

A Representational Similarity Analysis of Cognitive Control during Color-Word Stroop

Michael C. Freund,¹ Julie M. Bugg,¹ and Todd S. Braver^{1,2,3}

¹Department of Psychological & Brain Sciences, Washington University in St. Louis, St. Louis, Missouri 63130, ²Department of Radiology, Washington University in St. Louis School of Medicine, St. Louis, Missouri 63110, and ³Department of Neuroscience, Washington University in St. Louis School of Medicine, St. Louis, Missouri 63110

Progress in understanding the neural bases of cognitive control has been supported by the paradigmatic color-word Stroop task, in which a target response (color name) must be selected over a more automatic, yet potentially incongruent, distractor response (word). For this paradigm, models have postulated complementary coding schemes: dorsomedial frontal cortex (DMFC) is proposed to evaluate the demand for control via incongruity-related coding, whereas dorsolateral PFC (DLPFC) is proposed to implement control via goal and target-related coding. Yet, mapping these theorized schemes to measured neural activity within this task has been challenging. Here, we tested for these coding schemes relatively directly, by decomposing an event-related color-word Stroop task via representational similarity analysis. Three neural coding models were fit to the similarity structure of multivoxel patterns of human fMRI activity, acquired from 65 healthy, young-adult males and females. Incongruity coding was predominant in DMFC, whereas both target and incongruity coding were present with indistinguishable strength in DLPFC. In contrast, distractor information was strongly encoded within early visual cortex. Further, these coding schemes were differentially related to behavior: individuals with stronger DLPFC (and lateral posterior parietal cortex) target coding, but weaker DMFC incongruity coding, exhibited less behavioral Stroop interference. These results highlight the utility of the representational similarity analysis framework for investigating neural mechanisms of cognitive control and point to several promising directions to extend the Stroop paradigm.

Key words: dorsal ACC; dorsolateral PFC; executive function; fMRI; multivariate pattern analysis; response conflict

Significance Statement

How the human brain enables cognitive control — the ability to override behavioral habits to pursue internal goals — has been a major focus of neuroscience research. This ability has been frequently investigated by using the Stroop color-word naming task. With the Stroop as a test-bed, many theories have proposed specific neuroanatomical dissociations, in which medial and lateral frontal brain regions underlie cognitive control by encoding distinct types of information. Yet providing a direct confirmation of these claims has been challenging. Here, we demonstrate that representational similarity analysis, which estimates and models the similarity structure of brain activity patterns, can successfully establish the hypothesized functional dissociations within the Stroop task. Representational similarity analysis may provide a useful approach for investigating cognitive control mechanisms.

Received Nov. 18, 2020; revised May 23, 2021; accepted June 10, 2021.

Author contributions: M.C.F., J.M.B., and T.S.B. designed research; M.C.F. analyzed data; M.C.F. wrote the first draft of the paper; M.C.F. and T.S.B. wrote the paper; J.M.B. and T.S.B. edited the paper.

This work was supported by National Institutes of Health Grant R37 MH066078 to T.S.B. Mensh and Kording (2017) was a useful resource for organizing an initial draft of this manuscript. Computations were performed in part using the facilities of the Washington University Center for High Performance Computing, which were supported in part by National Institutes of Health Grants 1510RR022984-01A1 and 15100D018091-01. Surface images were prepared with *Connectome Workbench* (Marcus et al., 2011). The original copy of this manuscript was drafted with *papaja* (Aust and Barth, 2020) and *knitr* (Xie, 2015). We thank all former and current team members of the Dual Mechanisms of Cognitive Control Project for their efforts; Jo Etzel for general methodological wisdom; Atsushi Kikimoto for useful thoughts on our RSA models; the Cognitive Control and Psychopathology Laboratory for support and suggestions; and other Washington University in St. Louis colleagues in the J.M.B., Kool, and Zachs laboratories.

The authors declare no competing financial interests.

Correspondence should be addressed to Michael C. Freund at m.freund@wustl.edu.

<https://doi.org/10.1523/JNEUROSCI.2956-20.2021>

Copyright © 2021 the authors

Introduction

Goals, held in mind, can be used to overcome behavioral habits. Understanding how the human brain enables such cognitive control has been a fundamental interest of both basic and translational cognitive neuroscience. Toward this end, the use of response conflict tasks has been instrumental (e.g., Botvinick et al., 2001; Ridderinkhof et al., 2004). These tasks involve trials in which a less-automatic, but goal-relevant course of action, the target response, must be selected in the face of a habitual, but goal-irrelevant alternative, the distractor. The paradigmatic example is the color-word Stroop task (Stroop, 1935; Posner and Snyder, 1975; MacLeod, 1991): on each trial, the hue of a word must be named, despite the word expressing a potentially conflicting, that is, incongruent, color (see Fig. 1C). A major goal in

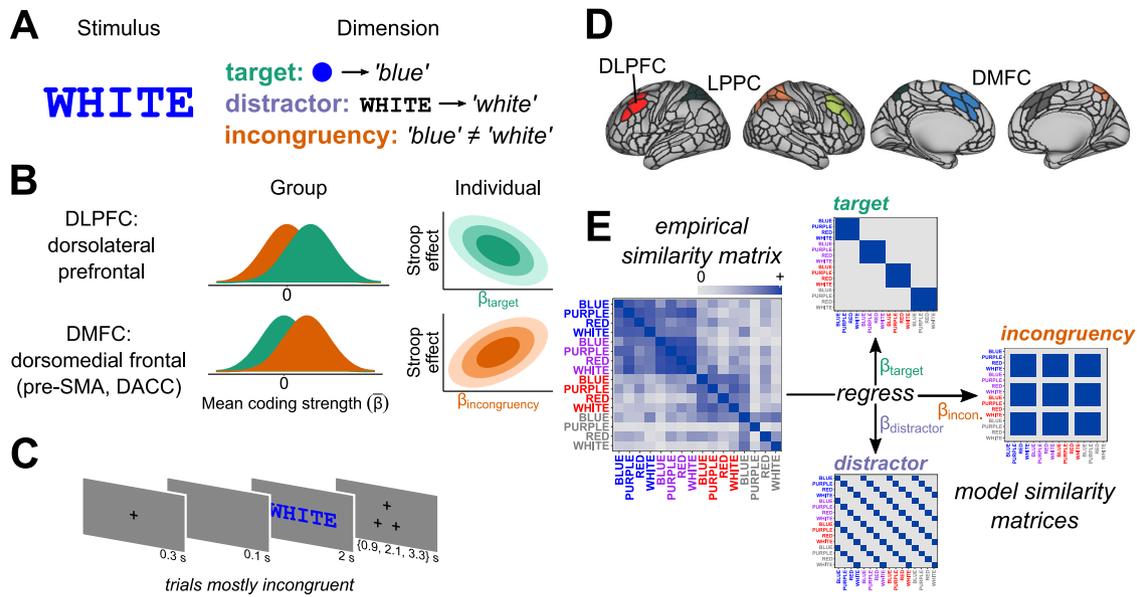


Figure 1. Schematic of framework and hypotheses. **A**, Conceptual framework. Cognitive theories and computational models of control have decomposed the classic color-word Stroop task into three task dimensions: the target (the goal-relevant mapping of the stimulus hue to the target response, i.e., color naming), the distractor (the prepotent but goal-irrelevant mapping of the stimulus word to the non-target or distractor response, i.e., word reading), and incongruency (whether the target and distractor responses match or mismatch). **B**, Hypotheses. Neuroscientific frameworks of cognitive control propose that representation of task dimensions is anatomically dissociated across medial and lateral frontoparietal cortices. DMFC, including dorsal anterior cingulate (dACC) and pre-supplementary motor area (pre-SMA), is proposed to “evaluate” demands for control, using information correlated with the incongruency dimension (**A**, bottom row), to signal when (and how) the current attentional or action selection policies are suboptimal (Ridderinkhof et al., 2004; Shenhav et al., 2013). Conversely, dorsolateral PFC (DLPFC), in concert with LPPC (LPPC), is proposed to guide, or “implement,” goal-driven attentional selection and mapping of hue–target-response processes, by way of representing information related to the goal-dependent target dimension (**A**, top row; Miller and Cohen, 2001; Buschman and Miller, 2007). Double dissociations are therefore predicted at multiple levels of analysis. At the group level, incongruency coding (orange univariate distributions) should predominate in DMFC, whereas target coding (green univariate distributions) should predominate in DLPFC. At the individual level, if the strength of target-related coding in DLPFC reflects the robustness of goal-driven selection, then subjects with stronger DLPFC target coding should resolve Stroop interference more efficiently (green bivariate distribution; Kane and Engle, 2002; Braver, 2012). Conversely, if the strength of incongruency-related coding in DMFC indicates a maladaptive selection policy, then subjects with stronger incongruency-related coding should resolve Stroop interference less efficiently (orange bivariate distribution; MacDonald et al., 2000; Braver, 2012). The β notation in axis titles corresponds to that used in **E** (and throughout this manuscript); β indicates mean over subjects. **C–E**, Analytic framework. Participants performed a color-word Stroop task while undergoing an fMRI scan (**C**). To derive neural correlates of the theorized task dimensions (in **A–B**), a general linear model estimated the BOLD response evoked by 16 unique Stroop stimuli (e.g., “WHITE” displayed in blue hue) independently for each voxel. We analyzed a mostly incongruent trial condition to obtain balanced estimates of each trial type. The Glasser et al. (2016) multimodal atlas was used to parcellate cortex (**D**, light silver borders). Contiguous sets of parcels that tiled our ROIs (“superparcels”) were defined and treated as analytic units (for list, see Extended Data Figure 1-1). Within each superparcel, linear correlations among response patterns from the 16 stimuli were estimated to form an empirical similarity matrix (**E**, left; stimuli that were white are presented in gray within this manuscript). Through rank regression, these matrices were fit to three representational models (**E**, right), which corresponded to the three hypothesized dimensions of the Stroop task (**A**). The resulting β coefficients summarized the extent to which a parcel emphasized, within its distributed activity patterns, the representation of each unique task dimension. These β coefficients were used as the primary dependent variables in group-level analyses, and primary independent variables in individual-level analyses (e.g., as in **B**). Critically, we verified the specificity of our design and analysis via simulation (Extended Data Figure 1-2; compare Cai et al., 2019).

this field has been to use measures of neural activity evoked by response conflict tasks, such as Stroop, to test models of cognitive control.

One broad, neurocomputational-level model ascribes particular roles to different frontoparietal regions in overcoming response conflict (Miller and Cohen, 2001; Shenhav et al., 2013). Central to this view is the type of task information these regions encode. The dorsomedial frontal cortex (DMFC) is proposed to “evaluate” demand for cognitive control, via encoding of incongruency-related information (see Fig. 1A, bottom row). Such information, according to this view, is used by dorsolateral PFC (DLPFC), in concert with lateral posterior parietal cortex (LPPC), to “implement” control, via encoding of goal and target-related information (see Fig. 1A, top row). Thus, this model predicts key functional dissociations between medial and lateral frontoparietal cortex (see Fig. 1B). But although this view has been influential, directly establishing these dissociations during the performance of standard color-word Stroop tasks has been difficult.

To date, the most traction on this problem has been gained via fMRI designs that temporally dissociate presentation of task-

rule and incongruency-related information, in which subjects were instructed before each Stroop trial about which task to perform (color-naming, word-reading; MacDonald et al., 2000; Floden et al., 2011). But, while these studies generally found supportive evidence for the key claims, results were subject to three notable limitations. First, these studies were likely underpowered for fMRI (e.g., $N = 12$ in Floden et al., 2011; $N = 9$ in MacDonald et al., 2000). This fact alone warrants a follow-up study. Second, it is unclear whether the results extend to the more-standard Stroop-task design, in which task rules are not explicitly instructed before each trial, but are instead internally maintained. For example, goal-relevant coding in DLPFC may depend on such explicit rule instruction. Third, the prior results do not speak to functional dissociations within a single Stroop trial, during which interference is actually experienced and resolved. It is therefore possible, for instance, that the role of DLPFC (or other frontoparietal regions) in Stroop is primarily preparatory, and is less critical during actual interference resolution.

To address these questions, a neuroanatomically precise technique is needed that does not rely on temporal dissociations, but can instead read out multiple, simultaneously encoded sources of

task information from individual brain regions of interest (ROIs). Multivariate (multivoxel) pattern analysis (MVPA) of fMRI, in popular use for over a decade (Edelman et al., 1998; Haxby et al., 2001; Cox and Savoy, 2003), accomplishes exactly this purpose. Surprisingly, however, these methods have not been brought to bear on the question of a functional dissociation between medial and lateral frontoparietal cortex in resolving Stroop conflict.

We fill this gap in the literature by using representational similarity analysis (RSA; Edelman et al., 1998; Kriegeskorte et al., 2008), a specific MVPA framework, to test for dissociations in frontoparietal coding during Stroop-task performance (see Fig. 1C–E). We conducted a retrospective analysis of data collected as part of the Dual Mechanisms of Cognitive Control project (Braver et al., 2020). Our primary goal was proof-of-principle: to demonstrate the potential of RSA for testing theorized distinctions in neural coding within cognitive control tasks, such as the Stroop (Freund et al., 2021).

Materials and Methods

We report how we determined our sample size, all data exclusions, all manipulations, and all measures in the study (Simmons et al., 2011).

Code, data, and task accessibility

Code (R Core Team, 2019) and data to reproduce all analyses, in addition to supplementary analysis reports, are publicly available (<https://doi.org/10.5281/zenodo.4784067>). As part of the planned data release of the Dual Mechanisms of Cognitive Control project, raw and minimally preprocessed fMRI data have been deposited on OpenNeuro (<https://doi.org/10.18112/openneuro.ds003465.v1.0.3>). Additionally, the authors will directly share the specific fMRI data used for this study on request. Task scripts are available at the project website (<http://pages.wustl.edu/dualmechanisms/tasks>). More detailed information regarding all aspects of the project can be found on the Project's OSF page (<https://osf.io/xfe32/>).

Participants

Individuals were recruited from the Washington University and surrounding St. Louis metropolitan communities for participation in the Dual Mechanisms of Cognitive Control project. The present study began with a subset ($N=66$; 38 women, 26 men, 1 “prefer not to answer”) of these subjects: those with a full set of imaging and behavioral data from the Stroop task during a particular scanning session (the “proactive” session), selected for methodological reasons (see *Selection of data*). One subject was excluded from all analyses because of a scanner error. We split the remaining sample into two sets of individuals: a primary analysis set ($N=49$; 27 women, 21 men, 1 “prefer not to answer”), which we used in all analyses, and a validation set ($N=16$; 11 women, 5 men), which was only used in the Model selection analysis (see below). This unbalanced partitioning was done to account for the familial structure present within our sample. Specifically, subjects within each set (primary validation) were all unrelated; however, subjects within the validation set were co-twins of 16 subjects within the primary analysis set. Two of these co-twins were selected for use in the primary analysis set as their respective co-twins had atypically high rates of response omission (>10%; >20% errors of any type); the remaining co-twins were randomly selected. Critically, partitioning the sample in this way ensured that the primary analysis set was a random sample of independent subjects.

The partitioning of the data into two subsets also afforded the opportunity to use the validation subset as held-out data for evaluation of the brain–behavior model within the Model selection analysis. As we performed this sorting of individuals into primary and validation sets only once and did not analyze the validation-set data (except to assess predictive accuracy of the final selected model) the validation set provides an unbiased assessment of predictive accuracy, in the sense that no statistical “double-dipping” could have occurred. But because the sets are

familially dependent, it is perhaps more accurate to consider the validation-set analyses as assessing a kind of test–retest reliability (i.e., while eliminating the potential confound of practice effects), rather than providing an estimate of out-of-sample predictive accuracy. To evaluate this matter, follow-up control analyses were conducted in which the co-twins were removed from the primary analysis set.

Experimental design and statistical analysis

Task. Participants performed the verbal color-word Stroop (1935) task. Names of colors were visually displayed in various hues, and participants were instructed to “say the name of the color, as fast and accurately as possible; do not read the word.”

The set of stimuli consisted of two subsets of color-word stimuli (randomly intermixed during the task): a *mostly incongruent* and an *unbiased* set. Each stimulus set was created by pairing four color words with four corresponding hues in a balanced factorial design, forming 16 unique color-word stimuli within each set. The *mostly incongruent* set consisted of stimuli with hues (and corresponding words) ‘blue’ (RGB = 0, 0, 255), ‘red’ (255, 0, 0), ‘purple’ (128, 0, 128), and ‘white’ (255, 255, 255); the *unbiased* set, of ‘black’ (0, 0, 0), ‘green’ (0, 128, 0), ‘pink’ (255, 105, 180), and ‘yellow’ (255, 255, 0). These words were centrally presented in uppercase, bold Courier New font on a gray background (RGB = 191, 191, 191). Of stimuli within the *mostly incongruent* set, incongruent stimuli were presented to subjects more often than congruent stimuli (per block, proportion congruent = 0.25). *Unbiased* stimuli were presented with a balanced frequency (proportion congruent = 0.5). These manipulations of incongruency statistics are standard manipulations to elicit proactive control (Bugg, 2014; e.g., Gonthier et al., 2016) and were performed to investigate questions outside the scope of the current study. Thus, as described further below, the *unbiased* stimulus set was excluded from all analyses.

Each trial (see, e.g., Fig. 1C) began with a central fixation cross, presented for 300 ms on a gray background (RGB = 191, 191, 191). The color-word stimulus, preceded by a blank screen following fixation offset (100 ms), was centrally presented for a duration of 2000 ms, fixed across trials. The duration of the intertrial interval (triangle of fixation crosses) was 900, 2100, or 3300 ms, selected randomly (with uniform probability). Each of two scanning runs consisted of three blocks of 36 trials, intermixed with four resting fixation blocks, during which a fixation cross appeared for 30 s. This formed a mixed block–event design (Petersen and Dubis, 2012). Each of the 16 *mostly incongruent* stimuli — that is, each unique colored word (e.g., “BLUE” displayed in red hue) — was presented in both runs. Within each run for each participant, mostly incongruent stimuli were presented an equal number of times within each block. Within each block, stimulus order was fully randomized.

Selection of data. We focused our fMRI pattern analyses solely on trials from the *mostly incongruent* stimulus set within a particular scanning session (the “proactive” session) of our Stroop task. This selection was made purely on the basis of methodological reasoning: these trials were the only set of trials within the larger Dual Mechanisms project in which each unique Stroop stimulus (i.e., one of the 16 color-word combinations) was presented an equal number of times (9) to each participant, constituting a balanced design. Balanced designs ensure that differences in the total number of trials per condition cannot explain any differences observed in pattern correlations among conditions.

Display and recording systems. The experiment was programmed in E-Prime 2.0 (2013, Psychology Software Tools), presented on a Windows 7 Desktop, and back-projected to a screen at the end of the bore for viewing via a mirror head-mount. Verbal responses were recorded for offline transcription and response time (RT) estimation. The first 45 participants spoke into a MicroOptics MR-compatible electronic microphone (MicroOptics Technologies); because of mechanical failure, however, we replaced this microphone with the noise-cancelling FOMRI III (OptoAcoustics), which subsequent participants used. A voice-onset processing script (from the MATLAB Audio Analysis Library) was used to derive RT estimates on each trial via spectral decomposition (the accuracy of which was verified by manually coding RTs from a subsample of subjects and ensuring the two methods gave

similar estimates). Code for this algorithm is available within the Dual Mechanisms GitHub repository (<https://github.com/ccplabwust/dualmechanisms/tree/master/preparationsAndConversions/audio>).

Importantly, we verified that the change in microphone did not induce confounding between-subject variance in RT measures of interest. While RT estimates recorded via the Micro-Optics microphone tended to be slower ($b = 102.59$, $p = 0.01$) and more variable ($\chi^2_5 = 3655$, $p = 0$), the magnitude of the Stroop effect was not observably impacted by the microphone change ($b = -5.88$, $p = 0.64$).

Image acquisition, preprocessing, and GLM. The fMRI data were acquired with a 3T Siemens Prisma (32 channel head-coil; CMRR multi-band sequence, factor = 4; 2.4 mm isotropic voxel, with 1200 ms TR, no GRAPPA, ipat = 0), and subjected to the minimally preprocessed functional pipeline of the Human Connectome Project (version 3.17.0), outlined by Glasser et al. (2013). More detailed information regarding acquisition can be found on the Project OSF site (<https://osf.io/tbhfq/>). All analyses were conducted in volumetric space; surface maps are displayed in figures only for ease of visualization. Before revision of this manuscript, the data were reprocessed with fMRIPrep, using the standard fMRIPrep pipelines (Esteban et al., 2019, 2020). At this point, the preprocessed results with the HCP pipelines were inadvertently removed. Thus, some follow-up control analyses were conducted with the fMRIPrep-preprocessed data (Extended Data Figs. 2-4, 3-2, 4-4). The fMRIPrep pipeline was implemented in a Singularity container (Kurtzer et al., 2017) with additional custom scripts used to implement file management (more detail on the pipeline is available at <https://osf.io/6p3en/>; container scripts are available at <https://hub.docker.com/u/ccplabwust/>).

After preprocessing, to estimate activation patterns, we fit a whole-brain voxelwise GLM to BOLD time-series in AFNI, version 17.0.00 (Cox, 1996). To build regressors of primary interest, we convolved with an HRF [via AFNI's *BLOCK(1,1)*] 16 boxcar time courses, each coding for the initial second of presentation of a *mostly incongruent* stimulus that resulted in a correct response. We also included two regressors [similarly created via *BLOCK(1,1)*] to capture signal associated with congruent and incongruent trials of noninterest (*unbiased* stimuli) that prompted correct responses, an error regressor coding for any trial in which a response was incorrect or omitted (via *BLOCK*), a sustained regressor coding for task versus rest (via *BLOCK*), a transient regressor coding for task-block onsets [as a set of piecewise linear spline functions via *TENTzero(0,16.8,8)*], six orthogonal motion regressors, five polynomial drift regressors (order set automatically) for each run, and an intercept for each run. These models were created via *3dDeconvolve* and solved via *3dREMLfit*. The data for each subject's model consisted of 2 runs \times 3 blocks \times 36 trials (144 from the *mostly incongruent* stimulus group, 72 from unbiased). Frames with $FD > 0.9$ were censored.

Definition of ROIs. Our primary hypotheses concerned a set of six anatomic regions: DMFC, DLPFC, and LPPC in each hemisphere. Consequently, our primary analyses used a targeted ROI-based analysis approach. Rather than defining functional ROIs via a whole-brain searchlight, which has known issues (Etzet et al., 2013), we defined ROIs via a cortical parcellation atlas. We selected the MMP atlas (Glasser et al., 2016) for two reasons: (1) the atlas was developed recently via multi-modal imaging measures; and (2) individual MMP parcels are relatively interpretable, as they are heterogeneously sized and have been explicitly connected to a battery of cognitive tasks (Assem et al., 2020), the canonical functional connectivity networks (Ji et al., 2019), and a large body of neuroanatomical research (Glasser et al., 2016). We used a volumetric version, obtained from https://figshare.com/articles/HCP-MMP1_0_projected_on_MNI2009a_GM_volumetric_in_Nifti_format/3501911?file=5534024 (also available on the project GitHub repository; see *Code, Data, and Task accessibility*). We then defined a set of six spatially contiguous sets of MMP parcels (three in each hemisphere), which we refer to as “superparcels,” that corresponded to each of our ROIs. For full superparcel definitions, see Extended Data Fig. 1-1. DMFC was defined as the four parcels covering SMA-pre-SMA and dACC. DLPFC was defined as the four parcels that cover middle frontal gyrus (i.e., mid-DLPFC). LPPC was defined as all parcels tiling IPS, from the occipital lobe to primary somatosensory cortex. The overwhelming majority of

parcels that met these anatomic criteria were assigned, within a previous report, to the cinguloopercular (most of DMFC), frontoparietal (most of DLPFC), and dorsal-attention (most of LPPC) control networks (Ji et al., 2019). Further, these ROI definitions contain several parcels that correspond to key nodes within the “multiple demand” network (Assem et al., 2020). To assess the robustness of our results to particular superparcel definitions, we additionally used alternative, more inclusive, superparcel definitions of DMFC and DLPFC (see Extended Data Fig. 1-1). For the brain-behavior model selection analysis (see *Model selection*), we compiled a larger set of anatomically clustered MMP parcels, covering regions across the cortex (Extended Data Fig. 3-3). Two additional, non-MMP ROIs were included in this set, to give better coverage of particular functional brain regions. A mask for ventral somatomotor cortex (the “SomatoMotor-Mouth” network) was obtained from the Gordon atlas (Gordon et al., 2016), as the MMP does not split somatomotor cortex into dorsal and ventral divisions. A mask for left ventral occipito-temporal cortex (encompassing the “visual word-form” area) was obtained using MNI coordinates $-54 < x < -30$, $-70 < y < -45$, $-30 < z < -4$, specified in a prior report (Twomey et al., 2011). To remove cerebellar voxels from this ROI, we used the Deidrichsen atlas (Diedrichsen, 2006) hosted by AFNI (https://afni.nimh.nih.gov/pub/dist/atlas/SUIT_Cerebellum/SUIT_2.6.1/).

Estimation of coding strength β . To estimate the regional strength of target, distractor, and incongruity coding, we used the RSA framework (Kriegeskorte et al., 2008). The RSA framework consists of modeling the observed similarity structure of activation patterns with a set of theoretically specified model similarity structures (see Fig. 1E). For a given subject and cortical region, fMRI GLM coefficient estimates for each of the 16 conditions of interest (four colors factorially paired with four words; e.g., the word “WHITE” presented in blue hue) were assembled into a condition-by-voxel activity pattern matrix **B**. The observed similarity structure was estimated as the condition-by-condition correlation matrix $\mathbf{R} = \text{Cor}(\mathbf{B})$. Cell R_{ij} of this matrix gives the linear correlation observed between activity patterns evoked by conditions i and j . Model similarity structures were specified in this same correlation matrix form. The target model assumed that conditions (stimuli) with the same hue will evoke identical patterns, regardless of whether the words or congruency match (or mismatch). That is, if the hue of condition $i = j$, this model predicts $R_{ij} = 1$, otherwise 0. (In the “target” matrix in Fig. 1E, the only cells equal to 1 — i.e., blue cells — are those in which the stimulus hues match.) The distractor model assumed that conditions with the same word will evoke identical patterns, regardless of the hue or congruency. (If the word of condition $i = j$, $R_{ij} = 1$, otherwise 0. In the “distractor” matrix in Fig. 1E, the only cells equal to 1 are those in which the stimulus words match.) The incongruity model assumed that only conditions that were incongruent would evoke identical patterns, regardless of the hue or word. (If i and j are both incongruent, $R_{ij} = 1$, otherwise 0. In the “incongruity” matrix in Fig. 1E, the only cells equal to 1 are those in which both stimuli are incongruent.) As a covariate of noninterest, we also included a model capturing similarity between congruent conditions (if i and j are both congruent, $R_{ij} = 1$, otherwise 0). We additionally examined alternative approaches to RSA modeling of incongruity, to see whether our results were robust to this parameterization of incongruity and congruency (see *Alternative RSA incongruity models*).

These four models were jointly fitted to the observed similarity structure from each region through multiple regression (ordinary least-squares), separately for each subject. The response vector \mathbf{y} and design matrix of this regression were assembled in a series of steps. (1) The 120 unique off-diagonal elements of each similarity matrix (one observed and four models) were extracted and unwrapped into vectors. (2) The four model similarity vectors were separately z -scored and assembled into columns. This formed the RSA design matrix. (3) The observed similarity vector was rank-transformed (Nili et al., 2014) then z score standardized, to form a vector \mathbf{r} . (4) The vector \mathbf{r} was prewhitened to remove a specific nuisance component. This component stemmed from the task design: though each *mostly incongruent* stimulus occurred an equal number of times throughout the course of a session, these stimuli were not fully balanced across the two scanning runs. Specifically, half of the

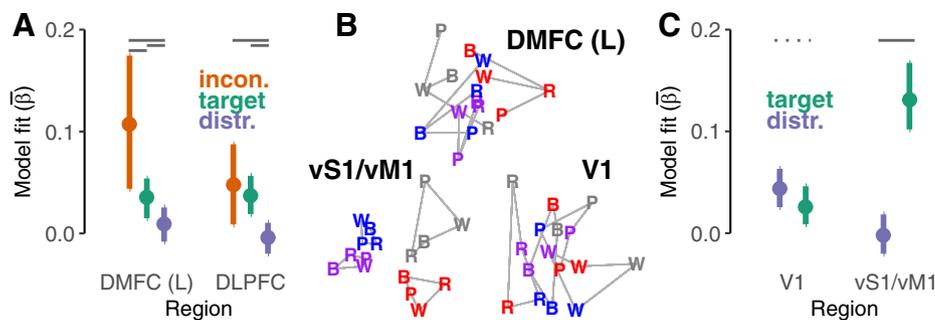


Figure 2. Group-level results. **A**, A single dissociation between DLPFC versus DMFC. For simplicity of display, DLPFC estimates are averaged across hemisphere (for per-hemisphere means for each region, see Table 1). Error bars indicate 95% CIs of between-subject variability (estimated via bias-corrected and accelerated bootstrap). Horizontal gray lines at the top indicate significance of within-subject significance tests (for contrast estimates, see Table 1). While all ROIs encoded target and incongruity information, and did so more strongly than distractor information (coding of which was not detected), incongruity coding was stronger in left DMFC versus DLPFC. These results were generally robust to different analysis decisions and implementations (Extended Data Figures 2-1, 2-2, 2-4, 2-5). Further, a univariate analysis failed to detect significantly higher activation on incongruent versus congruent trials in these regions (Extended Data Figure 2-3). **B**, Across-voxel activity patterns from three select regions, embedded within a two-dimensional space (via nonmetric multidimensional scaling). These geometries exemplify regional dissociations in coding of Stroop-task dimensions. Letters represent the first letter of the corresponding distractor word (BLUE, PURPLE, RED, WHITE), and hues represent the target color (stimuli that were white are presented here in gray). Connecting lines are drawn to highlight the various structures present. The geometry of (left) DMFCs marked by a radial separation of patterns evoked by congruent and incongruent conditions (lines connect congruent and incongruent stimuli of each target color), such that incongruent patterns tend to be located more centrally (i.e., more similar to each other), whereas congruent patterns peripherally diverge (i.e., less similar to each other). This suggests that incongruent trials drove a common component of activation in this region, regardless of target or distractor features. In contrast, within primary ventral somatomotor cortices (vS1/vM1), patterns strongly cluster by target level (color), whereas in primary visual cortex (V1), they tend to cluster by distractor level (word). This pattern formed a double dissociation (C), and indicated our distractor model was adequately powered to detect distractor coding in sensory regions (compare A). **C**, Plotting conventions follow those in A. Horizontal dotted lines indicate $0.1 > p > 0.05$ (for estimates, see text).

stimuli were presented 3 times in the first run versus 6 times in the second (vice versa for the other half). As each scanning run contains a large amount of run-specific noise (Mumford et al., 2014; Alink et al., 2015), this imbalance across runs could lead to a bias in the resulting β coefficients, in which pattern similarity of stimuli that mostly occurred within the same run would be inflated. We formalized this component of bias as another model similarity vector, \mathbf{v} , with elements equal to 1 if the run in which condition i most frequently occurred = the run in which condition j most frequently occurred, otherwise 0. The magnitude of this bias was estimated as the slope term b_1 in a linear regression $\mathbf{r} = \mathbf{v}b_1 + b_0 + \mathbf{e}$, where b_0 is the intercept coefficient and \mathbf{e} is the residual vector. The model \mathbf{v} was scaled by its magnitude and then subtracted from \mathbf{r} , forming the RSA response vector $\mathbf{y} = \mathbf{r} - \mathbf{v}b_1$. We additionally used an alternative, downsampling technique to verify that our primary findings were robust to this issue (see *Downsampling analysis*).

Thus, the RSA regression yielded three β coefficients of interest: β_{target} , $\beta_{\text{distr.}}$, $\beta_{\text{incon.}}$. These coefficients can be understood as a (standardized) contrast on (rank-transformed) correlations of activity patterns, between conditions in which only one task dimension was shared (e.g., the target dimension for β_{target}), versus those in which no dimensions were shared (i.e., different levels of target, distractor, and congruency).

Dimensionality reduction. We used non-metric multidimensional scaling (Kruskal, 1964), a flexible, nonparametric dimensionality reduction technique, to visualize the structures of activity patterns within selected regions (see Fig. 2): ventral somatomotor cortex (corresponding to “mouth” homunculi), primary visual cortex (V1), and (left) DMFC. These parcels were selected to highlight coding of each task dimension. For each selected region, we averaged observed correlation matrices across subjects and then subtracted these values from 1 to obtain a dissimilarity matrix. Before averaging, we z -transformed (inverse hyperbolic tangent, artanh) correlations, and inverted this transform after averaging. (In contrast to the RSA regression above, we did not rank-transform correlation matrices, as non-metric multidimensional scaling incorporates a monotonic regression). Similar to our RSA, we prewhitened each similarity matrix before conducting this procedure (see Step 4 in *Estimation of coding strength β*). Each mean dissimilarity matrix was submitted to an implementation of Kruskal’s nonmetric multidimensional scaling, *vegan::metaMDS()* in R, to generate a two-dimensional configuration (Oksanen et al., 2019).

Group-level dissociation analysis. To test for regional dissociations in coding preferences, we fit a hierarchical linear model on RSA model fits

(see *Estimation of coding strength β*) obtained from our three ROIs within each hemisphere, and for our three RSA models. Fixed effects were estimated for the interaction of RSA model, ROI, and hemisphere. Random effects by subject were estimated for the interaction of RSA model and ROI, with a full covariance structure (9×9 ; Barr et al., 2013). This model was fit with *lme4::lmer()* in R (Bates et al., 2014).

Planned contrasts on the fixed effects were performed to test our hypotheses. p values were estimated using an asymptotic z -test, as implemented by the *multcomp::glht()* function in R (Hothorn et al., 2008). We performed three types of contrasts: (1) to compare coding strengths within-region (e.g., DMFC: incon.–target); (2) to compare between regions, within-model (e.g., target: DLPFC – DMFC); and (3) to test their interaction [e.g., (target – incon.) · (DLPFC – DMFC)]. These contrasts were performed first by collapsing across hemisphere, then within each hemisphere separately. As we did not have any hypotheses regarding lateralization, p values from hemisphere-specific contrasts were FDR-corrected across hemispheres.

Selection of behavioral measures for individual-level analyses. Audio recordings of verbal responses were transcribed and coded for errors offline by two researchers independently. Discrepancies in coding were resolved by a third. Errors were defined as any nontarget color word spoken by a subject before utterance of the correct response (e.g., including distractor responses, but not disfluencies) or as a response omission. Trials in which responses were present but unintelligible (e.g., because of high scanner noise or poor enunciation) were coded as such.

We fit two hierarchical models on these data: one on errors and one on RTs. Several observations were excluded from these models. From the error model, only trials with responses coded as “unintelligible” were excluded (54). From the RT model, several types of trials were excluded: trials with RTs >3000 ms (1) or <250 ms (53; 52 of which were equal to zero). A cluster of fast and unrealistically invariable RTs from 2 subjects (23/216, 26/216) that were likely because of an artifact of insufficient voice-onset signal within the recording. Trials with a residual RT that was more extreme than three interquartile ranges from an initial multilevel model fitted to all subjects data (of the structure of the model equation below; as in Baayen and Milin, 2010). All were trials with incorrect (137), unintelligible (54), or no response (52). In total, 232 trials were excluded (0–62 per subject), leaving 10,352 trials for analysis (154–216 per subject).

RTs for subject s were modeled (following Laird and Ware, 1982 notation) as follows:

$$\begin{aligned} RT_s &= (\mathbf{1}, \mathbf{stroop}_s)\mathbf{b} + (\mathbf{1}, \mathbf{stroop}_s)\mathbf{u}_s + \mathbf{e}_s \\ \mathbf{u}_s &\sim N_2(\mathbf{0}, \mathbf{T}_{2 \times 2}^2) \\ \mathbf{e}_s &\sim N(0, c_s \sigma^2) \end{aligned}$$

where RT_s is a column vector of the RTs from subject s , $\mathbf{1}$ is an all-ones vector (intercepts), \mathbf{stroop}_s is a vector indicating incongruent trials, and c_s is a subject-specific parameter by which their residual variance was scaled. Critically, \mathbf{b} and \mathbf{u}_s contained coefficients corresponding to the classical Stroop interference effect contrast (incongruent – congruent), for the group (\mathbf{b}) and subject (\mathbf{u}_s) levels. The scaling parameters c_s relaxed the assumption that each subject had a common residual variance, a well-warranted complexity, given the vastly improved fit of the heterogeneous-variance model ($\chi^2_{48} = 4299, p = 0, \Delta BIC = -3856, \Delta AIC = -4203$). To accomplish this, the RT model was fitted with *nlme::lme()* in R (Pinheiro et al., 2019). The behavioral error model had similar predictors, but assumed binomially distributed error, with logit link function (fitted in *lme4*).

This hierarchical modeling framework enabled us to estimate the amount and internal consistency of individual variability in the Stroop interference effect within both RTs and errors while accounting for trial-level error (e.g., Haines et al., 2020). We used these subject-level estimates to validate that our behavioral measures met prerequisite properties for individual differences analyses. To assess the amount of between-subject variability in the Stroop interference effect, a nested model comparison was conducted, in which the models fitted above were compared with a random-intercept model. Stroop effects differed significantly across subjects in RT ($\chi^2_2 = 87.53, p = 0.00; \Delta BIC = -69.08; \Delta AIC = -83.53$), but not in accuracy ($\chi^2_2 = 0.69, p = 0.71; \Delta BIC = 18; \Delta AIC = 3$). We therefore did not further analyze accuracy. To assess internal consistency (defined here as cross-run correlation), we fit a model with separate congruency factors (fixed and random) per run, and a full 4×4 covariance structure, which was used to obtain the cross-run correlation in Stroop effect. Individuals' Stroop interference effects in RT were estimated to be highly consistent across scanning runs ($r = 0.95$). We therefore considered our measurements of RT to be adequate for individual differences analysis.

Individual-level dissociation analysis. Similar to our Group-level dissociation analysis, we tested our individual-level hypotheses within a hierarchical modeling framework. Preliminary analyses suggested that error measures were inadequate for individual differences analyses (see above), so we focused solely on RT measures.

We began with the RT model described in the preceding section. However now, for a given RSA model and ROI, we incorporated into the fixed effects each subject's estimated coding strength, β_s , by interacting this coding-strength term with the congruency factor. This formed a cross-level, continuous-by-categorical interaction, $\beta_s \cdot \mathbf{stroop}_s$. The coefficient on this interaction term described how the Stroop interference effect (within-subject) varied across subjects as a function of their coding strength (between-subject). To test our hypotheses, we performed contrasts on these interaction coefficients, which are outlined within Results describing Figure 3.

Model selection. To complement our hypothesis-driven brain-behavior analyses, we used a more data-driven model-selection approach. An expanded set of 24 superparcels (in addition to our six ROIs) was defined (see Fig. 3B; for list, see Extended Data Fig. 3-3). Some superparcels were included as ROIs, others were included as negative controls (i.e., regions that were not predicted to be important for explaining behavioral performance). Subject-level coefficients of the Stroop interference effect contrast were extracted from the behavior-only RT model and used as the response vector (in the model equation, the slope elements of \mathbf{u}_s ; also known as BLUPs or conditional modes). RSA was conducted on each superparcel (see *Estimation of coding strength β*), furnishing three β coefficients per superparcel. These 72 measures were fitted to the response vector via elastic net regression, implemented via *glmnet::glmnet()* in R (Friedman et al., 2010). Parameter α was set to 0.5 (balancing lasso and ridge penalties). The parameter λ was tuned via 10-fold cross-validation (default) using the "1-SE rule": the smallest $\lambda > 1$

SE of the minimum λ across folds was saved (Hastie et al., 2009); this routine was repeated 1000 times, and the minimum mode (there were ties) of the saved λ s was selected. We selected the minimum mode because the maximum suppressed all variables from the model.

To assess validation-set accuracy, the selected model coefficients were applied to a design matrix from validation-set subjects, generating a predicted Stroop effect vector. The linear correlation was estimated between this predicted Stroop effect and observed Stroop effects (estimated as conditional modes via a hierarchical model separately fitted to validation-set data). The significance of this correlation was assessed by randomly permuting the training-set response vector, refitting the model, generating new predicted validation-set values, and re-estimating the predicted–observed correlation 10,000 times. The p value was given by the proportion of resamples in which the null correlation was greater than the observed correlation.

Exploratory whole-cortex RSA. The RSA-model fitting procedure, as outlined in Estimation of coding strength β was separately conducted on each MMP cortical parcel. Inferential statistics followed those suggested by Nili et al. (2014). One-sided signed-rank tests were conducted for significance testing (>0). p values were FDR-corrected over all 360 parcels, separately for each task dimension (target, distractor, incongruency).

Univariate activation analyses. We additionally conducted a standard "univariate activation" analysis on these data. This was not meant to evaluate whether univariate activity was a plausible confounding variable in our analysis, but rather to provide some basis for comparing our data to most extant neuroimaging studies of Stroop. For a given ROI (or MMP parcel), β coefficients from the first-level fMRI GLM were averaged over voxels by stimuli, then over stimuli by congruency. These mean values were then contrasted, analogous to the behavioral Stroop interference effect (incongruent – congruent). This statistic gives an estimate of the overall (across-voxel) difference in fMRI activity within a given brain region on incongruent versus congruent trials.

Follow-up control analyses

To establish the robustness of our results, we conducted several control analyses that examined a number of confounds and concerns: potential differences in signal-to-noise ratio (SNR) across prefrontal ROIs, the effects of head motion, different RSA models for incongruency coding, the presence of bias imposed by the experimental design, within-run versus between-run RSA estimation, and the effects of downsampling to account for different trial numbers across runs. These are each described next.

Comparison of SNR ratios. To test for differences in SNR between DMFC and DLPFC, we estimated "noise ceilings" within each region and contrasted them across regions. Noise ceilings indicate the maximum observable group-level effect size (RSA model fit) given the level of between-subject variability in similarity structure (Nili et al., 2014). Lower (smaller) average noise ceilings indicate poorer SNR for group-level tests. We used the cross-validated "lower-bound" noise ceiling estimator of Nili et al. (2014), as this yields a lower-variance estimate, and therefore more powerful contrast across regions, than the non-cross-validated "upper-bound" (Hastie et al., 2009). For a given region, the lower-bound noise ceiling is defined for each subject s in $1, \dots, N$ as $\text{Cor}(\mathbf{y}_s, \bar{\mathbf{y}}_{-s})$: the linear correlation between a subject s ' observed similarity vector, \mathbf{y}_s , and the group-mean vector excluding subject s , $\bar{\mathbf{y}}_{-s}$ (for definition of \mathbf{y} , see *Estimation of coding strength β*). As each of the N estimates is interdependent, we used a percentile bootstrap to contrast noise ceilings across regions. In each resample (of 10,000), noise ceilings were artanh transformed, contrasted across regions within-subject, then averaged across subjects. A two-sided p value was provided by computing the proportion of resampled means greater than zero, p , then taking the minimum of $2p$ and $2(1 - p)$. Finally, we conducted "two one-sided tests" for equivalence (Schuirmann, 1987; Lakens, 2017) to affirm a null hypothesis of no difference between regions in noise ceiling. This consists of defining a threshold effect size, the minimum effect size of interest, and testing whether the observed difference in noise ceilings is significantly less extreme than the threshold. For an objective threshold, we used the smallest standardized effect size at which our bootstrap

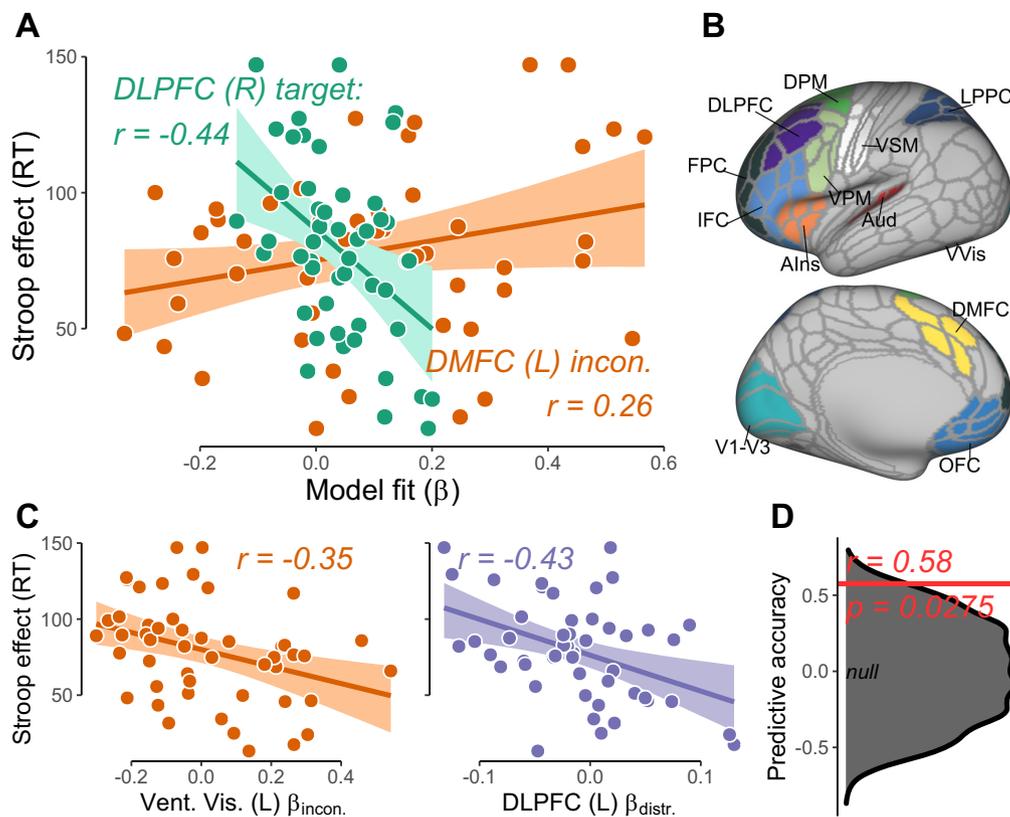


Figure 3. Individual-level results. **A**, The coding strength of task variables (β) within right DLPFC and left DMFC accounted for individual differences in the size of the Stroop interference effect (RT) in hypothesized ways. Inset text represents linear correlation coefficients (r) and 95% CIs (estimated via bias-corrected and accelerated bootstrap). Lines indicate a regression line of RT onto Model Fit (β), and confidence bands indicate 95% CIs of predicted values from this regression obtained via percentile bootstrap. The direction of these observed brain–behavior relationships was robust to alternative, more conservative RSA estimation techniques (cross-run RSA and cross-validated RSA; Extended Data Figure 3-2). For display of relationships between target and incongruity coding in all ROIs (DLPFC, DMFC, LPPC) and Stroop effect, see Extended Data Figure 3-1. For contrasts on slopes (interaction tests between ROI and RSA model), see Table 2. **B**, A larger set of 24 cortical regions were defined for use in a data-driven model selection analysis. Surface maps display in various colors all left-hemisphere regions (excluding left lateral occipito-temporal cortex; for full list and definitions, see Extended Data Figure 2-1). Three RSA models (target, distractor, incongruity) were fitted per region, yielding a set of 72 coding-strength estimates. These 72 estimates were used to predict Stroop interference effects within a cross-validated model selection procedure. **C**, Additional brain–behavior associations identified via cross-validated model selection. Left ventral visual incongruity coding (left panel) and left DLPFC distractor coding (right panel) were selected — in addition to right DLPFC target, right LPPC target, and left DMFC incongruity. For description of inset text, lines, and bands, see **A**. **D**, Estimate of prediction accuracy of the selected model. Within a validation (held-out) set of 17 subjects (monozygotic co-twins of subjects within the analysis sample), observed Stroop effects were moderately correlated (red line) to predictions from the selected model. A permutation test (gray distribution) indicated the selected model significantly explained validation-set variance, suggesting that the selected model results are relatively stable.

procedure was expected to retain 80% power, Cohen’s $d = 0.35$, determined via Monte Carlo simulation.

RSA on head motion estimates. As a negative control analysis for our exploratory whole-cortex RSA, we attempted to decode task variables (target, distractor, incongruity) via RSA from framewise estimates of head motion. The 6 motion regressors that were used in the fMRI GLM as nuisance covariates (corresponding to translation and rotation in 3 dimensions) were regressed on the design matrix containing predicted BOLD timecourses of our 16 conditions of interest. The coefficient matrix resulting from this regression was then submitted to the RSA procedures described in *Estimation of coding strength β* and *Exploratory whole-cortex RSA*. To check whether more aggressive movement denoising within the fMRI GLM was warranted (i.e., in addition to the 6 nuisance regressors), we conducted this same movement-based RSA, however, using 12 motion regressors (the 6 bases and their temporal derivatives). RSA model fits between the 6 and 12-basis motion-based RSAs were compared via paired-sample signed-rank test.

Alternative RSA incongruity models. The RSA incongruity model parameterized the congruent–congruent correlations (i.e., R_{ij} where i and j are both patterns from congruent trials, denoted here simply as CC) with a separate nuisance regressor. That is, these cells were effectively excluded, and the model instead computed the contrast $II - IC$. This exclusion was done because we have no specific hypotheses regarding how congruent trials should be encoded relative to one

another. Other parameterizations are possible, however, including models that (a) incorporate CC correlations within the “baseline” or intercept term, by omitting the congruent nuisance regressor [i.e., $II - ave(IC + CC)$], or (b) that omit IC correlations from the contrast (i.e., $II - CC$). We note, however, that (b) is a suboptimal parameterization as it effectively excludes 40% (48/120) of observations per subject (i.e., all IC cells). To verify that our results were not dependent on the particular modeling approach we chose, we compared the observed RSA model fits to these alternative parameterizations. In brief, subject’s model fits were highly similar between the parameterizations (see *Group and Exploratory whole-cortex RSA*; Extended Data Fig. 2-5). Thus, to help streamline the Results, we only conducted and reported analyses using the original RSA incongruity parameterization (in which the congruity model was a nuisance covariate).

Design bias. For the current study, a typical within-run form of RSA estimation was implemented, in which correlations were computed among activation patterns estimated within the same scanning run and first-level GLM. Within-run RSA has been criticized because it is susceptible to design biases that occur when trial orders are insufficiently randomized within the experiment (Cai et al., 2019). *A priori*, this was not a strong concern in the current design, as trial orders were fully randomized both within and between participants. Nevertheless, we conducted several diagnostic and robustness analyses to validate that our

results, and conclusions were not impacted by this potential bias. First, we estimated the extent of the potential bias by running through our RSA pipeline data simulated under a “worst-case” scenario, that is, when SNR = 0 (details described within Extended Data Fig. 1–2) and across a wide range of autocorrelation strengths. At worst, the three models were weakly biased (within 0.02–0.05 of $\alpha = 0.05$; Extended Data Fig. 1–2). This result indicated that while the bias was present, it was relatively minimal (compare Cai et al., 2019). Second, we validated that these simulated estimates were realistic, by using our actual fMRI data to estimate the false positive rate empirically. To do this, we conducted RSA on the first-level GLM coefficients from each subjects’ ventricles, as these voxels should contain no brain activity signal but similar noise characteristics as those of interest (a ventricle mask of 2431 voxels was obtained from AFNI servers: https://afni.nimh.nih.gov/pub/dist/tgz/suma_MNI_N27.tgz). By treating the group-level mean and SD of these RSA model fits as the parameters of a non-central null distribution, $X \sim N(\bar{\beta}, SD(\beta))$, we computed the empirical false positive rate as $P(X > z_{\alpha=0.05})$, the proportion of this distribution greater than the customary one-tailed z criterion. All rates were within 0.03 of $\alpha = 0.05$ (target = 0.07, distractor = 0.04, incongruency = 0.08), confirming the bias was quite minimal.

Between-run RSA. As an alternative to within-run RSA, various between-run estimation approaches have been proposed which have been shown to be less sensitive to potential design biases. We opted not to use between-run RSA for our primary analyses, both because of the reduced effects of design bias established above, but also because between-run RSA is noted to be considerably more conservative than within-run RSA (Cai et al., 2019). Moreover, several particulars of the present design are known to further hamper its power. Namely, between-run RSA makes incomplete use of the data, an issue that is exacerbated to the maximum extent possible in the present case, as our design has only the minimum number (two) of cross-validation folds (runs; Diedrichsen et al., 2020). Additionally, because the image acquisition sequence involved a reversal of phase-encoding direction across the two runs, this effectively adds a strong nonlinear component of noise if between-run RSA is used.

Nevertheless, to examine further the extent of design bias in our data, as well as the robustness of our results to the drop in power imposed by between-run RSA, we conducted a follow-up analysis of our primary results using between-run RSA approaches. We conducted two forms of between-run RSA: the first used “cross-run RSA”, which operates on the cross-correlation of patterns between scanning runs (see Alink et al., 2015), and the second used “cross-validated RSA”, which operates on the inner product of pattern contrasts between runs (see Walther et al., 2016). We selected these two forms of RSA as they have complementary benefits. Cross-run correlation is most comparable to our original within-run correlation, as they are both linear correlations. However, using this method within our data set also necessitated using downsampling (as the numbers of trials per condition were not perfectly balanced at the run-level; see *Downsampling analysis*), which increases the variance of resulting estimates because of discarding data. In contrast, cross-validated RSA is insensitive (in terms of expected value) to the issue of trial numbers per condition (Diedrichsen et al., 2020). Using this method, therefore allowed us to conduct the RSA using all the data at once, without downsampling. But cross-validated RSA tests a more constrained hypothesis than cross-run RSA. Whereas cross-run RSA can be sensitive to nonlinear differences between conditions, cross-validated RSA tests linear discriminability between conditions. Thus, when a nonlinear boundary separates conditions, cross-validated RSA will fail to detect an effect, whereas cross-run RSA could succeed. A nonlinear boundary would occur, for example, when one condition (e.g., incongruent stimuli) drives a reliable, common, response while the other conditions (e.g., congruent stimuli) either drive unreliable, or stable but heterogeneous, responses. Nevertheless, to make cross-validated RSA as comparable as possible to our primary RSA, we z -score standardized patterns before computation and omitted spatial prewhitening (Walther et al., 2016).

Downsampling analyses. Last, we checked whether the primary results were robust to the prewhitening method of data preprocessing,

which was introduced to handle the imbalance of trials across runs (see *Estimation of coding strength β* , step 4). In this analysis, we instead handled this issue by performing RSA after equating the number of trials per run*condition, by iteratively downsampling conditions with random subsets of trials. Specifically, we first fitted GLMs on fMRI time-series, separately for each scanning run, that contained a single regressor per trial (LS-A method of Mumford et al., 2012). The minimum number of times we presented each unique Stroop stimulus in a single run was 3; this was the number to which we downsampled all conditions with >3 occurrences. For these conditions, we randomly sampled three trials, and averaged GLM coefficients voxelwise over these trials. (For 3-trial conditions, we simply averaged all trials that were present.) This formed 16 separate condition-level coefficient vectors (activity patterns) per run. We then averaged these coefficient vectors across run, then estimated condition \times condition correlation matrices from these patterns (averaging across runs was omitted from the downsampled cross-run RSA; see preceding section). We repeated this resampling, averaging, and correlation process 1000 times, and averaged the resulting correlation matrices across iterations (after an artanh transform). These correlation matrices were then submitted to the same RSA as outlined in *Estimation of coding strength β* , with the omission of the prewhitening step (4).

Results

Influential theories of cognitive control have proposed specific dissociations in the type of task information encoded by human medial and lateral frontoparietal cortex (Fig. 1A,B). But previous studies have largely approached this question indirectly, by using tasks designed to recruit these regions differentially in time, then testing for temporal dissociations in regional-mean levels of fMRI activity. Here, we used a more direct approach, by using the similarity structure of neural activity patterns evoked within these regions to estimate their informational content. In particular, through RSA, we compared neural coding of three distinct types of Stroop-task information — target, distractor, and incongruency (Fig. 1D,E) — within each ROI simultaneously, while Stroop interference was being experienced and resolved.

We describe three sets of analyses. First, we examined group-level effects, to test for neuroanatomical dissociations in representation of task information (Fig. 1B, middle). Second, we examined individual-level effects (i.e., individual differences), to test for dissociations in brain–behavior relationships (Fig. 1B, right). These two analyses were ROI-based and primarily focused on dorsomedial frontal and lateral frontoparietal regions (Fig. 1D), but also included sensorimotor ROIs for comparison purposes. The last set of analyses was conducted in whole-cortex exploratory fashion, to provide a more comprehensive picture of the anatomic profile of each task dimension.

Group

DMFC and DLPFC exhibit distinct coding profiles

Primary group-level results are summarized in Figure 2. Statistical estimates corresponding to results outlined within this section are contained in Table 1.

The DMFC has been strongly associated with the coding of incongruency information in response conflict tasks, such as the Stroop. RSA approaches were used to directly test the specificity of that hypothesis. This region was indeed found to encode the incongruency dimension of the Stroop task (left: $\bar{\beta} = 0.11, p = 0.00$, right: $\bar{\beta} = 0.08, p = 0.00$). Additionally, within left DMFC, a preference for this dimension was observed: incongruency information was encoded more strongly than either target or distractor information (vs target: $\Delta\bar{\beta} = 0.07, p = 0.05$; vs distractor: $\Delta\bar{\beta} = 0.10, p = 0.00$; p values adjusted across hemisphere; Fig. 2A). This DMFC

Table 1. Group-level results from RSA in frontoparietal ROIs^a.

Contrast	$\bar{\beta}$	σ	t	p
DMFC (L) incon. ^b	0.11	0.03	3.65	0.00026
DMFC (R) incon. ^b	0.08	0.03	2.82	0.00485
DMFC (L) target ^b	0.04	0.01	2.61	0.00918
DMFC (R) target ^b	0.05	0.01	3.94	0.00008
DMFC (L) distr.	0.01	0.01	0.78	0.43667
DMFC (R) distr.	0.00	0.01	0.19	0.84907
DLPFC (L) incon. ^b	0.04	0.02	2.04	0.04155
DLPFC (R) incon. ^b	0.05	0.02	2.39	0.01682
DLPFC (L) target ^b	0.03	0.01	2.56	0.01062
DLPFC (R) target ^b	0.04	0.01	3.18	0.00147
DLPFC (L) distr.	−0.01	0.01	−0.98	0.32561
DLPFC (R) distr.	0.00	0.01	0.31	0.75797
LPPC (L) incon. ^b	0.07	0.02	3.07	0.00211
LPPC (R) incon. ^b	0.06	0.02	2.64	0.00818
LPPC (L) target ^b	0.03	0.01	2.74	0.00617
LPPC (R) target ^b	0.05	0.01	4.38	0.00001
LPPC (L) distr.	0.00	0.01	0.28	0.77722
LPPC (R) distr. ^c	0.02	0.01	1.86	0.06359
incon.—target DMFC (L) ^b	0.07	0.03	2.26	0.04737
incon.—target DMFC (R)	0.03	0.03	0.91	0.36047
incon.—distr. DMFC (L) ^b	0.10	0.03	3.22	0.00256
incon.—distr. DMFC (R) ^b	0.08	0.03	2.64	0.00826
target — distr. DMFC (L)	0.03	0.02	1.63	0.10230
target — distr. DMFC (R) ^b	0.05	0.02	3.18	0.00297
incon.—target DLPFC (L)	0.01	0.02	0.46	0.65809
incon.—target DLPFC (R)	0.01	0.02	0.44	0.65809
incon.—distr. DLPFC (L) ^b	0.06	0.02	2.41	0.03212
incon.—distr. DLPFC (R) ^b	0.05	0.02	2.09	0.03671
target — distr. DLPFC (L) ^b	0.04	0.02	2.76	0.01145
target — distr. DLPFC (R) ^b	0.04	0.02	2.33	0.01989
incon.—target LPPC (L)	0.04	0.03	1.62	0.21226
incon.—target LPPC (R)	0.01	0.03	0.43	0.66707
incon.—distr. LPPC (L) ^b	0.07	0.03	2.75	0.01202
incon.—distr. LPPC (R)	0.04	0.03	1.63	0.10282
target — distr. LPPC (L) ^c	0.03	0.02	1.87	0.06185
target — distr. LPPC (R) ^c	0.03	0.02	1.96	0.06185
DMFC (L) — DLPFC incon. ^b	0.06	0.02	2.38	0.01737
DMFC (L) — DLPFC target	0.00	0.01	−0.09	0.92467
DMFC (L) — LPPC incon.	0.04	0.03	1.47	0.14112
DMFC (L) — LPPC target	−0.01	0.01	−0.56	0.57868
[incon.—target] · [DMFC (L) — DLPFC] ^b	0.06	0.03	2.20	0.02797
[incon.—target] · [LPPC (L) — DLPFC]	0.03	0.03	1.06	0.29066

^aGroup-mean RSA model fits and contrasts between RSA models (target, distractor, incongruency) and between ROIs (DMFC, DLPFC, LPPC in each hemisphere). Each row displays the statistical estimates from a given RSA model and ROI, or from a contrast between RSA models or ROIs. Contrasts were conducted either between RSA models (within ROI) or between ROIs (within RSA model). For example, target — distr. | DLPFC (R) indicates the contrast between target and distractor RSA model fits (“coding strengths”) within right DLPFC. $\bar{\beta}$ indicates the mean RSA model fit over subjects. Unless noted, contrasts were averaged across hemisphere. For contrasts performed separately in each hemisphere, p values were FDR-adjusted for two comparisons. Products of bracketed terms indicate interactions between RSA model and ROI. incon., Incongruency; distr., distractor.

^bContrasts with $p < 0.05$.

^cContrasts with $0.05 < p < 0.1$.

preference was prominent enough that it could be seen in the structure of a low-dimensional embedding (Fig. 2B). Finally, incongruency coding was neuroanatomically dissociated, as this coding scheme was reflected more strongly in left DMFC than in DLPFC activity patterns ($\Delta\bar{\beta} = 0.06, p = 0.02$; Table 1).

We next focused on lateral frontoparietal regions and the coding of target information. DLPFC indeed encoded target information (left: $\bar{\beta} = 0.03, p = 0.01$ right: $\bar{\beta} = 0.04, p = 0.00$), and more strongly than distractor information ($\Delta\bar{\beta} = 0.04, p = 0.00$). We did not find a preference, however, for target over incongruency information in DLPFC ($\Delta\bar{\beta} = -0.01, p = 0.62$), nor did we

observe significantly stronger target coding than in DMFC ($\bar{\beta} = 0.00, p = 0.92$). A qualitatively similar pattern of results was observed in LPPC: significant target and incongruency coding, significantly enhanced target versus distractor coding, but no observed preference among target or incongruency nor detectable difference from DMFC target coding (Table 1).

Thus, at the group level, we observed a single rather than double dissociation between medial and lateral frontoparietal cortex, in the form of enhanced sensitivity to incongruency relative to target coding in left DMFC. In all ROIs, however, these two sources of task information were more strongly encoded than distractor information.

Sensitivity and control analyses

We next tested a series of hypotheses to scrutinize and extend our results.

First, we conducted a positive control analysis to bolster confidence in the statistical power of RSA methods within the present design. In particular, we sought to determine whether our methods could detect dissociations in task coding that are strongly expected to exist. For this, we focused on primary somatomotor and visual cortical ROIs, the responses of which can be assumed to reflect, relatively selectively, response-related (i.e., motoric) and visual form-related coding. As the distractor (word) defines the visual form of the Stroop stimulus, coding of form-related features should be captured by our distractor model. In parallel, as our analysis included only correct-response trials (i.e., in which the target response was spoken), coding of motoric features should be captured by our target model. Consistent with this logic, within early visual cortex, evidence of preferential distractor coding was observed (Fig. 2C; distractor: $\bar{\beta} = 0.04, p < 0.001$, target: $\bar{\beta} = 0.02, p = 0.17$, distractor vs target: $\Delta\bar{\beta} = 0.03, p = 0.05$), whereas, within primary ventral somatomotor cortices (encompassing the “mouth” homunculi), a relatively selective pattern of target coding was found (Fig. 2C; distractor: $\bar{\beta} = 0.00, p = 0.9$, target: $\bar{\beta} = 0.13, p < 0.001$, distractor vs target: $\Delta\bar{\beta} = -0.13, p < 0.001$). Further, representation of these dimensions were strong enough to predominate the overall structures of patterns within low-dimensional embeddings (Fig. 2B). Thus, in primary visual and somatomotor regions, a relatively clear-cut group-level double-dissociation emerged. This suggests that our models were adequately powered, at least in primary sensorimotor cortices, to detect functional distinctions, and that the failure to observe distractor coding in DMFC and DLPFC was not because of a general deficiency in our distractor coding model.

Second, we conducted sensitivity analyses to assess the robustness of our results to the particular ROI definitions used. In one analysis, we tested more expansive ROI definitions, by using alternatively-defined superparcels (Extended Data Fig. 1-1). These definitions included additional, more rostral PFC parcels (1 in DMFC, 3 in DLPFC), which begin to encroach into ventromedial PFC and frontopolar cortex (e.g., the rostral DMFC parcel was assigned to the Default-Mode network within the Cole-Anticevic divisions). Nevertheless, the previously observed dissociations were robust to these more liberal definitions (Extended Data Fig. 2-1). In the other analysis, we examined whether the overall superparcel coding profiles were representative of individual parcels. While DMFC and DLPFC results generally reflected that of constituent parcels (Extended Data Fig. 2-2A,B), interestingly, there was substantial heterogeneity within left LPPC (Extended Data Fig. 2-2C). Similar to left

DMFC, a collection of left LPPC regions spanning the length of intraparietal sulcus strongly encoded the incongruity dimension (i.e., IP1, IP2, IPS1, AIP, LIP, MIP).

Third, we examined whether the lack of observed discrimination between target and incongruity coding dimensions within DLPFC could be explained by increased error variance potentially present in fMRI activity patterns within this region. Prior work has suggested that PFC regions might be particularly susceptible to this confound (Bhandari et al., 2018). It is possible that we might have observed a dissociation in left DMFC but not DLPFC because of differential levels of statistical power across the two regions. We therefore derived an SNR analysis to determine whether this was a viable explanation (see Materials and Methods). A paired-sample bootstrap test did not indicate a systematic difference between DLPFC versus left DMFC group-level SNR ($\Delta z = -0.01, p = 0.60$). However, the SNRs in these regions were also not confirmably similar, as indicated by an equivalence test (which provides a confirmatory test of the null hypothesis; $p = 0.10$). Therefore, while we cannot rule out the possibility that the single dissociation observed between left DMFC and DLPFC was driven by better group-level SNR in DMFC, any potential SNR differences between the two regions were not substantial enough for our methods to detect.

Fourth, to provide a basis for comparison to most extant neuroimaging research of the Stroop task, we conducted a “univariate activation” analysis, examining whether these brain regions were generally more active during incongruent versus congruent conditions. No regions were found to respond more strongly overall to incongruent versus congruent conditions (Extended Data Fig. 2-3A), although the mean contrast in DMFC (L) was positive (i.e., incongruent > congruent). This null result was not surprising, however, because of the high frequency of incongruent trials within the experiment — which is known to reduce both the behavioral and neural univariate Stroop effect (Logan and Zbrodoff, 1979; Carter et al., 2000; De Pisapia and Braver, 2006). While this null result demonstrates the utility of using RSA in this case, it should not, however, be seen as direct evidence for the increased sensitivity of RSA versus univariate methods, as the univariate and RSA-based tests as implemented here are subject to different constraints and are thus incomparable (Allefeld et al., 2016). Finally, the magnitude of the univariate Stroop effects was only weakly correlated to incongruity coding model fits (Extended Data Fig. 2-3B), suggesting that these measures were nonredundant.

Finally, we tested whether these patterns of results were robust to alternative RSA techniques, including a downsampling technique to equate trial counts across runs, two “between-run” RSA methods (cross-run RSA and cross-validated RSA), and alternative parameterizations of the incongruity coding model (see Materials and Methods). Findings were robust to downsampling and to between-run RSA (Extended Data Fig. 2-4), and were highly similar across different parameterizations (Extended Data Fig. 2-5). Interestingly, however, the detection of incongruity coding in DMFC depended on whether a linear or nonlinear RSA method was used. When using an RSA method that tests linear discriminability between conditions (cross-validated RSA), the incongruity coding effect was abolished in DMFC; whereas when using a comparable method that is sensitive to nonlinear pattern differences (cross-run RSA), the effect remained quite strong. This pattern of results suggests that incongruity information was encoded nonlinearly within DMFC activation patterns. Indeed, this can be seen within the

Table 2. Parameter estimates from hierarchical brain–behavior models that explained individual differences in behavioral performance (RT) with variability in the strength of neural coding (β , see Method)^a

Contrast	b	σ	t	p
DMFC (L) incon. ^c	47.66	25.97	1.84	0.06647
DMFC (R) incon.	−0.10	33.15	0.00	0.99762
DMFC (L) target	48.40	90.63	0.53	0.59330
DMFC (R) target	−24.04	63.75	−0.38	0.70613
DLPFC (L) incon.	14.08	41.17	0.34	0.73227
DLPFC (R) incon.	38.42	33.66	1.14	0.25367
DLPFC (L) target	−105.93	78.34	−1.35	0.17637
DLPFC (R) target ^b	−239.13	71.68	−3.34	0.00085
LPPC (L) incon.	−48.95	38.13	−1.28	0.19921
LPPC (R) incon.	−32.40	32.22	−1.01	0.31467
LPPC (L) target	−1.90	93.65	−0.02	0.98385
LPPC (R) target ^b	−235.34	71.23	−3.30	0.00096
incon.−target DMFC (L)	20.92	97.23	0.22	0.82962
incon.−target DLPFC (R) ^b	299.21	80.62	3.71	0.00021
incon.−target LPPC (R) ^b	204.56	79.94	2.56	0.01050
DMFC (L) − DLPFC (R) incon.	22.85	52.92	0.43	0.66587
DMFC (L) − DLPFC (R) target ^b	500.63	139.14	3.60	0.00032
DMFC (L) − LPPC (R) incon. ^b	134.19	49.85	2.69	0.00710
DMFC (L) − LPPC (R) target ^b	468.99	137.46	3.41	0.00065

^aSeparate models were fit per region and RSA model (target, incongruity) combination. The displayed estimates correspond to the interaction $\text{stroop} \cdot \beta$, which indicates how subjects' coding strengths relate to the magnitude of their Stroop interference effects. Also displayed are contrasts on $\text{stroop} \cdot \beta$ estimates, between RSA models (within ROI) or between ROIs (within RSA model). For example, target − incon. | DLPFC (R) indicates, for right DLPFC, whether target and incongruity coding were differentially related to the size of the Stroop interference effect.

^bContrasts with $p < 0.05$.

^cContrasts with $0.05 < p < 0.1$.

two-dimensional embedding (Fig. 2), as a radial, rather than linear, separation of incongruent (central) and congruent (peripheral) stimuli.

Individual

Primary individual-level results are summarized in Figure 3. Statistical estimates corresponding to results outlined within this section are contained in Table 2 (for scatter plots of all associations, see Extended Data Fig. 3-1).

Better-performing subjects have stronger lateral frontoparietal target coding

The fidelity of target-related information in lateral frontoparietal cortex — DLPFC, in particular — is thought to be closely linked to the efficiency with which an individual resolves response conflict in tasks such as Stroop. By using subject-level target-coding estimates (β_{target}) to model behavioral performance (RT), we tested this fundamental prediction relatively directly. Indeed, subjects with stronger target coding in both right DLPFC and right LPPC resolved Stroop interference effects more quickly (DLPFC: $b = -239.13, p = 0.00, r = -0.44 [-0.66, -0.15]$, $\rho = -0.39 [-0.63, -0.08]$; LPPC: $b = -235.34, p = 0.00, r = -0.44 [-0.63, -0.17]$, $\rho = -0.35 [-0.59, -0.06]$; p values corrected across hemisphere; linear correlation, r ; rank correlation, ρ ; bootstrapped 95% CI, [lower, upper]; Fig. 3A; Table 2). Subjects' target coding estimates were moderately correlated ($r = 0.53$) between these two lateral frontoparietal regions, as expected based on their strong neuroanatomical and functional connectivity (e.g., Petrides and Pandya, 1999; Buschman and Miller, 2007). In neither of these lateral regions was incongruity coding significantly related to the Stroop interference effect (DLPFC: $b = 38.42, p = 0.25, r = -0.17 [-0.14, 0.45]$, $\rho = -0.10 [-0.19, 0.37]$; LPPC $b = -32.40, p = 0.31, r = -0.15$

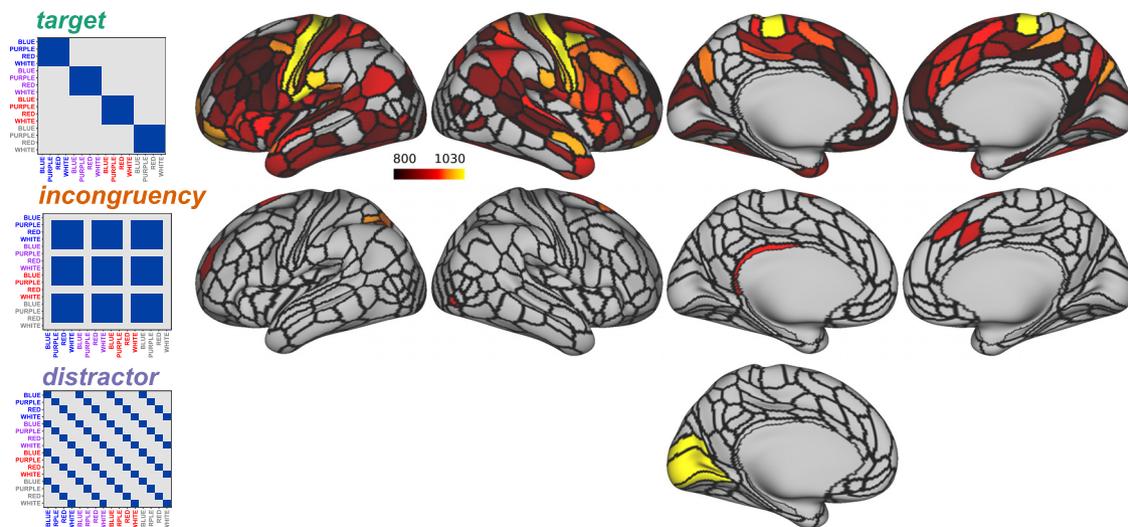


Figure 4. Cortical distributions of target (top), incongruity (middle), and distractor (bottom) coding, identified via an exploratory, whole-cortex analysis. Color bar represents test statistic from one-sided sign-ranked test. For all statistical estimates, see Extended Data Figures 4-1, 4-2, and 4-3. Sensitivity and robustness analyses suggested that the core results were robust to various analysis decisions and implementations (handling of head motion, see Results; parameterization of incongruity coding, Extended Data Figure 2-5B; RSA estimation method, Extended Data Figure 4-4).

$[-0.40, 0.13]$, $\rho = -0.18 [-0.45, 0.12]$; p values uncorrected), and within right DLPFC, target coding was a considerably stronger explanatory variable than incongruity coding ($\Delta b = -299.21$, $p = 0.00$; Table 2).

Conversely, incongruity-related responses of DMFC are thought to be positively associated with maladaptive policies of response selection. In line with this notion, subjects with stronger incongruity coding in left DMFC tended to exhibit greater Stroop interference, although this interaction was nonsignificant ($b = 47.66$, $p = 0.13$, $r = -0.26 [-0.01, 0.49]$, $\rho = -0.17 [-0.12, 0.44]$; p value corrected across hemisphere; Fig. 3A). Similarly, incongruity coding was numerically more strongly associated with Stroop interference effects than target coding in this region, although this statistical difference was also nonsignificant ($\Delta b = -20.92$, $p = 0.83$; Table 2). While weak, this DMFC–Stroop association notably emerged within the same hemisphere that displayed a group-level preference for incongruity information (Fig. 2B). Importantly, the sign of the correlations we detected within our ROIs were robust to alternative RSA techniques (Extended Data Fig. 3-2; general attenuation in effect size is expected with higher variance techniques; see Downsampling analysis and Between-run RSA).

Considering the collective pattern of results, we conclude here that our findings support the hypothesis that target coding in (right) DLPFC and in LPPC reflected a common process of implementing control.

Model selection affirms a lateral–medial dissociation and identifies unexpected relationships

Because the preceding hypothesis-driven analysis exclusively focused on a limited set of regions, important brain–behavior relationships may have been missed. A more accurate model may even omit target and incongruity coding from DLPFC, DMFC, and LPPC altogether. Consequently, we conducted a more comprehensive test to identify regions and task dimensions that could better account for Stroop performance variability across individuals.

A data-driven model selection analysis was conducted to address this question (see Materials and Methods). We defined

an expanded set of 24 cortical regions (superparcels), including the six defined and used earlier, which covered various areas that may be important for performing the Stroop task (Fig. 2B; Extended Data Fig. 3-3). Conducting RSA on each superparcel furnished three coding estimates (one per coding model) per superparcel. These 72 estimates were then used as features in a cross-validated model selection procedure.

Strikingly, the selected model contained all three hypothesized measures: (right) DLPFC and LPPC target coding, and (left) DMFC incongruity coding. In addition, two unexpected measures were identified, both with negative slopes: left DLPFC distractor coding ($r = -0.43 [-0.66, -0.14]$, $\rho = -0.33 [-0.60, -0.04]$) and left ventral visual incongruity coding ($r = -0.35 [-0.52, -0.13]$, $\rho = -0.43 [-0.62, -0.20]$; bivariate correlations shown in Fig. 3C). We tested the predictive accuracy of the selected model by using a held-out validation set of 17 subjects (co-twins of the primary analysis set). The selected model explained significant variance in Stroop interference effects within the validation set (Fig. 3D). While not fully independent, these validation set data were from a true hold-out and obtained from distinct individuals. This result therefore bolsters claims regarding the stability of the selected model.

Nevertheless, to provide a cursory test of a truly independent validation set (i.e., with no familial dependency to the training set), we excluded all subjects in the training set who were co-twins of those in the validation set, and reconducted this model selection procedure. This amounted to discarding 16/49 (33%) of training-set observations. The selected model contained only one variable, which was not in our ROIs (but in early visual cortex), and was unable to predict held-out Stroop effects ($r = -0.12$). This is perhaps unsurprising, however, given the substantial reduction in the size of the training dataset for an already high-dimensional model. To reduce the dimensionality, we reconducted this analysis, focusing now instead only on ROIs and coding schemes of interest — target coding in DLPFC and LPPC, and incongruity coding in DMFC (within each hemisphere) — and additionally ensured that all variables were used in prediction (via ridge regression). This model was better able to predict the held-out Stroop effect ($r = 0.20$), in particular,

relative to a comparable model that contained theoretically “mismatched” ROI×coding scheme combinations (incongruency coding in DLPFC and LPPC, target coding in DMFC; $r = -0.17$; bootstrapped $p = 0.088$).

Results from these model selection analyses affirm a functional dissociation across the medial–lateral axis of frontoparietal cortex, and further demonstrate that Stroop-task representations within DLPFC and DMFC hold relatively privileged relationships with behavior

Exploratory whole-cortex RSA

In a final, exploratory analysis, we estimated RSA models separately for each MMP parcel, to determine more comprehensively how target, distractor, and incongruency coding are distributed across cortex. These three task dimensions were encoded across cortex according to different neuroanatomical profiles. Target coding was widespread (observed in 207/360 parcels), covering substantial portions of the frontal and temporal lobes, including many perisylvian regions (Fig. 4, top; Extended Data Fig. 4-1). Notably, the strongest target coding was observed within regions that receive strong sensory and (or) motor-related input (Extended Data Fig. 4-1). Contrastingly, incongruency coding was detected predominantly within prefrontal and intraparietal sulcal parcels — including DMFC, but also left LPPC, bilateral superior frontal gyrus, and left lateral frontopolar cortex (rostral DLPFC) — but additionally within left retrosplenial and right lateral occipital cortex (Fig. 4, middle; Extended Data Fig. 4-3). Aside from this occipital area, these incongruency-coding parcels notably belonged to control networks (frontoparietal, cinguloopercular, and dorsal attention; Extended Data Fig. 4-3). In a third, distinct profile, distractor coding was only observed within early visual cortex (left V1 and V2; Fig. 4, bottom; Extended Data Fig. 4-2).

As a negative control analysis, we tested whether we could decode these three task variables (target, distractor, incongruency) from framewise head motion estimates, using the same RSA procedures as above. No task variable was significantly encoded within patterns of head movements ($bs < 0.01$, $ps > 0.36$). Using a larger basis set of motion estimates (12 vs 6) did not yield significant decoding ($bs < 0.01$, $ps > 0.17$), nor any increased sensitivity to task variables ($bs < 0.00$, $ps > 0.17$), suggesting that our motion removal procedures (scrubbing and 6 motion regressors) were adequate.

Finally, we repeated this exploratory analysis using alternative RSA methods involving downsampling and between-run RSA (see *Downsampling analysis*, and *Between-run RSA*; Extended Data Fig. 4-4). Across these analyses, the core results were quite robust: we found incongruency coding in parcels within DMFC (though this depended on the use of nonlinear RSA methods, as in Extended Data Figs. 2-4, 3-2), target coding in mid-DLPFC, and distractor coding in visual cortex. Our findings were therefore not specific to a particular estimation method.

Collectively, these exploratory results confirm and extend our prior findings. (1) As with the reported brain–behavior associations, target coding was emphasized relative to coding of other task dimensions. (2) Yet despite this emphasis, important and expected dissociations in anatomic profiles were identified across our three coding models, further suggesting that these models were successful in measuring coding of distinct task dimensions (Fig. 1A).

Discussion

We analyzed the similarity structure, or representational geometry (Kriegeskorte and Kievit, 2013), of frontoparietal activity

patterns associated with cognitive control, during performance of the classic color-word Stroop task. In left DMFC, incongruency coding predominated. While DLPFC and LPPC encoded both target and incongruency-related information, distractor coding was not detected in these regions but was instead identified in early visual cortex. Further, these neural coding estimates were important and specific indicators of individual differences in magnitude of the behavioral Stroop interference effect. Individuals with stronger target coding in right DLPFC and right LPPC, but weaker incongruency coding in left DMFC, exhibited enhanced cognitive control, in terms of a reduced Stroop effect. Further, in a more comprehensive predictive model that included coding measures from a wide set of cortical regions, coding measures specifically from lateral frontoparietal and dorsomedial frontal regions were privileged in their link to behavior.

On one level, this study is a specific extension of research that has drawn dissociations between control-related functions of DLPFC and DMFC (MacDonald et al., 2000; Floden et al., 2011). Most prominently, MacDonald et al. (2000) used a modified Stroop-task design, in which task rules (delivered via precues) randomly alternated between color naming and word reading across trials, to demonstrate that DLPFC and DMFC encoded different types of information during cognitive control engagement. In particular, DLPFC was selectively recruited following cues for the more demanding color-naming task, whereas DMFC was instead driven by incongruent color-naming trials. This pattern of recruitment suggested that DLPFC encodes task-set and rule-related information in a preparatory manner, whereas DMFC encodes incongruency and conflict-related information in a stimulus-evoked manner. In the current study, we leveraged the high spatial dimensionality of fMRI to test whether this functional dissociation can be observed within a common time window of response selection, and further with a more traditional Stroop-task design, which does not involve task cues or switches. Our findings reinforce the conclusions of these relatively low-powered studies ($N = 12$ in Floden et al., 2011; $N = 9$ in MacDonald et al., 2000) and indicate that the dissociations were not dependent on the use of cued task-switching designs. Synthesizing these prior findings with those of the present study hints at a continuity in the putative role of DLPFC during target selection. Rather than exclusively contributing to preparation, DLPFC coding may evolve from proactively representing abstract rule or set-related information, toward more concrete targets and behavioral choices as relevant stimulus information becomes available in the environment. This view accords with work in monkey neurophysiology (Mante et al., 2013; Rigotti et al., 2013; Stokes et al., 2013), yet further work is needed to determine whether similar dynamics occur within human DLPFC and how such dynamics may reflect or interact with specific processes of cognitive control in Stroop-like tasks.

More broadly, the results of this study highlight the utility of RSA and the general representational geometric framework for investigating cognitive control. Previous work has used MVPA decoding in the Stroop task to study the impact of control demand on posterior representations (e.g., Banich et al., 2019). Here, we used the RSA framework to explicitly model and decompose control-related frontoparietal representations. Indeed, a major motivation of the current study was to assess how well RSA measures of frontoparietal coding map to theorized mechanisms of control. For this purpose, the medial–lateral functional dissociation in frontoparietal cortex was a useful test-bed, as it features in several theoretical accounts (Botvinick et al., 2001; Miller and Cohen, 2001; Ridderinkhof et al., 2004; Shenhav et al., 2013). Our results were

generally in line with these accounts, joining with a growing body of research in suggesting that RSA provides a convenient yet powerful framework from which neural measures can be used to test cognitive control theory (for review, see Freund et al., 2021).

Nevertheless, the current work represents only an initial step in using the RSA framework to investigate cognitive control within Stroop-like tasks. As such, our study raises a number of unaddressed questions. But promisingly, there are ample opportunities for improving and extending the RSA framework highlighted here. For instance, a key limitation of the current study was the finding of widespread coding of the target dimension, suggesting a lack of specificity in the target RSA model. This is perhaps not surprising, however, as the model would capture not only coding of attentional-template and choice-related information, but also hue and response-related information. We mitigated this issue by demonstrating that target coding was selectively related to behavior within DLPFC and LPPC. Yet, this limitation could be addressed more powerfully by experimental design. Adding specific factorial manipulations, such as a task rule manipulation (see MacDonald et al., 2000; Hall-McMaster et al., 2019) or a response modality manipulation (see Minxha et al., 2020; see also Barch et al., 2001), would enable a richer, more precise set of cognitive control-relevant coding variables to be estimated (Fig. 5).

Future work could also address some of the complexities revealed by our data that were not entirely accounted for by the theoretical frameworks we used. For one, although predicted coding profiles emerged in some frontoparietal ROIs, all regions encoded incongruency and target information. This incongruency-coding finding is consistent with prior univariate fMRI research (e.g., Nee et al., 2007; Niendam et al., 2012) and a more recent finding that the responses of single neurons in human dACC and DLPFC are robustly modulated by conflict (Smith et al., 2019). With respect to target coding, however, one speculative interpretation is that, during the relatively late phase of response selection and execution, control networks may lose modular structure as the circuitry collectively converges on a behavioral choice. This interpretation accords with the fact that “choice axes” are encoded within multiple key nodes of frontoparietal decision circuitry, for example: in macaque LIP (Roitman and Shadlen, 2002), in macaque caudal DLPFC (Mante et al., 2013), and in human dACC and pre-SMA (Minxha et al., 2020; see also Okazawa et al., 2021). This account could be addressed using the enriched experimental design described above, identifying when and where choice coding is emphasized over the course of a trial.

Another unexpected finding was a relatively robust negative relationship between the strength of incongruency coding in left ventral visual cortex and the magnitude of the Stroop effect (Fig. 3B, left). One interpretation of this finding is provided by the biased competition framework, as an effect of selective visual attention. Prior work has demonstrated that certain ventral visual

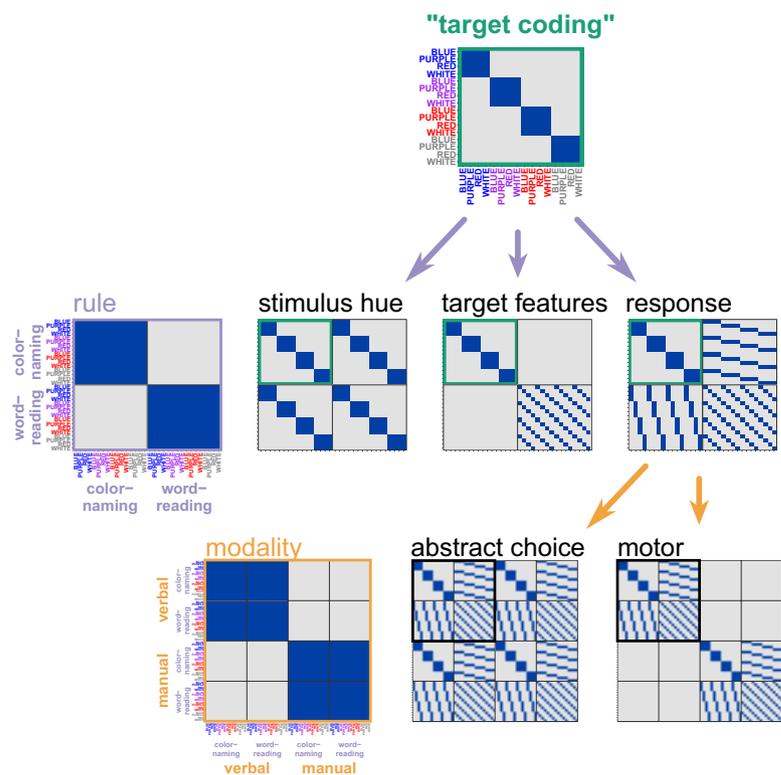


Figure 5. Expanded experimental designs. The current target coding model (top, green text; see also Fig. 1) conflates distinct coding schemes, including those associated with relatively “early” vision, relatively “late” motoric commands, or more “central” rule-dependent schemes. By adding a factorial manipulation of task rule (i.e., participants respond to the present set of stimuli, but under both color-naming and word-reading conditions; MacDonald et al., 2000; light purple arrows and text), the target coding model can be expanded into three more precise models (middle row of matrices), which are identical to our “target coding” model within the top-left quadrant (green square; naming–naming), but are now distinguished among the other quadrants (reading–reading, reading–naming). “Stimulus hue” and “response” models would identify coding that depends on features of either the stimulus or response, independent of the task rules. Conversely, “target features” would identify a more flexible “attentional template” coding scheme, which depends critically on the current task goals (bottom right quadrant resembles our “distractor” coding model in Fig. 1, as word features become the target in the word-reading condition). Importantly, too, abstract coding of task rule would also now be distinguishable (“rule” matrix on left). For an elegant example of a similar design within cued-task switching, see Hall-McMaster et al. (2019). The “response” model could be further elaborated by incorporating another manipulation of response modality (verbal, manual; orange arrows). Now, coding of abstract choice options, independent of effector, could be separated from effector-dependent motoric coding. For an elegant example of a similar manipulation, see Minxha et al. (2020).

regions, those that are strongly tuned to target features, activate as a function of Stroop incongruency (Egner and Hirsch, 2005). In our task, mid-ventral stream areas may have received biasing input, selectively on incongruent trials, which enhanced stimulus-related coding and communication with downstream regions. Using the expanded RSA design sketched above, we might expect that such an effect would be limited to color naming conditions, when selective attention processes would be most prominent.

Perhaps the most surprising result was the robust negative correlation observed between left DLPFC distractor coding and the behavioral Stroop effect (Fig. 3B, right). At face, accounting for this finding within the framework of top-down biased competition is difficult. But, given the statistics of our task in which incongruent trials were frequent and congruent were rare, distractor information could have been used to facilitate performance. An association between distractor features and incongruency could have been learned and used to influence response selection, for example, by retrieving and implementing a stimulus-appropriate attentional setting (Melara and Algom,

2003; Bugg and Crump, 2012). Indeed, subjects clearly do exploit these associations, as indicated by reduced Stroop effects for stimuli that are “mostly incongruent,” also known as “item-specific proportion congruency” effects (ISPC; Bugg et al., 2011; Jiang et al., 2015; see also Crump and Milliken, 2009). The prediction that distractor coding might reflect ISPC effects could be tested by varying ISPC levels across different stimuli. For stimuli in which the specific color or word is not predictive of congruency, the relationship between DLPFC distractor coding and improved Stroop performance should not be present.

Finally, the present study sets the stage for using RSA to test the dual-mechanisms framework of cognitive control (Braver, 2012). This framework explains much within- and between-individual variability in cognitive control function by the existence of two operational “modes” of cognitive control: proactive and reactive. These modes are proposed to have dissociable signature neural coding schemes. Proactive control should rely heavily on goal-relevant coding schemes that originate in LPFC before target-stimulus onset as abstract rule or context coding, but which may morph into target coding after stimulus onset. In contrast, reactive control should rely on an incongruity-based coding scheme (including coding of whichever task dimensions are predictive of incongruency), originating post-target onset, with potential loci in DMFC or subcortical structures (Jiang et al., 2015; Chiu et al., 2017). As suggested here, it may be possible to measure correlates of these neural coding schemes via RSA. Experimental factors that encourage subjects to adopt one mode over another (e.g., strategy training, expectancy of difficulty) should correspondingly shift frontoparietal coding schemes along these proactive and reactive dimensions. Further, their behavioral relevance should predictably change, as well: for example, in task contexts in which a proactive control mode is theoretically maladaptive, subjects with stronger proactive coding should perform worse. Thus, the dual-mechanisms framework suggests a broad range of hypotheses amenable to testing with RSA methods (see, e.g., Hall-McMaster et al., 2019). Such hypotheses can be addressed in the broader Dual Mechanisms of Cognitive Control dataset (Braver et al., 2020), of which the data used here are a small subset.

References

- Alink A, Walthers A, Krugliak A, Bosch JJ, van den, Kriegeskorte N (2015) Mind the drift: improving sensitivity to fMRI pattern information by accounting for temporal pattern drift. *bioRxiv* 032391.
- Allefeld C, Gørgen K, Haynes JD (2016) Valid population inference for information-based imaging: from the second-level *t*-test to prevalence inference. *Neuroimage* 141:378–392.
- Assem M, Glasser MF, Van Essen DC, Duncan J (2020) A domain-general cognitive core defined in multimodally parcellated human cortex. *Cereb Cortex* 30:4361–4380.
- Aust F, Barth M (2020) papaja: create APA manuscripts with R Markdown. <https://github.com/crsh/papaja>.
- Baayen RH, Milin P (2010) Analyzing reaction times. *Int J Psychol Res* 3:12–28.
- Banich MT, Smolker HR, Snyder HR, Lewis-Peacock JA, Godinez DA, Wager TD, Hankin BL (2019) Turning down the heat: neural mechanisms of cognitive control for inhibiting task-irrelevant emotional information during adolescence. *Neuropsychologia* 125:93–108.
- Barch DM, Braver TS, Akbudak E, Conturo T, Ollinger J, Snyder A (2001) Anterior cingulate cortex and response conflict: effects of response modality and processing domain. *Cereb Cortex* 11:837–848.
- Barr DJ, Levy R, Scheepers C, Tily HJ (2013) Random effects structure for confirmatory hypothesis testing: keep it maximal. *J Mem Lang* 68:255–278.
- Bates D, Maechler M, Bolker B, Walker S (2014) lme4: linear mixed-effects models using Eigen and S4. R Package Version 1:1–23.
- Bhandari A, Gagne C, Badre D (2018) Just above chance: is it harder to decode information from prefrontal cortex hemodynamic activity patterns? *J Cogn Neurosci* 30:1473–1498.
- Botvinick MM, Braver TS, Barch DM, Carter CS, Cohen JD (2001) Conflict monitoring and cognitive control. *Psychol Rev* 108:624–652.
- Braver TS (2012) The variable nature of cognitive control: a dual mechanisms framework. *Trends Cogn Sci* 16:106–113.
- Braver TS, Kizhner A, Tang R, Freund MC, Etzel JA (2020) The Dual Mechanisms of Cognitive Control (DMCC) Project. *bioRxiv*.
- Bugg JM (2014) Conflict-triggered top-down control: default mode, last resort, or no such thing? *J Exp Psychol Learn Mem Cogn* 40:567–587.
- Bugg JM, Crump MJ (2012) In support of a distinction between voluntary and stimulus-driven control: a review of the literature on proportion congruent effects. *Front Psychol* 3:367.
- Bugg JM, Jacoby LL, Chanani S (2011) Why it is too early to lose control in accounts of item-specific proportion congruency effects. *J Exp Psychol Hum Percept Perform* 37:844–859.
- Buschman TJ, Miller EK (2007) Top-down versus bottom-up control of attention in the prefrontal and posterior parietal cortices. *Science* 315:1860–1862.
- Cai MB, Schuck NW, Pillow JW, Niv Y (2019) Representational structure or task structure? Bias in neural representational similarity analysis and a Bayesian method for reducing bias. *PLoS Comput Biol* 15:e1006299.
- Carter CS, MacDonald AM, Botvinick M, Ross LL, Stenger VA, Noll D, Cohen JD (2000) Parsing executive processes: strategic vs. evaluative functions of the anterior cingulate cortex. *Proc Natl Acad Sci USA* 97:1944–1948.
- Chiu YC, Jiang J, Egner T (2017) The caudate nucleus mediates learning of stimulus–control state associations. *J Neurosci* 37:1028–1038.
- Cox DD, Savoy RL (2003) Functional magnetic resonance imaging (fMRI) ‘brain reading’: detecting and classifying distributed patterns of fMRI activity in human visual cortex. *Neuroimage* 19:261–270.
- Cox RW (1996) AFNI: software for Analysis and Visualization of Functional Magnetic Resonance Neuroimages. *Comput Biomed Res* 29:162–173.
- Crump MJ, Milliken B (2009) The flexibility of context-specific control: evidence for context-driven generalization of item-specific control settings. *Q J Exp Psychol (Hove)* 62:1523–1532.
- De Pisapia N, Braver TS (2006) A model of dual control mechanisms through anterior cingulate and prefrontal cortex interactions. *Neurocomputing* 69:1322–1326.
- Diedrichsen J (2006) A spatially unbiased atlas template of the human cerebellum. *Neuroimage* 33:127–138.
- Diedrichsen J, Berlot E, Mur M, Schütt HH, Kriegeskorte N (2020) Comparing representational geometries using the unbiased distance correlation. *arXiv* 2007.02789.
- Edelman S, Grill-Spector K, Kushnir T, Malach R (1998) Toward direct visualization of the internal shape representation space by fMRI. *Psychobiology* 26:309–231.
- Egner T, Hirsch J (2005) Cognitive control mechanisms resolve conflict through cortical amplification of task-relevant information. *Nat Neurosci* 8:1784–1790.
- Esteban O, Markiewicz CJ, Blair RW, Moodie CA, Isik AI, Erramuzpe A, Kent JD, Goncalves M, DuPre E, Snyder M, Oya H, Ghosh SS, Wright J, Durnez J, Poldrack RA, Gorgolewski KJ (2019) fMRIPrep: a robust preprocessing pipeline for functional MRI. *Nat Methods* 16:111–116.
- Esteban O, Ciric R, Finc K, Blair RW, Markiewicz CJ, Moodie CA, Gorgolewski KJ (2020) Analysis of task-based functional MRI data preprocessed with fMRIPrep. *NatProtoc* 15:2186–2202.
- Etzel JA, Zacks JM, Braver TS (2013) Searchlight analysis: promise, pitfalls, and potential. *Neuroimage* 78:261–269.
- Floden D, Vallesi A, Stuss DT (2011) Task context and frontal lobe activation in the Stroop task. *J Cogn Neurosci* 23:867–879.
- Freund MC, Etzel JA, Braver T (2021) Neural coding of cognitive control: the representational similarity analysis approach. *Trends Cogn Sci* 25:622–638.
- Friedman J, Hastie T, Tibshirani R (2010) Regularization paths for generalized linear models via coordinate descent. *J Stat Softw* 33:1–22.
- Glasser MF, Sotiropoulos SN, Wilson JA, Coalson TS, Fischl B, Andersson JL, Xu J, Jbabdi S, Webster M, Polimeni JR, Van Essen DC, Jenkinson M, WU-Minn HCP Consortium (2013) The minimal preprocessing pipelines for the Human Connectome Project. *Neuroimage* 80:105–124.

- Glasser MF, Coalson TS, Robinson EC, Hacker CD, Harwell J, Yacoub E, Ugurbil K, Andersson J, Beckmann CF, Jenkinson M, Smith SM, Van Essen DC (2016) A multi-modal parcellation of human cerebral cortex. *Nature* 536:171–178.
- Gonthier C, Braver TS, Bugg JM (2016) Dissociating proactive and reactive control in the Stroop task. *Mem Cognit* 44:778–788.
- Gordon EM, Laumann TO, Adeyemo B, Huckins JF, Kelley WM, Petersen SE (2016) Generation and evaluation of a cortical area parcellation from resting-state correlations. *Cereb Cortex* 26:288–303.
- Haines N, Kvam PD, Irving LH, Smith C, Beauchaine TP, Pitt MA, Turner B (2020) Learning from the reliability paradox: how theoretically informed generative models can advance the social, behavioral, and brain sciences. *PsyArXiv*.
- Hall-McMaster S, Muhle-Karbe PS, Myers NE, Stokes MG (2019) Reward boosts neural coding of task rules to optimize cognitive flexibility. *J Neurosci* 39:8549–8561.
- Hastie T, Tibshirani R, Friedman J (2009) *The elements of statistical learning: data mining, inference, and prediction*, Ed 2. New York: Springer.
- Haxby JV, Gobbini MI, Furey ML, Ishai A (2001) Distributed and overlapping representations of faces and objects in ventral temporal cortex. *Science* 293:2425–2430.
- Hothorn T, Bretz F, Westfall P (2008) Simultaneous inference in general parametric models. *Biomed J* 50:346–363.
- Ji JL, Spronk M, Kulkarni K, Repovš G, Anticevic A, Cole MW (2019) Mapping the human brain's cortical-subcortical functional network organization. *Neuroimage* 185:35–57.
- Jiang J, Brashier NM, Egnor T (2015) Memory meets control in hippocampal and striatal binding of stimuli, responses, and attentional control states. *J Neurosci* 35:14885–14895.
- Kane MJ, Engle RW (2002) The role of prefrontal cortex in working-memory capacity, executive attention, and general fluid intelligence: an individual-differences perspective. *Psychon Bull Rev* 9:637–671.
- Kriegeskorte N, Kievit RA (2013) Representational geometry: integrating cognition, computation, and the brain. *Trends Cogn Sci* 17:401–412.
- Kriegeskorte N, Mur M, Bandettini P (2008) Representational similarity analysis: connecting the branches of systems neuroscience. *Front Syst Neurosci* 2:4–28.
- Kruskal JB (1964) Multidimensional scaling by optimizing goodness of fit to a nonmetric hypothesis. *Psychometrika* 29:1–27.
- Kurtzer GM, Sochat V, Bauer MW (2017) Singularity: scientific containers for mobility of compute. *PLoS One* 12:e0177459.
- Laird NM, Ware JH (1982) Random-effects models for longitudinal data. *Biometrics* 38:963–974.
- Lakens D (2017) Equivalence tests. *Soc Psychol Personal Sci* 8:355–362.
- Logan GD, Zbrodoff NJ (1979) When it helps to be misled: facilitative effects of increasing the frequency of conflicting stimuli in a Stroop-like task. *Mem Cogn* 7:166–174.
- MacDonald AW, Cohen JD, Stenger VA, Carter CS (2000) Dissociating the role of the dorsolateral prefrontal and anterior cingulate cortex in cognitive control. *Science* 288:1835–1838.
- MacLeod CM (1991) Half a century of research on the Stroop effect: an integrative review. *Psychol Bull* 109:163–203.
- Mante V, Sussillo D, Shenoy KV, Newsome WT (2013) Context-dependent computation by recurrent dynamics in prefrontal cortex. *Nature* 503:78–84.
- Marcus DS, Harwell J, Olsen T, Hodge M, Glasser MF, Prior F, Jenkinson M, Laumann T, Curtiss SW, Van Essen DC (2011) Informatics and data mining tools and strategies for the human connectome project. *Front Neuroinform* 5:4.
- Melara RD, Algom D (2003) Driven by information: a tectonic theory of Stroop effects. *Psychol Rev* 110:422–471.
- Mensh B, Kording K (2017) Ten simple rules for structuring papers. *PLoS Comput Biol* 13:e1005619.
- Miller EK, Cohen JD (2001) An integrative theory of prefrontal cortex function. *Annu Rev Neurosci* 24:167–202.
- Minxha J, Adolphs R, Fusi S, Mamelak AN, Rutishauser U (2020) Flexible recruitment of memory-based choice representations by the human medial frontal cortex. *Science* 368:eaba3313.
- Mumford JA, Davis T, Poldrack RA (2014) The impact of study design on pattern estimation for single-trial multivariate pattern analysis. *Neuroimage* 103:130–138.
- Mumford JA, Turner BO, Ashby FG, Poldrack RA (2012) Deconvolving BOLD activation in event-related designs for multivoxel pattern classification analyses. *Neuroimage* 59:2636–2643.
- Nee DE, Wager TD, Jonides J (2007) Interference resolution: insights from a meta-analysis of neuroimaging tasks. *Cogn Affect Behav Neurosci* 7:1–17.
- Niendam TA, Laird AR, Ray KL, Dean YM, Glahn DC, Carter CS (2012) Meta-analytic evidence for a superordinate cognitive control network subserving diverse executive functions. *Cogn Affect Behav Neurosci* 12:241–268.
- Nili H, Wingfield C, Walther A, Su L, Marslen-Wilson W, Kriegeskorte N (2014) A toolbox for representational similarity analysis. *PLoS Comput Biol* 10:e1003553.
- Okazawa G, Hatch CE, Mancoo A, Machens CK, Kiani R (2021) The geometry of the representation of decision variable and stimulus difficulty in the parietal cortex. *bioRxiv* 2021–2001.
- Oksanen J, Blanchet FG, Friendly M, Kindt R, Legendre P, McGlenn D, Wagner H (2019) *Vegan: community ecology package*. <https://CRAN.R-project.org/package=vegan>.
- Petersen SE, Dubis JW (2012) The mixed block/event-related design. *Neuroimage* 62:1177–1184.
- Petrides M, Pandya D (1999) Dorsolateral prefrontal cortex: comparative cytoarchitectonic analysis in the human and the macaque brain and cortico-cortical connection patterns. *Eur J Neurosci* 11:1011–1036.
- Pinheiro J, Bates D, DebRoy S, Sarkar D, R Core Team (2019) *nlme: Linear and nonlinear mixed effects models*. <https://CRAN.R-project.org/package=nlme>.
- Posner MI, Snyder CRR (1975) Attention and cognitive control. In: *Information processing and cognition*, Solso R (ed), pp 669–682. Chicago: Loyola Symposium.
- R Core Team (2019) *R: a language and environment for statistical computing*. Vienna: R Foundation for Statistical Computing.
- Ridderinkhof KR, Ullsperger M, Crone EA, Nieuwenhuis S (2004) The role of the medial frontal cortex in cognitive control. *Science* 306:443–447.
- Rigotti M, Barak O, Warden MR, Wang XJ, Daw ND, Miller EK, Fusi S (2013) The importance of mixed selectivity in complex cognitive tasks. *Nature* 497:585–590.
- Roitman JD, Shadlen MN (2002) Response of neurons in the lateral intraparietal area during a combined visual discrimination reaction time task. *J Neurosci* 22:9475–9489.
- Schuirman DJ (1987) A comparison of the two one-sided tests procedure and the power approach for assessing the equivalence of average bioavailability. *J Pharmacokinetic Biopharm* 15:657–680.
- Shenhav A, Botvinick MM, Cohen JD (2013) The expected value of control: an integrative theory of anterior cingulate cortex function. *Neuron* 79:217–240.
- Simmons JP, Nelson LD, Simonsohn U (2011) False-positive psychology: undisclosed flexibility in data collection and analysis allows presenting anything as significant. *Psychol Sci* 22:1359–1366.
- Smith EH, Horga G, Yates MJ, Mikell CB, Banks GP, Pathak YJ, Schevon CA, McKhann GM, Hayden BY, Botvinick MM, Sheth SA (2019) Widespread temporal coding of cognitive control in the human prefrontal cortex. *Nat Neurosci* 22:1883–1891.
- Stokes MG, Kusunoki M, Sigala N, Nili H, Gaffan D, Duncan J (2013) Dynamic coding for cognitive control in prefrontal cortex. *Neuron* 78:364–375.
- Stroop JR (1935) Studies of interference in serial verbal reactions. *J Exp Psychol* 18:643–662.
- Twomey T, Kawabata Duncan KJ, Price CJ, Devlin JT (2011) Top-down modulation of ventral occipito-temporal responses during visual word recognition. *Neuroimage* 55:1242–1251.
- Walther A, Nili H, Ejaz N, Alink A, Kriegeskorte N, Diedrichsen J (2016) Reliability of dissimilarity measures for multi-voxel pattern analysis. *Neuroimage* 137:188–200.
- Xie Y (2015) *Dynamic documents with R and knitr*, Ed 2. Boca Raton, FL: Chapman; Hall/CRC.