenbaum J, Black M. 2009. Explaining human multiple object tracking as
2 resource-constrained approximate inference in a dynamic probabilistic model In: Bengio Y,
3 Schuurmans D, Lafferty J, Williams C, Culotta A, editors. *Advances in Neural Information*
4 *Processing Systems*. Curran Associates, Inc.
- 5 Walther A, Nili H, Ejaz N, Alink A, Kriegeskorte N, Diedrichsen J. 2016. Reliability of
6 dissimilarity measures for multi-voxel pattern analysis. *Neuroimage* **137**:188–200.
7 doi:10.1016/j.neuroimage.2015.12.012
- 8 Weichart ER, Turner BM, Sederberg PB. 2020. A model of dynamic, within-trial conflict
9 resolution for decision making. *Psychol Rev*. doi:10.1037/rev0000191
- 10 Winkler AM, Ridgway GR, Webster MA, Smith SM, Nichols TE. 2014. Permutation inference
11 for the general linear model. *Neuroimage* **92**:381–397.
12 doi:10.1016/j.neuroimage.2014.01.060
- 13 Wisniewski D, Reverberi C, Momennejad I, Kahnt T, Haynes J-D. 2015. The Role of the Parietal
14 Cortex in the Representation of Task–Reward Associations. *J Neurosci* **35**:12355–12365.
15 doi:10.1523/JNEUROSCI.4882-14.2015
- 16 Woolgar A, Afshar S, Williams MA, Rich AN. 2015a. Flexible Coding of Task Rules in
17 Frontoparietal Cortex: An Adaptive System for Flexible Cognitive Control. *J Cogn Neurosci*
18 **27**:1895–1911. doi:10.1162/jocn_a_00827
- 19 Woolgar A, Hampshire A, Thompson R, Duncan J. 2011a. Adaptive coding of task-relevant
20 information in human frontoparietal cortex. *J Neurosci* **31**:14592–14599.
21 doi:10.1523/JNEUROSCI.2616-11.2011
- 22 Woolgar A, Thompson R, Bor D, Duncan J. 2011b. Multi-voxel coding of stimuli, rules, and
23 responses in human frontoparietal cortex. *Neuroimage* **56**:744–752.
24 doi:10.1016/j.neuroimage.2010.04.035
- 25 Woolgar A, Williams MA, Rich AN. 2015b. Attention enhances multi-voxel representation of
26 novel objects in frontal, parietal and visual cortices. *Neuroimage* **109**:429–437.
27 doi:10.1016/j.neuroimage.2014.12.083
- 28 Yantis S, Schwarzbach J, Serences JT, Carlson RL, Steinmetz MA, Pekar JJ, Courtney SM.
29 2002. Transient neural activity in human parietal cortex during spatial attention shifts. *Nat*
30 *Neurosci* **5**:995–1002. doi:10.1038/nn921
- 31 Yantis S, Serences JT. 2003. Cortical mechanisms of space-based and object-based attentional
32 control. *Curr Opin Neurobiol* **13**:187–193. doi:10.1016/s0959-4388(03)00033-3
- 33 Yarkoni T, Barch DM, Gray JR, Conturo TE, Braver TS. 2009. BOLD correlates of trial-by-trial
34 reaction time variability in gray and white matter: a multi-study fMRI analysis. *PLoS One*
35 **4**:e4257. doi:10.1371/journal.pone.0004257
- 36 Yarkoni T, Poldrack RA, Nichols TE, Van Essen DC, Wager TD. 2011. Large-scale automated
37 synthesis of human functional neuroimaging data. *Nat Methods* **8**:665–670.
38 doi:10.1038/nmeth.1635
- 39 Yeo BTT, Krienen FM, Sepulcre J, Sabuncu MR, Lashkari D, Hollinshead M, Roffman JL,
40 Smoller JW, Zöllei L, Polimeni JR, Fischl B, Liu H, Buckner RL. 2011. The organization of

1 the human cerebral cortex estimated by intrinsic functional connectivity. *J Neurophysiol*
2 **106**:1125–1165. doi:10.1152/jn.00338.2011

3 Zarr N, Brown JW. 2016. Hierarchical error representation in medial prefrontal cortex.
4 *Neuroimage* **124**:238–247. doi:10.1016/j.neuroimage.2015.08.063

5