# Constructing neural network models from brain data reveals representational transformations underlying adaptive behavior

Takuya Ito[1,2], Guangyu Robert Yang[3], Patryk Laurent[4], Douglas H. Schultz[5], Michael W. Cole[1]

[1]Center for Molecular and Behavioral Neuroscience, Rutgers University, Newark, NJ
[2]Behavioral and Neural Sciences PhD Program, Rutgers University, Newark, NJ
[3]Center for Theoretical Neuroscience, Columbia University, New York, NY
[4]Independent Researcher
[5]Center for Brain, Biology and Behavior, University of Nebraska-Lincoln, Lincoln, NE

Contact: taku.ito1@gmail.com

# Abstract

The human ability to adaptively implement a wide variety of tasks is thought to emerge from the dynamic transformation of cognitive information. We hypothesized that these transformations are implemented via conjunctive representations in *conjunction hubs* – brain regions that selectively integrate sensory, cognitive, and motor representations. We used recent advances in using functional connectivity to map the flow of activity between brain regions to construct a task-performing neural network model from fMRI data during a cognitive control task. We verified the importance of conjunction hubs in cognitive computations by simulating neural activity flow over this empirically-estimated functional connectivity model. These simulations produced above-chance task performance (motor responses) by integrating sensory and task rule information in conjunction hubs. These findings reveal the role of conjunction hubs in supporting flexible cognitive computations, while demonstrating the feasibility of using empirically-estimated neural network models to gain insight into cognitive computations in the human brain.

# Introduction

The human brain exhibits remarkable cognitive flexibility. This cognitive flexibility enables humans to perform a wide variety of cognitive tasks, ranging from simple visual discrimination and motor control tasks, to highly complex context-dependent tasks. Key to this cognitive flexibility is the ability to use cognitive control, which involves goal-directed implementation of task rules to specify cognitive and motor responses to stimuli[1–3]. Previous studies have investigated how task-relevant sensory, motor, and rule features are represented in the brain, finding that sensory stimulus features are represented in sensory cortices[4,5], motor action features are represented in motor cortices[6], while task rule features are represented in prefrontal and other association cortices[3,7–10]. However, exactly how and where in the brain different task representations mix to convert incoming stimuli to motor responses remains unclear[11]. In contrast, artificial neural network models (ANNs) can provide computationally rigorous accounts of how different task representations mix to implement cognitive computations[12,13]. Inspired by the formalization of ANNs, we constructed an empirically-estimated neural network (ENN) model from task fMRI data to provide insight into the representational transformations in the brain during a cognitive control task.

The Flexible Hub theory provides a network account of how large-scale cognitive control networks implement flexible cognition by updating task rule representations[14,15]. While Flexible Hub theory primarily focuses on the importance of flexible rule updating for complex task performance, it does not specify how rules interact with incoming sensory stimulus activity. However, Flexible Hub theory was built upon the Guided Activation theory of prefrontal cortex – a seminal theory of the neural mechanisms underlying cognitive control – which posits that successful performance of a cognitive control task requires the selective mixing of task context with sensory stimulus activity[3]. The selective mixing of task context and sensory stimulus encoding activations would produce conjunctive (conditional association) activations that implement task rules on sensory stimuli. These conjunctive activations are thought to form through inter-area guided activations in "hidden units" located somewhere in association cortex, which we term *conjunction hubs* (Fig. 1a). The outputs of conjunction hubs then produce motor activations to produce task-appropriate behavior. Thus, by leveraging the notion in Guided Activation theory of interacting rule- and stimulus-guided neural activations (i.e., conjunctions), we built upon Flexible Hub theory to provide a brain implementation for flexible task control.

We recently developed a method – activity flow mapping – that provides a framework for testing Guided Activation theory with empirical brain data[16]. Activity flow mapping involves several steps. First, a network model is derived from empirically-estimated connectivity weights. Second, empirical task activations (e.g., activity patterns from sensory regions) are used as inputs to simulate the activity flow (i.e., propagating activity) within the brain network model. Finally, the predictions generated by simulated activity flow are tested against independent empirical brain activations for model validation. Here we used activity flow mapping to test whether putative conjunction hubs could implement the context-dependent transformations of task-rule and stimulus activations necessary to produce accurate behavioral (motor) activations in a 64-context cognitive paradigm.

2

We sought a principled approach to identify brain areas that form the conjunctive activations to produce flexible behavior. Recent studies have successfully used ANNs to probe the emergence of representational transformations in cognitive tasks[12,13,17]. Importantly, the representational geometry of ANNs has often converged with the geometry of neural representations[18–20], suggestive of the utility of ANNs in investigating task representations in the brain. Inspired by these previous studies, we constructed a feedforward ANN to investigate how conjunctive representations emerged through the transformation of task context and stimulus input activations during the 64-context cognitive paradigm. After identifying the representational geometry of task-context and stimulus conjunctions in the ANN, we identified brain regions – conjunction hubs – with similar conjunctive representations in fMRI data. The identification of brain regions selective for task rules, sensory stimuli, motor responses, and conjunctions, made it possible to construct an ENN and empirically test Guided Activation theory with activity flow mapping over data-constrained functional connections. We found that behavioral activations (in motor cortices) could be predicted through the formation of conjunctive activations through activity flow guided by task rule and sensory stimulus activations.

To summarize, we empirically tested Guided Activation theory by constructing a task-performing ENN during a 64-task cognitive paradigm. This ENN was constructed directly from fMRI data, and illustrated the importance of conjunction hubs in facilitating representational transformations. This contrasts with many possible alternative hypotheses, such as the possibility that representations are transformed directly from task input areas (e.g., sensory systems) to motor cortices, bypassing association areas. We first identified brain areas selective to different task components, namely task rules, sensory stimuli, motor responses, and conjunctions. These areas formed the spatial areas/layers of the ENN, which are conceptually similar to layers in a feedforward ANN. Next, in contrast to ANNs, which often use supervised learning to estimate connectivity weights between layers, we show that activations in ENNs can be transformed via activity flow over functional connectivity (FC) weights estimated from resting-state fMRI (Fig. 1d). This resulted in a task-performing, ENN model that transforms stimulus and task-rule fMRI activations into response activations in motor cortex during a flexible cognitive control task. Critically, the transformations implemented by the ENN were carried out without classic optimization approaches such as gradient learning, demonstrating that the intrinsic architecture of the resting brain is suitable for implementing representational transformations. Together, these findings illustrate the computational relevance of functional network organization and the importance of conjunctive representations in supporting flexible cognitive computations in the human brain.
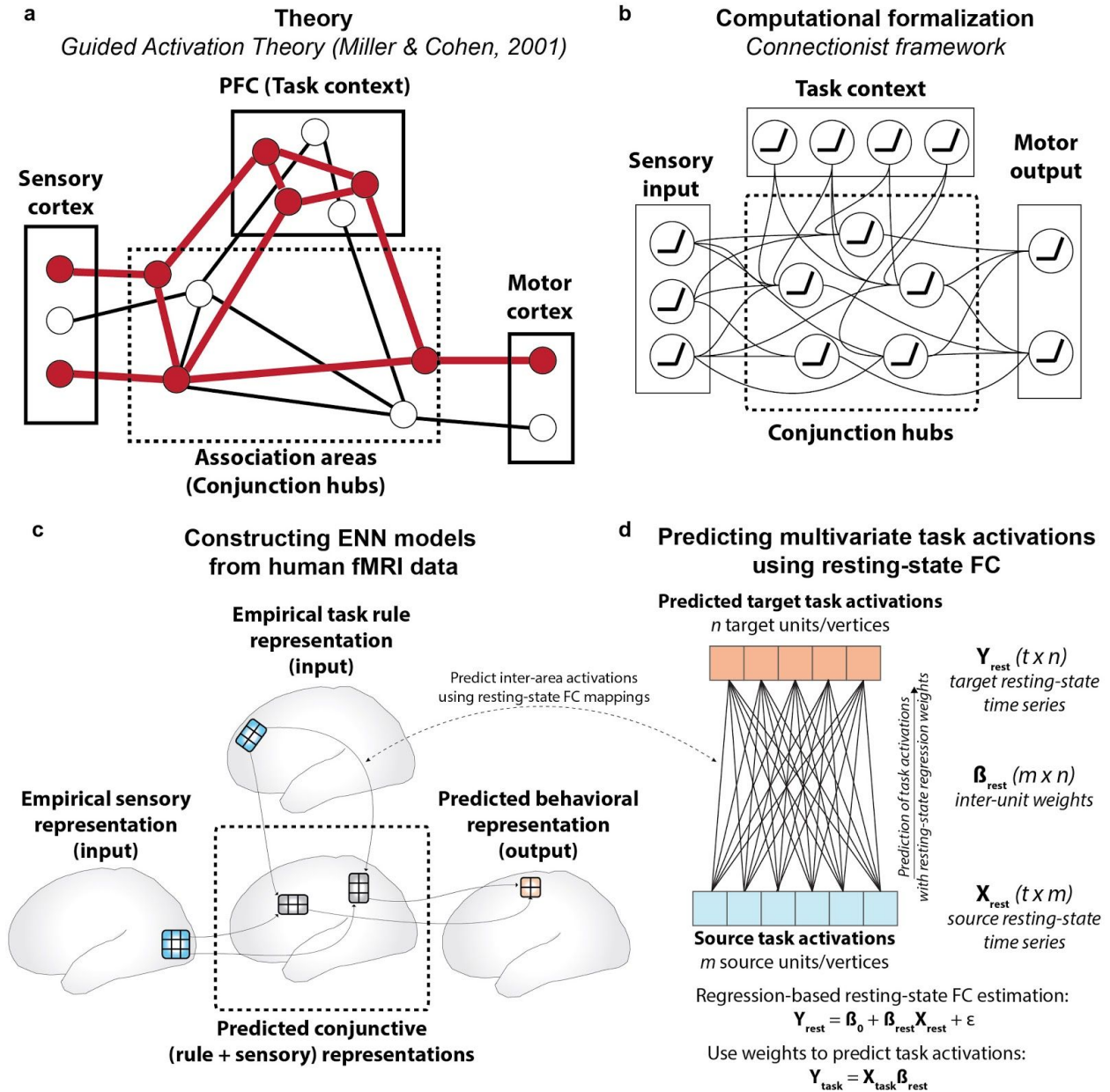
**Figure 1. Leveraging Guided Activation theory to inspire ENN models of cognitive computation during task-based fMRI. a)** A modified version of the Guided Activation theory of prefrontal cortex, highlighting a potential key role for conjunction hubs. Guided Activation theory posits that sensory cortices (left), which contain sensory stimulus-related activations, and prefrontal areas (top), which contain task context activations, integrate in association cortex to produce conjunctive activations through patterns of guided activations. Conjunctive activations are then guided to motor areas to generate motor response activations for task behavior. **b)** Guided Activation theory can be reconceptualized in a connectionist framework. This provides a formalization of how flexible sensorimotor transformations may be implemented computationally. The formalization involves the task context and sensory stimuli representing the input layer, the association units representing a hidden layer, and the behavioral (motor) responses as the output layer. **c)** Testing Guided Activation theory using task fMRI data collected in humans during context-dependent tasks. Using quantitative methods, we empirically test how different task activations (e.g., sensory stimuli and task context) form conjunctive activations to produce motor

4

response activations using activity flow mapping[16]. **d)** Guided Activation theory can be empirically tested by projecting multivariate task activations between brain areas by estimating inter-area FC weight mappings obtained from resting-state fMRI data**.** Based on the activity flow principle[16], we estimated inter-vertex mappings using regression (see Methods) on resting-state fMRI data. This approach identifies a projection that maps across distinct spatial units in empirical data, similar to how inter-layer weights propagate activity across layers in a feedforward ANN.

# Results

## Identifying brain areas containing task-relevant information

Flexible Hub theory posits that rapid updates to rule representations facilitate flexible behavior, while Guided Activation theory[3] states that sensory stimulus and task rule activations integrate in association cortex to form conjunctive representations (Fig. 1a,c). Thus, due to its comprehensive assessment of rule-guided sensorimotor behavior across 64 task contexts, we used the Concrete Permuted Rule Operations (C-PRO) task paradigm[5] to test both theories (Fig. 2a). Briefly, the C-PRO paradigm is a highly context-dependent cognitive control task, with 12 distinct rules that span three rule domains (four rules per domain; logical gating, sensory gating, motor selection). These rules were permuted within rule domains to generate 64 unique task contexts, and up to 16384 unique trials possibilities (with various stimulus pairings; see Methods).
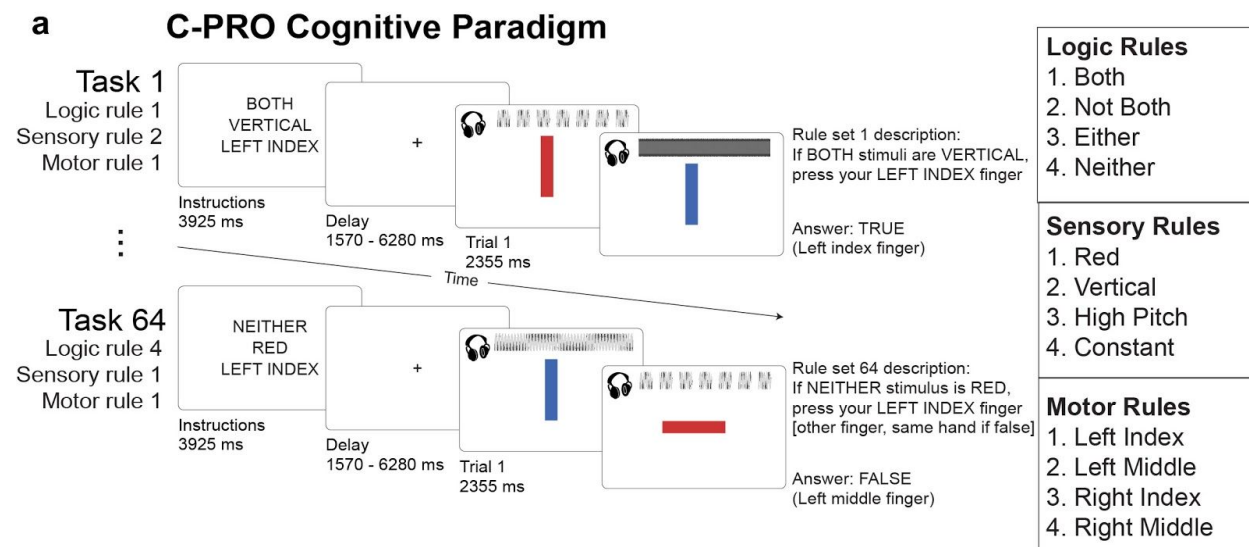


**Figure 2. The Concrete Permuted Rule Operations (C-PRO) task paradigm**[8]. For a given trial, subjects were presented with a task rule set (context), in which they were presented with three rules sampled from three different rule domains (i.e., logical gating, sensory gating, and motor selection domains). After a delay period, subjects applied the task rule set to two consecutively presented sensory stimuli (simultaneous audio-visual stimuli) and responded accordingly with button presses (index and middle fingers on either hand). We employed a miniblock design, in which for a given task rule set, three trials were presented separated by an inter-trial interval (1570ms). See Methods for additional details.

To test both Flexible Hub and Guided Activation theories, we needed to identify the set of regions responsive to different task components (sensory stimuli, task context, motor responses, and conjunctions). We first identified the set of cortical areas that contained decodable sensory stimulus representations (Fig. 3a). Because our stimuli were multimodal (audiovisual), this involved the identification of surface vertices that contained the relevant visual (color and orientation) and auditory (pitch and continuity) dimensions. We performed a four-way multivariate pattern analysis[21] (using a minimum-distance classifier[22]) to decode stimulus pairs for each of the four stimulus dimensions (e.g., red-red vs. red-blue vs. blue-red vs. blue-blue). Decoding analyses were performed within each brain parcel using the Glasser et al. atlas[23], using vertices within each parcel as decoding features. For all decoding analyses, statistical thresholding was performed using a one-sided binomial test (greater than chance=25%), and corrected for multiple comparisons using an FDR-corrected $p<0.05$ threshold. We collectively defined the units in the ENN (i.e., vertices) that contained sensory stimulus information to be the set of all vertices within the parcels that contained decodable stimulus information (Fig. 3f; Supplementary Tables 1-4).
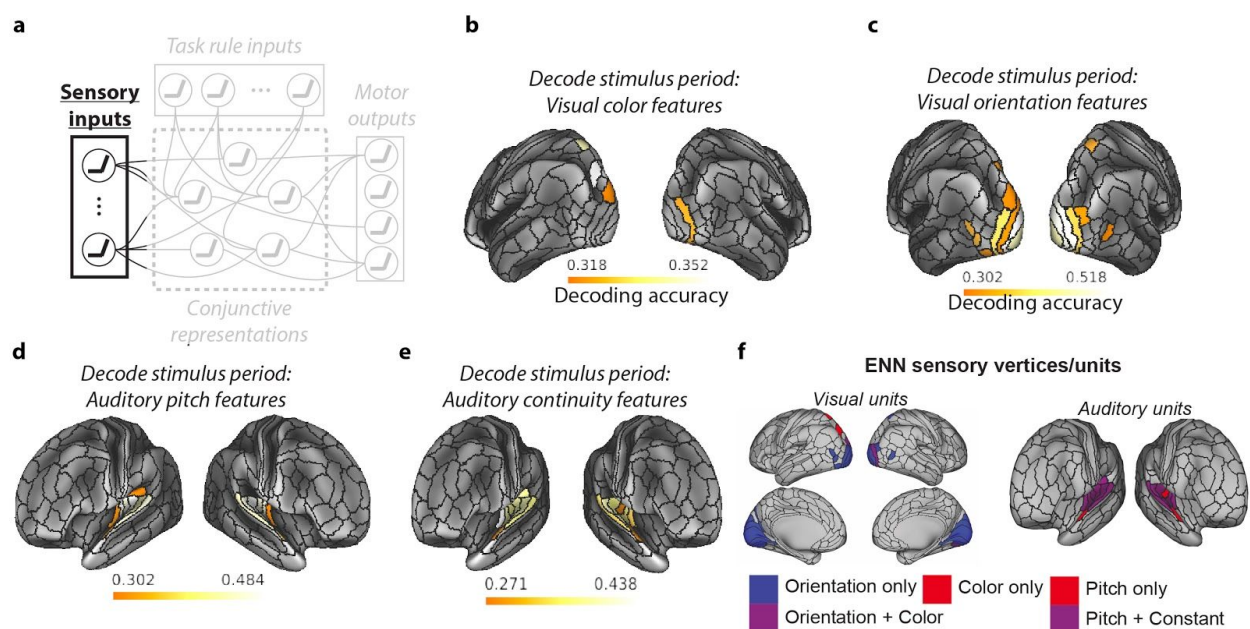


**Figure 3. Identifying sensory stimulus input units (vertices) of the ENN using multivariate pattern classification analysis. a)** We identified the sensory stimulus representations in empirical data using multivariate pattern decoding of stimulus activations. This corresponded to the sensory input component of Guided Activation theory. To decode visual features (i.e., color and orientation stimulus features) we decoded the vertices within each parcel in the visual network using a recent functional network atlas[24]. To decode auditory features (i.e., pitch and continuity) we decoded the vertices within each parcel in the auditory network (see Methods). **b)** Decoding of color features using task activation estimates (from a task GLM) during the stimulus presentation period of the C-PRO task. Chance was 25%; cortical maps were thresholded using an FDR-corrected threshold of $p<0.05$. **c)** 4-way decoding of orientation features. **d)** 4-way decoding of auditory pitch features. **e)** 4-way decoding of auditory continuity features. **f)** The

ENN sensory units, which were derived from a mask of the vertices that could successfully decode stimulus features.

Next, we performed a 12-way decoding analysis – isolated to the fMRI activation during the task instruction period – across all 12 task rules to identify the set of vertices that contained task rule information. Our previous study illustrated that rule representations are widely distributed across cortex[8], such that we tested for rule representations in every parcel in the Glasser et al. atlas (360 total parcels[23]). We again found that task rule representations were widely distributed across cortex (Fig. 4b; FDR-corrected p<0.05 threshold; Supplementary Table 6). The set of vertices that survived statistical thresholding were included as "task rule" input units in the ENN (Fig. 1c).

The C-PRO task paradigm required button presses (using index and middle fingers on either hand) to indicate task responses. Thus, to isolate finger representations in empirical neural data, we performed a univariate contrast of the vertex-wise response-evoked activation estimates during index and middle finger response windows (see Methods). For each hand, we performed a two-sided paired t-test (paired across subjects) for middle versus index finger responses in the somatomotor network[24]. Contrast maps were corrected for multiple comparisons (comparisons across vertices) using an FDR-corrected threshold of p<0.05 (Fig. 4c). Vertices that survived statistical thresholding were then selected for use as output units in the ENN (Fig. 1c).
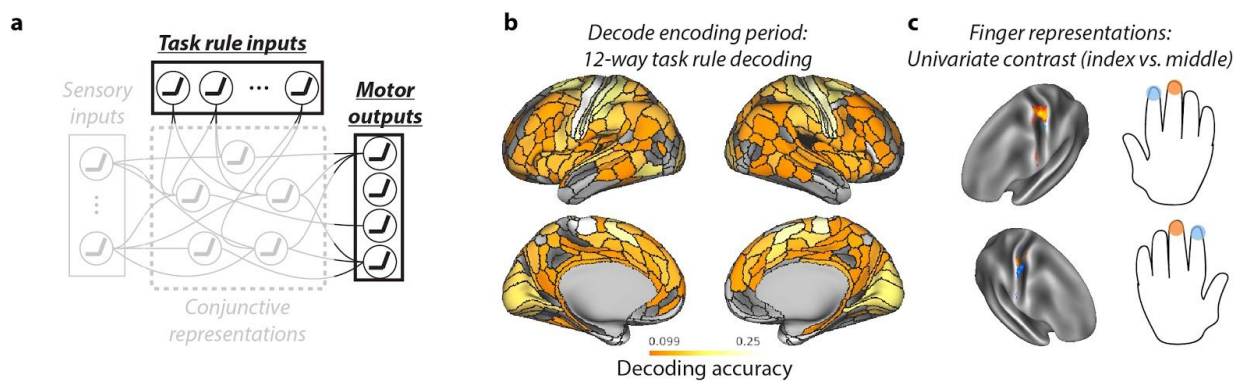


**Figure 4. Identifying ENN units (i.e., fMRI vertices) containing relevant task rule (context) and motor response (behavior) representations. a)** We identified the task rule input and motor output representations in empirical data using MVPA and univariate task activation contrasts. **b)** A 12-way decoding of each of the task rules (across the 3 rule domains) using task activations (estimated from a task GLM) during the encoding period of the C-PRO task. We applied this 12-way decoding to every parcel, given that task rule representations have been previously shown to be widely distributed across cortex[8]. Chance decoding was 8.33%; statistical maps were thresholded using an FDR-corrected p<0.05 threshold. **c)** To identify the motor/output representations, we performed a univariate contrast, contrasting the middle versus index finger response activations for each hand separately. Finger response activations were estimated during the response period, and univariate contrasts were performed on a vertex-wise basis using all vertices within the somatomotor network[24]. Contrast maps were statistically thresholded

using an FDR-corrected p<0.05 threshold. The resulting finger representations matched the placement of finger representations in the well-established somatomotor homunculus in the human brain.

## Identifying conjunction hubs

We next sought to identify conjunctive representations that could plausibly implement the transformation of input to output activations across the 64 task contexts (Fig. 5a). However, we were uncertain as to what sorts of activation patterns (i.e., representations) we would expect in putative conjunction hubs. Thus, we began by building an ANN that formalizes Guided Activation theory (Fig. 1b). We trained the ANN model on an analogous version of the C-PRO task until the model achieved 99.5% accuracy (see Methods). We were specifically interested in characterizing the representations in the hidden layer, since these activations necessarily integrated task rule and sensory stimulus activations (i.e., conjunctions). To identify the task rule and sensory stimulus conjunctive representations, we performed a representational similarity analysis (RSA) on the hidden layer of the ANN[22]. The representational similarity matrix (RSM) of the hidden layer consisted of 28 task activation features: 12 task rules (which spanned the 3 rule domains), and 16 stimulus pairings (which spanned each sensory dimension). We then compared the RSM of the ANN's hidden units (Fig. 5b) to RSMs of each brain region in the empirical fMRI data (Fig. 5c). This provided a map of brain regions with similar representations to those of the ANN's hidden units, which contain the conjunction of task rule and sensory stimulus activations.



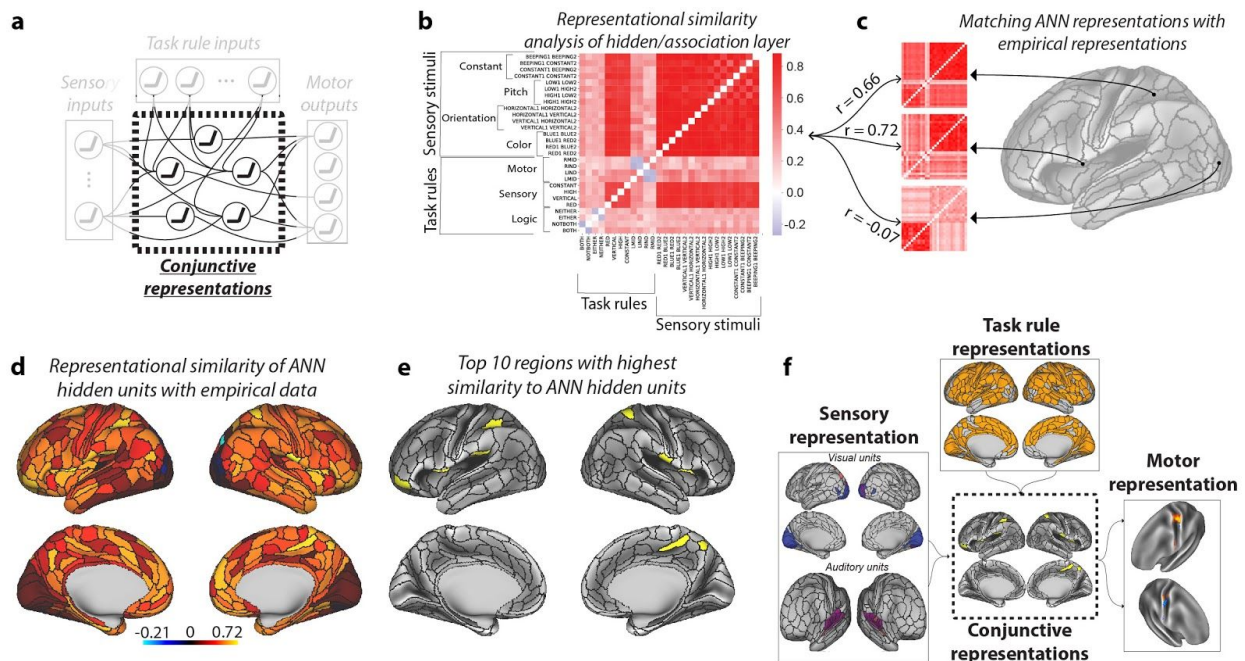**Figure 5. Identifying conjunction hubs: brain areas (vertices) that contain task-relevant conjunctions of sensory stimulus and task rule activations. a)** Guided Activation theory states that there exist a specific set of association (or hidden) areas that integrate sensory stimulus and task context activations to select appropriate motor response activations. Computationally, this corresponded to the

2nd hidden layer in our ANN implementation (see Methods). **b)** We therefore used the representational similarity matrix (RSM) of the ANN's hidden layer as a blueprint to identify analogous conjunctive activations in empirical data. **c)** We constructed RSMs for each brain parcel (using the vertices within each parcel as features). We evaluated the correspondence between the representational geometry of the ANN's 2nd hidden layer and each brain parcel's representational geometry. Correspondence was assessed by taking the correlation of the upper triangle of the ANN and empirical RSMs. **d)** The representational similarity of ANN hidden units and each brain parcel. **e)** We showed the top 10 regions with highest similarity to the ANN hidden units. **f)** The full ENN architecture for the C-PRO task. We identified the vertices that contained task-relevant rule, sensory stimulus, conjunctive, and motor output activations.

To evaluate the similarity of the ANN's hidden representational geometry with each brain parcel, we computed the similarity (using Spearman's correlation) of the ANN's RSM with the brain parcel's RSM (Fig. 5c). This resulted in a cortical map, which showed the representational similarity between each brain region and the ANN's hidden representations (Fig. 5d). For our primary analysis, we selected the top 10 parcels with highest similarity to the ANN's hidden units to represent the set of spatial units that contain putative conjunctive activations in the ENN (Fig. 5e). The conjunction hubs were most strongly represented by the cingulo-opercular network, a network previously reported to be involved in task set maintenance (Supplementary Fig. 2; Supplementary Table 5)[25]. However, we also performed ENN simulations using the top 20, 30, and 40 regions with highest similarity to the ANN hidden units below.

## Task-performing neural network simulations via empirical connectivity

The previous sections provided the groundwork for constructing an ENN model from empirical data. After estimating the connectivity weights between the surface vertices between ENN layers using resting-state fMRI (see Methods), we next sought to evaluate whether we could use this ENN to produce representational transformations sufficient for performing the C-PRO paradigm. This would demonstrate that the empirical input activations (task rule and sensory stimulus activations) and the estimated connectivity patterns between ENN layers are sufficient to approximate the cognitive computations involved in task performance.

The primary goal was to predict the motor response activation pattern (i.e., behavior) yielding correct task performance. The only inputs to the model were a combination of activation patterns for a specific task context (rule combination) and sensory stimulus pair sampled from empirical data (Fig. 6a). The outputs of the model were the predicted motor response activation pattern in motor cortex that should correspond to the correct button press (Fig. 6c). High correspondence between the predicted and actual motor activation patterns would constitute an empirical identification of representational transformation in the brain, where task rule and sensory stimulus activity is transformed into task-appropriate motor response activation patterns.

9

$$Y = f\left( X_{rule} \, W_{rule \to hidden} + X_{stimulus} \, W_{stimulus \to hidden} \right) W_{hidden \to output}$$

Y : *predicted motor activation pattern*
W : *Estimated connectivity weights from A ~ B*
f : *Rectified linear function (threshold negative values)*
X : *Input activation patterns (either rule activations or sensory stimulus activations)*
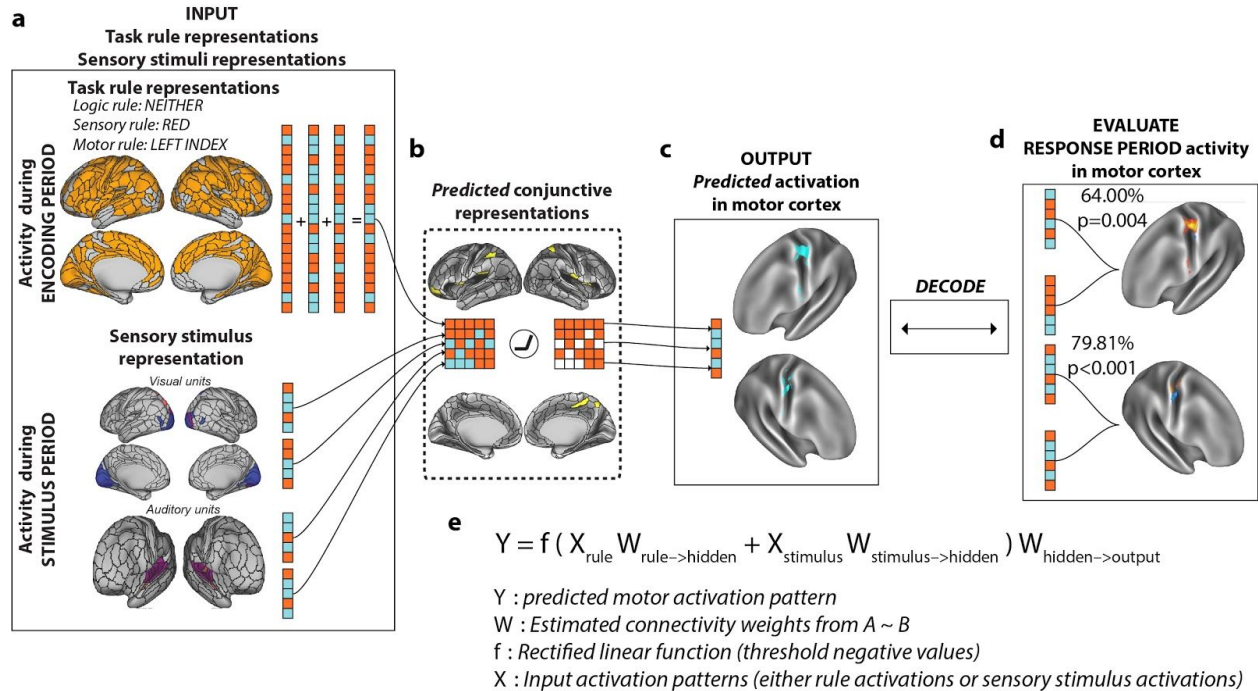
**Figure 6. Simulating context-dependent sensorimotor transformations with empirically-estimated task activations and inter-unit functional connectivity estimates.** We constructed the ENN by identifying the vertices that contained task rule, sensory stimulus, and motor response activations (via decoding) and by estimating the resting-state FC weights between them. **a)** The input layer, consisting of vertices with decodable task rule and sensory stimulus activations. **b)** Through activity flow mapping, input representations were mapped onto surface vertices in conjunction hubs. The activity flow-mapped vertices were passed through a nonlinearity, which removed any negative values. This threshold was chosen given the difficulty in interpreting *predicted* negative BOLD values. **c)** The predicted conjunctive representations were then activity flow-mapped onto the motor output vertices, generating a predicted motor activation pattern. **d)** These predicted motor activations were then tested against the actual motor response activations of other subjects using a leave-8-subject out cross validation scheme. A decoder was trained on the predicted motor response activations and tested on the actual motor response activations of the held-out cohort (see Methods and Supplementary Fig. 1). **e)** An equation summarizing the ENN model's computations.

Simulating activity flow in the ENN involved first extracting the task rule activation patterns (inputs) for a randomly generated task context (see Methods and Supplementary Fig. 1). Independently, we sampled sensory stimulus activation patterns for each stimulus dimension (color, orientation, pitch, continuity) (Fig. 3). Then, using activity flow mapping with resting-state FC weights, we projected the activation patterns from the input vertices onto the conjunction hub vertices (Fig. 6b). The predicted conjunction hub activation pattern was then passed through a simple rectified linear function, which removed any negative values (i.e., any values lower than resting-state baseline; see Methods). Thresholded values were then projected onto the output layer vertices in motor cortex (Fig. 6c), yielding a predicted response activation pattern. The sequence of computations performed to generate a predicted motor activation pattern (Fig. 6a-c) is encapsulated by the equation in Fig. 6e. Thus, predicted motor activation

10

patterns can be generated by randomly sampling different task context and sensory stimuli activations for each subject.

While the above procedure yielded a predicted activation pattern in the motor output layer, these predictions may not actually yield meaningful activation patterns. Thus, we evaluated whether the predicted motor activation patterns would accurately predict the *actual* motor response activation pattern extracted (via GLM) during the response period. Activity flow simulations generated predicted motor responses for each subject (Supplementary Fig. 1). This yielded four predicted motor response activations per subject, one for each behavioral response. Importantly, the predicted motor response activations were generated using only input task activations from the task encoding period and stimulus presentation period (Fig. 6a). Independently, each subject also had four corresponding real motor response activations, which were estimated from the task GLM during the response period. Using a leave-8-subjects out cross-validation scheme, we trained a decoder on the four predicted motor responses and decoded the four actual motor responses (Fig. 6c,d). Training a decoder on the predicted activations and decoding the actual activations (rather than vice versa) made this analysis more in line with a prediction perspective – we could test if, in the absence of any motor task activation information, the ENN could predict actual motor response activation patterns that correspond to correct behavior.

We note that this decoding analysis is highly non-trivial, given that the predicted motor responses are independent from the test set (actual motor responses) in three ways: 1) The predicted motor responses were generated from task rule and stimulus activation patterns, which (due to temporal separation in the task paradigm and counterbalancing) were statistically independent from the motor responses; 2) The motor response predictions were generated via activity flow mapping, and thus from a spatially independent set of vertices (see Methods); 3) The actual motor responses in the test set were sampled from independent subjects. By simulating neural network computations from stimulus and task context activations to predict motor response, we accurately decoded the correct finger response on each hand separately: decoding accuracy of right hand responses = 64.00%, non-parametric p=0.004; decoding accuracy of left hand responses = 79.81%, non-parametric p<0.001. These results demonstrate that task rule and sensory stimulus activations can be transformed into motor output activations by simulating multi-step neural network computations using activity flow mapping on empirical fMRI data.

## The importance of the conjunctive representations

We next evaluated whether specific components of the ENN model were necessary to produce accurate stimulus-response transformations. We first sought to evaluate the role of the conjunction hubs (hidden layer) in model performance. This involved re-running the ENN with the conjunction hubs removed (Fig. 7c), which required resting-state FC weights to be re-estimated between the input and motor output layer directly. We found that the removal of conjunction hubs severely impaired task performance to chance accuracy (RH accuracy=49.05%, p=0.54; LH accuracy=50.14%, p=0.46; Fig. 7h,i). This illustrated the

importance of conjunction hub computations in producing the conjunctive activations required to perform context-dependent stimulus-response mappings[15].



**Figure 7. Systematic alteration of ENN model architecture verifies validity of "full S-R model" results. a)** We first benchmarked the motor response decoding accuracy for each hand separately using a standard cross-validation scheme on motor activation patterns for each hand (tested across subjects). This standard motor decoding was done independently of modeling sensorimotor transformations. **b)** The full stimulus-response model, taking stimulus and context input activations to predicting motor response patterns in motor cortex. **c)** The ENN model after entirely removing the hidden layer. **d)** The ENN model, where we randomly sampled regions in the hidden layer (conjunction hubs) 1000 times and estimated task performance. **e)** The ENN model after removing the nonlinearity (ReLU) function in the hidden layer.

**f)** The ENN model after lesioning connections from the task context input activations. **g)** The ENN model, where we shuffled the connectivity patterns from the stimulus and context layers 1000 times. **h)** Benchmarking the performances of all model architectures. Accuracy distributions were obtained by bootstrapping samples (leave-8-out cross-validation scheme and randomly sample within the training set). Boxplot whiskers reflect the 95% confidence interval. Grey distributions indicate the null distribution generated from permutation tests (permuting labels 1000 times). (*** = p<0.001; ** = p<0.01; * = p<0.01) **i)** Summary statistics of model performances. Reported accuracies are the mean of the bootstrapped samples.

We next replaced conjunction hubs with randomly sampled parcels in empirical data. This assessed the importance of using the ANN's hidden layer RSM to identify conjunction hubs in data (Fig. 7d). We sampled random parcels 1000 times, recomputing the inter-layer vertex-wise FC each time. The distribution of randomly selected conjunction hubs did not yield task performance accuracies that were statistically different than chance for both hands (RH mean accuracy=50.87%, p= 0.47; LH mean accuracy 50.85%, p=0.44; Fig. 7h,i). However, the overall distribution had high variance, indicating that there may be other sets of conjunction hubs that would yield above-chance (if not better) task performance. However, compared to the conjunction hubs we identified by matching empirical brain representations with ANN representations, we found that the ANN-matched conjunction hubs performed better than 85.2% of all randomly selected conjunction hubs for RH responses, and greater than 97.7% of all randomly selected conjunction hubs for LH responses.

In addition, we evaluated whether the precise number of hidden regions was critical to task performance. We ran the full ENN model, but instead of using only the top 10 regions with highest similarity to the ANN's hidden layer's representations, we constructed ENN variants containing the top 20, 30, and 40 hidden regions. We found that we were able to reproduce correct task performance using 20 hidden regions (RH accuracy=63.90%, p<0.001; LH accuracy=76.95%, p<0.001). Using 30 hidden regions yielded reduced yet above-chance accuracies for RH responses, but not for LH responses (RH accuracy=59.83%, p=0.024; LH accuracy=43.54%, p=0.917). Inclusion of an additional 10 hidden regions (totaling 40 hidden regions) did not yield above-chance predictions of motor responses for either hand. These results suggest that conjunction hubs were better identified the greater the similarity of a region's representational geometry was to that of the ANN's hidden layer.

## The importance of nonlinearities when combining rule and stimulus activations

We next removed the thresholding of negative BOLD values (i.e., those lower than resting baseline) in the hidden layer. This is conceptually similar to removing nonlinearities in an ANN (Fig. 7e). We found that the removal of the ReLU function significantly impaired model performance (RH accuracy=47.74%, p=0.70; LH=47.90%, p=0.692; Fig. 7h and 7i). This is likely due to the fact that context-dependent sensorimotor transformations require a nonlinear mapping between stimulus-response pairs, as predicted by prior computational studies[26,27].

13

## Removing task context impairs task performance

We next sought to evaluate the importance of including task rule activations in model performance. To remove context information, we lesioned all connections from the rule input layer to the hidden layer. This was achieved by setting all resting-state FC connections from the context input layer to 0 (Fig. 7f). We ran the model on the exact same set of tasks, and found that as hypothesized, model performance was at chance without task context information (RH accuracy=50.00%, p=0.44; LH=50.00%, p=0.47; Fig. 7h,i). This illustrated that the model implemented a representational transformation from task context and sensory stimulus representations to the correct motor responses.

## The influence of specific functional network topography

We next evaluated whether the empirically-estimated connectivity topography was critical to successful task performance. This involved shuffling the connectivity weights within the context and stimulus input layers 1000 times (Fig. 7g). While we hypothesized that the specific resting-state FC topography would be critical to task performance, we found that shuffling connectivity patterns yielded a very high variance distribution of task performance (Fig. 7h). While the mean across all connectivity shuffles were approximately at chance for both hands (RH mean accuracy=50.90%, p=0.45; LH mean accuracy=50.39%, p=0.48), we found that there were some connectivity configurations that would significantly improve task performance, and other connectivity configurations that would yield significant below chance task performance. Notably, the FC topography that was estimated from resting-state fMRI (the full S-R model, without shuffling; Fig. 7b) performed greater than 85.3% of all connectivity reconfigurations in RH responses, and greater than 97.7% of all connectivity reconfigurations for LH responses. This indicates that while there may exist better connectivity patterns for task performance, the weights derived from resting-state fMRI were sufficient to model correct task performance. We note that while the distribution of performance accuracies when shuffling FC weights and randomly sampling hidden layers are quite similar, these two permutation analyses control for fundamentally distinct properties of the ENN: specificity of FC topography versus specificity of conjunction hubs.

# Discussion

Characterizing how different cognitive activations are transformed throughout the brain would fill a critical gap in understanding how the brain implements cognitive computations[28–30]. To address this gap, we built a task-performing ENN from empirical data to characterize representational transformations during a cognitive control task. First, we identified brain vertices that were selective for task rules, sensory stimuli, motor responses, and conjunctions. Second, we mapped resting-state FC weights between these areas using multiple linear regression. Finally, using activity flow mapping, we found that incoming sensory and task rule activations were transformed via conjunction hubs to produce above-chance behavioral

predictions of outgoing motor response activations. These findings suggest that flexible cognitive control is implemented by guided activations, as originally suggested by Guided Activation theory[3].

The present results build on Flexible Hub theory and other findings emphasizing the role of cognitive control networks (CCNs) in highly flexible cognition[1,25,31,32]. Consistent with previous accounts, we found that the task rule layer and conjunction hubs are most strongly affiliated with CCNs (e.g., cingulo-opercular and frontoparietal networks) (Supplementary Fig. 2)[25,31]. (However, we note that other functional networks also represented task rules, though to a lesser extent.) In addition, several studies of rapid instructed task learning found that CCNs represent rules compositionally in activity[7,10,32] and FC[14,15] patterns, which are considered essential for flexible reuse of task components[10,12,14]. The present results also demonstrate that the CCN and other networks use compositional rule representations, since the ENN rule activation inputs contained three rules whose fMRI activity patterns were added compositionally to create the full task context. Critically, however, we found that these compositional codes were not enough to implement flexible task performance – rather, conjunctive representations were required to interact non-linearly with these compositional representations. Moreover, our results showed that without conjunctive representations producing conditional interactions (e.g., through conjunction hub lesioning), the task performance of the ENN was substantially impaired. It will be important for future research to determine the exact relationship between compositional and conjunctive representations in implementing flexible cognitive programs.

The ENN characterized the representational transformations required to transform task input activations to output activations (in motor cortex) directly from data. Model parameters, such as unit identification and inter-unit connectivity estimation, were estimated *without optimizing for task performance*. This contrasts with mainstream machine learning techniques that iteratively train ANNs that directly optimize for behavior[12,17,18,33,34]. Our approach enabled the construction of functioning ENNs with above-chance task performance without optimizing for behavior; instead, we were able to derive parameters from empirical neural data alone. These results suggest that the human brain's intrinsic network architecture, as estimated with human fMRI data, is informative regarding the design of task-performing functioning models of cognitive computation.

We showed that the specific FC topography could predict inter-area transformations. In contrast, shuffling these specific inter-area FC topographies yielded ENNs with highly variable task performances, suggesting the computational utility of the empirically-estimated FC patterns. Previous work has illustrated that the functional network architecture of the brain emerges from a structural backbone[35–39]. Building on this work, we recently proposed that the functional network architecture of the brain can be used to build network coding models – models of brain function that describe information encoding and decoding processes constrained by empirically-estimated connectivity[40]. Related proposals have also been suggested in the electron microscopy connectomics literature, suggesting that structural wiring diagrams of the brain (e.g., in drosophila) can inform functional models of biological systems (e.g., the drosophila's visual system)[44,47]. Consistent with these proposals, our findings establish that the intrinsic functional network architecture in humans provides a meaningful foundation from which to implement cognitive computations.

Despite strong evidence that the estimated functional network model can perform tasks, there are several theoretical and methodological limitations. First, though we perform numerous control analyses by either lesioning or altering the ENN architecture (Fig. 7), the space of alternative possible models that can potentially achieve similar (if not better) task performances is large. For example, here we assumed only a single hidden layer (one layer of 'conjunction hubs'). However, it is possible – if not probable – that such transformations actually involve a large sequence of transformations, similar to how the ventral visual stream transforms visual input into object codes, from V1 to inferior temporal cortex[18,19]. It is therefore likely that the identification of conjunction hubs is likely dependent on both specific task demands and the targeted level of analysis (e.g., neuronal circuits versus large-scale functional networks). Here we opted for the simplest possible network model that involved conjunction hubs at the level of large-scale functional networks. Starting from this simple model allowed us to reduce potential extraneous assumptions and model complexity (such as modeling the extraction of stimulus features from early visual areas) which likely would have been necessary in more complex and detailed models. However, the current findings provide a strong foundation for future studies to unpack the mechanisms of finer-grained computations important for adaptive behavior.

Another assumption in the ENN was that activations were guided by additive connectivity weights. Additive connectivity weights assume inter-area predicted activations are the sum of source activations weighted by connections. One potential alternative (among others) would have been multiplicative guided activations; weighted activations that are multiplied (rather than summed) from incoming areas, which has been previously proposed as a potential alternative to designing ANNs[42]. However, several recent studies have suggested that inter-area activations are predicted via additive connectivity weights in both human fMRI[8,16], the primate visual system[20], and the drosophila's visual system[39]. Nevertheless, it will be important for future work to systematically test alternative network architectures and dynamics in producing functional ENN models.

Finally, another limitation is that we constructed an ENN model that did not model realistic dynamics. Typical experimental paradigms include separate intervals for encoding, delay, stimulus, and response periods, since cognitive processing occurs over time. Here, we did not explicitly model temporal dynamics based on the empirical data when simulating the ENN. (However, we note that activation estimates for different task components, such as task encoding and stimulus presentation, were obtained from temporally distinct intervals.) Nevertheless, though it is likely that temporal dynamics (with recurrent feedback) likely play a role in shaping cognitive computations, we illustrate here that simple dynamics (i.e., rules + sensory inputs → conjunction hubs → motor outputs) involving the interplay of static activations are sufficient to model representational transformations. It will be important for future studies to construct task-performing brain models that can simulate temporal and recurrent dynamics constrained by empirical data, as this can provide a more detailed computational account of the representational transformations that contribute to behavioral variability.

In conclusion we constructed an ENN model capable of performing adaptive cognitive control tasks. This model provides strong evidence for the well-known Guided Activation theory by providing a computational implementation of the theory that is directly estimated from empirical data. We first identified the relevant brain representations associated with different

task features. We then used an ANN to identify conjunction hubs that were critical to the selective integration of task input information for motor response selection. Finally, by estimating FC patterns from resting-state fMRI data, we parameterized a network model to generate predictive stimulus-to-response transformations using activity flow mapping. We expect that these findings will drive new investigations into characterizing the neural implementation of cognitive computations, providing dual insight into how the brain implements cognitive processes and how such knowledge can inform the design of ANN architectures.

# Methods

## Participants

Data were collected from 106 human participants across two different sessions (a behavioral and an imaging session). Participants were recruited from the Rutgers University-Newark community and neighboring communities. Technical error during MRI acquisition resulted in removing six participants from the study. Four additional participants were removed from the study because they did not complete the behavior-only session. fMRI analysis was performed on the remaining 96 participants (54 females). All participants gave informed consent according to the protocol approved by the Rutgers University Institutional Review Board. The average age of the participants that were included for analysis was 22.06, with a standard deviation of 3.84. Additional details regarding this participant cohort have been previously reported[43].

## C-PRO task paradigm

We used the Concrete Permuted Operations (C-PRO) paradigm (Fig. 2a) during fMRI acquisition, and used a computationally analogous task when training our ANN model. The details of this task are described below, and are adapted from a previous study[8].

The C-PRO paradigm is a modified version of the original PRO paradigm introduced in Cole et al., (2010)[44]. Briefly, the C-PRO cognitive paradigm permutes specific task rules from three different rule domains (logical decision, sensory semantic, and motor response) to generate dozens of novel and unique task contexts. This creates a context-rich dataset in the task configuration domain akin in some ways to movies and other condition-rich datasets used to investigate visual and auditory domains[5]. The primary modification of the C-PRO paradigm from the PRO paradigm was to use concrete, sensory (simultaneously presented visual and auditory) stimuli, as opposed to the abstract, linguistic stimuli in the original paradigm. Visual stimuli included either horizontally or vertically oriented bars with either blue or red coloring. Simultaneously presented auditory stimuli included continuous (constant) or non-continuous (non-constant, i.e., "beeping") tones presented at high (3000Hz) or low (300Hz) frequencies. Fig. 2a demonstrates two example task-rule sets for "Task 1" and "Task 64". The paradigm was presented using E-Prime software version 2.0.10.353[45].

Each rule domain (logic, sensory, and motor) consisted of four specific rules, while each task context was a combination of one rule from each rule domain. A total of 64 unique task contexts (4 logic rules x 4 sensory rules x 4 motor rules) were possible, and each unique task set was presented twice for a total of 128 task miniblocks. Identical task sets were not presented in consecutive blocks. Each task miniblock included three trials, each consisting of two sequentially presented instances of simultaneous audiovisual stimuli. A task block began with a 3925 ms instruction screen (5 TRs), followed by a jittered delay ranging from 1570 ms to 6280 ms (2 – 8 TRs; randomly selected). Following the jittered delay, three trials were presented for 2355 ms (3 TRs), each with an inter-trial interval of 1570 ms (2 TRs). A second jittered delay followed the third trial, lasting 7850 ms to 12560 ms (10-16 TRs; randomly selected). A task block lasted a total of 28260 ms (36 TRs). Subjects were trained on four of the 64 task contexts for 30 minutes prior to the fMRI session. The four practiced rule sets were selected such that all 12 rules were equally practiced. There were 16 such groups of four task sets possible, and the task sets chosen to be practiced were counterbalanced across subjects. Subjects' mean performance across all trials performed in the scanner was 84% (median=86%) with a standard deviation of 9% (min=51%; max=96%). All subjects performed statistically above chance (25%).

## fMRI acquisition and preprocessing

The following fMRI acquisition details is taken from a previous study that used the identical protocol (and a subset of the data)[8].

Data were collected at the Rutgers University Brain Imaging Center (RUBIC). Whole-brain multiband echo-planar imaging (EPI) acquisitions were collected with a 32-channel head coil on a 3T Siemens Trio MRI scanner with TR=785 ms, TE=34.8 ms, flip angle=55°, Bandwidth 1924/Hz/Px, in-plane FoV read=208 mm, 72 slices, 2.0 mm isotropic voxels, with a multiband acceleration factor of 8. Whole-brain high-resolution T1-weighted and T2-weighted anatomical scans were also collected with 0.8 mm isotropic voxels. Spin echo field maps were collected in both the anterior to posterior direction and the posterior to anterior direction in accordance with the Human Connectome Project preprocessing pipeline[49]. A resting-state scan was collected for 14 minutes (1070 TRs), prior to the task scans. Eight task scans were subsequently collected, each spanning 7 minutes and 36 seconds (581 TRs). Each of the eight task runs (in addition to all other MRI data) were collected consecutively with short breaks in between (subjects did not leave the scanner).

## fMRI Preprocessing

The following details are adapted from a previous study that used the same preprocessing scheme on a different data set[50].

Resting-state and task-state fMRI data were minimally preprocessed using the publicly available Human Connectome Project minimal preprocessing pipeline version 3.5.0. This pipeline included anatomical reconstruction and segmentation, EPI reconstruction, segmentation, spatial normalization to standard template, intensity normalization, and motion correction[64]. After minimal preprocessing, additional custom preprocessing was conducted on CIFTI 64k grayordinate standard space for vertex-wise analyses using a surface based atlas[23].

18

This included removal of the first five frames of each run, de-meaning and de-trending the time series, and performing nuisance regression on the minimally preprocessed data[51]. We removed motion parameters and physiological noise during nuisance regression. This included six motion parameters, their derivatives, and the quadratics of those parameters (24 motion regressors in total). We applied aCompCor on the physiological time series extracted from the white matter and ventricle voxels (5 components each extracted volumetrically)[52]. We additionally included the derivatives of each component time series, and the quadratics of the original and derivative time series (40 physiological noise regressors in total). This combination of motion and physiological noise regressors totaled 64 nuisance parameters, and is a variant of previously benchmarked nuisance regression models[51].

## fMRI task activation estimation

We performed a standard task GLM analysis on fMRI task data to estimate task-evoked activations from different conditions. Task GLMs were fit for each subject separately, but using the fully concatenated task data set (concatenated across 8 runs). We obtained regressors for each task rule (during the encoding period), each stimulus pair combination (during stimulus presentation), and each motor response (during button presses). For task rules, we obtained 12 regressors that were fit during the encoding period, which lasted 3925ms (5 TRs). For logic rules, we obtained regressors for "both", "not both", "either", and "neither" rules. For sensory rules, we obtained regressors for "red", "vertical", "high", and "constant" rules. For motor rules, we obtained regressors for "left middle", "left index", "right middle", and "right index" rules. Note that a given encoding period contained overlapping regressors from each of the logic, sensory, and motor rule domains. However, the regressors were not collinear since specific rule instances were counterbalanced across all encoding blocks.

To obtain activations for sensory stimuli, we fit regressors for each stimulus pair. For example, for the color dimensions of a stimulus, we fit separate regressors for the presentation of red-red, red-blue, blue-red, and blue-blue stimulus pairs. This was done (rather than fitting regressors for just red or blue) due to the inability to temporally separate individual stimuli with fMRI's low sampling rate. Thus, there were 16 stimulus regressors (four conditions for each stimulus dimension: color, orientation, pitch, continuity). Stimulus pairs were presented after a delay period, and lasted 2355ms (3 TRs). Note that a given stimulus presentation period contained overlapping regressors from four different conditions, one from each stimulus dimension. However, the stimulus regressors were not collinear since stimulus pairings were counterbalanced across all stimulus presentation periods (e.g., red-red stimuli were not exclusively presented with vertical-vertical stimuli).

Finally, to obtain activations for motor responses (or finger button presses), we fit regressors for each motor response. There were four regressors for motor responses, one for each finger (i.e., left middle, left index, right middle, right index fingers). Responses overlapped with the stimulus period, so we fit regressors for each button press during the 2355ms (3 TR) window during stimulus presentations. Note, however, that while response regressors overlapped with stimulus regressors, response regressors were not collinear with stimulus presentations. This is because a response is statistically independent from a stimulus pair,

19

enabling the extraction of meaningful response activation patterns. A strong validation was that the finger representations could be reliably extracted according to the appropriate topographic organization in somatomotor cortex (Fig. 4c).

(For a schematic of how task GLMs were performed, see Supplementary Fig. 3. For the task design matrix of an example subject, see Supplementary Fig. 4.)

## fMRI decoding: Identifying sensory stimulus representations

Decoding analyses were performed to identify the brain areas that contained relevant task context and sensory stimulus representations. To identify the brain areas that contained relevant sensory stimulus representation, we performed four, four-way decoding analyses on each stimulus dimension: color (vision), orientation (vision), pitch (audition), constant (audition). For color stimulus information, we decoded activation patterns where the stimulus pairs were red-red, red-blue, blue-red, and blue-blue. For orientation stimulus information, we decoded activation patterns where the stimulus pairs were vertical-vertical, vertical-horizontal, horizontal-vertical, horizontal-horizontal. For pitch stimulus information, we decoded activation patterns where the stimulus pairs were high-high, high-low, low-high, and low-low. Finally, for constant (beeping) stimulus information, we decoded activation patterns where the stimulus pairs were constant-constant, constant-beeping, beeping-constant, beeping-beeping.

Decoding analyses were performed using the vertices within each parcel as decoding features. We limited decoding to visual network parcels for decoding visual stimulus features, and auditory network parcels for decoding auditory stimulus features. Visual parcels were defined as the VIS1 and VIS2 networks in Ji et al. (2019)[24], and auditory networks as the AUD network. We performed a group-level decoding analysis, with a leave-8-subjects out cross-validation scheme. The choice of leaving 8 (out of 96) subjects out was due to recent studies suggesting that test sets should contain roughly 10% of the entire data set to yield stable predictive estimates of the test-set[25]. Moreover, of the 88 subjects that remained in the train set pool (for each cross-validation fold), the training set was randomly sampled (with replacement, number of bootstrapped samples per fold = 88). We used a minimum-distance classifier (based on Pearson's correlation score), where a test set sample would be classified as the condition whose centroid is closest to in the multivariate activation pattern space[22]. P-values were calculated using a binomial test. Statistical significance was assessed using a False Discovery Rate (FDR) corrected threshold of p<0.05 across all 360 regions.

## fMRI decoding: Identifying task rule representations

To identify the brain areas that contained task rule information, we performed a 12-way decoding analysis on the activation patterns for each of the 12 task rules. We used the same decoding and cross-validation scheme as above (for identifying sensory stimulus representations). However, we ran the decoding analyses on all 360 parcels, given previous evidence that task rule information is widely distributed across cortex[8]. P-values were calculated using a binomial test. Statistical significance was assessed using an FDR-corrected threshold of p<0.05 across all 360 regions.

20

## fMRI activation analysis: Identifying motor response activations

To identify the brain areas/vertices that contained motor response information, we performed univariate analyses to identify the finger press activations in motor cortex. We performed two univariate activation contrasts, identifying index and middle finger activations on each hand. For each hand, we performed a two-sided group paired (by subject) t-test contrasting index versus middle finger representations. We constrained our analyses to include only vertices in the somatomotor network. Statistical significance was assessed using an FDR-corrected $p<0.05$ threshold, resulting in a set of vertices that were selective to button press representations in motor cortex (see Fig. 4c).

We subsequently performed a decoding analysis on these sets of vertices (see Fig. 7h). We decoded finger presses on each hand separately. Note that this decoding analysis is circular, since we had already determined that the selected vertices contained relevant information with regards to motor responses (via a univariate t-test). However, this provided an important benchmark to evaluate how well we could predict motor button responses using only context and stimulus activations (described below) relative to cross-validation of motor button response activations (i.e., a noise ceiling). Similar to the previous decoding analyses, we performed a leave-8-out cross validation scheme using a minimum-distance classifier, bootstrapping training samples for each fold. Moreover, because the decoding analysis was limited to a single ROI (as opposed to across many parcels/ROIs), we were able to compute confidence intervals (by bootstrapping cross-validation folds) and run nonparametric permutation tests since it was computationally tractable. We ran each cross-validation scheme 1000 times to generate confidence intervals. Null distributions were computed by randomly permuting labels 1000 times. P-values were computed by comparing the null distribution against the mean of the bootstrapped accuracy values.

## Identifying conjunctive representations: ANN construction

We trained a simple feedforward ANN on a computationally analogous form of the C-PRO task. This enabled us to investigate how task rule and stimulus activations integrate into conjunctive representations in an ANN's hidden layer.

To model the task context input layer, we designated an input unit for each task rule across all rule domains. Thus, we had 12 units in the task context layer. A specific task context (or rule set) would selectively activate three of the 12 units; one logic rule, one sensory rule, and one motor rule. Input activations were either 0 or 1, indicating an active or inactive state.

To model the stimulus input layer, we designated an input unit for a stimulus pair for each sensory dimension. To isolate visual color stimulus pairings, we designated input units for a red-red pairing, red-blue pairing, blue-red pairing, and blue-blue pairing. (Note that each unit represented a stimulus pair because the ANN had no temporal dynamics to present consecutive stimuli.) To isolate visual orientation stimulus pairings, we designated inputs for a vertical-vertical, vertical-horizontal, horizontal-vertical, and horizontal-horizontal stimulus pairing. To isolate auditory pitch stimulus pairings, we designated input units for high-high, high-low, low-high, and low-low frequency combinations. Finally, to isolate auditory continuity stimulus

pairings (i.e., whether an auditory tone was constant or beeping), we designated input units for constant-constant, constant-beeping, beeping-constant, and beeping-beeping. Altogether, across the four sensory domains, we obtained 16 different sensory stimulus pair input units. For a given trial, four units would be activated to simulate a sensory stimulus combination (one unit per sensory domain). For example, a single trial might observe red-red (color), vertical-horizontal (orientation), low-high (pitch), constant-beeping (continuity) stimulus combination. Thus, to simulate an entire trial including both context and sensory stimuli, 7/28 possible input units would be activated.

We constructed our ANN with two hidden layers containing 1280 units each. This choice was due to recent counterintuitive evidence suggesting that the learning dynamics of extremely high-dimensional ANNs (i.e., those with many network parameters to tune) naturally protect against overfitting, supporting generalized solutions[46]. Moreover, we found that across many initializations, the representational geometry identified in the ANN's hidden layer was highly replicable. Finally, our output layer contained four units, one for each motor response (corresponding to left middle, left index, right middle, right index finger presses).

The ANN transformed a 28-element input vector (representing a specific trial instance) into a 4-element response vector, and obeyed the equation

$$Y = f_s(X_{hidden2}W_{out} + b) \qquad (1)$$

where $Y$ corresponds to the 4-element response vector, $f_s$ is a sigmoid function, $W_{out}$ corresponds to the connectivity weight matrix between the hidden and output layer, $b$ is a bias term, and $X_{hidden2}$ is the activity vector of the 2nd hidden layer. $X_{hidden2}$ was obtained by the equation

$$X_{hidden2} = f_r((X_{hidden1} + I)W_{hidden} + b) \qquad (2)$$
$$X_{hidden1} = f_r((X_{input})W_{input} + b) \qquad (3)$$

Where $f_r$ is a rectified linear function (ReLU), $W_{hidden}$ is the connectivity matrix between the hidden layers, $X_{hidden1}$ corresponds to the 1st hidden layer activations that contain trial information, $X_{input}$ is the input layer, $W_{input}$ is the connectivity matrix between the input and 1st hidden layer, and $I$ is a noise vector sampled from a normal distribution with 0-mean and $\frac{1}{n}$ -variance, where $n$ refers to the number of hidden units.

## Identifying conjunctive representations: ANN training

The ANN was trained by minimizing the mean squared error between the network's outputs and the correct target output. The mean squared error was computed using a mini-batch approach, where each mini-batch comprised of 192 distinct trials. (Each of the 64 unique task contexts were presented three times (with randomly sampled stimuli) in each mini-batch. Training was optimized using Adam, a variant of stochastic gradient descent[47]. We used the default parameters in PyTorch (version 1.0.1), with a learning rate of 0.0001. Training was stopped when the last 1000 mini-batches achieved over 99.5% average accuracy on the task. This performance was achieved after roughly 10,000 mini-batches (or 1,920,000 trials). Weights and biases were initialized with a uniform distribution $U(-\sqrt{k}, \sqrt{k})$, where $k = \frac{1}{targets}$,

where 'targets' represents the number of units in the next layer. Note that no cross-validation was performed (nor was it necessary), since we were only interested in representational geometry of the hidden layer. We also note that the representational geometry we observed in the hidden layer was robust to different initializations and hyperparameter choices.

## Identifying conjunctive representations: ANN representational analysis

We extracted the representational geometry of the ANN's 2nd hidden layer using representational similarity analysis (RSA)[48]. This was done to understand how task rule and stimulus activations were transformed in the hidden layer. To extract the representational geometry of the hidden layer, we systematically activated a single unit in the input layer (which corresponded to either a task rule or sensory stimulus pair), and estimated the corresponding hidden layer activations (using trained connectivity weights). This resulted in a total of 28 (12 task rules and 16 sensory stimuli combinations) activation patterns. The representational similarity matrix (RSM) was obtained by computing the Pearson's correlation between the hidden layer activation patterns for all 28 conditions.

## Identifying conjunctive representations: fMRI analysis

We compared the representational geometry of the ANN's hidden layer to the representational geometry of each brain parcel. This was possible because we extracted the exact same set of activation patterns (e.g., activations for task rules and sensory stimuli) in empirical data as our ANN model, enabling a direct comparison of representations. The representational geometry was estimated as the representational similarity matrix (RSM) of all task rules and sensory stimuli conditions.

We first estimated the empirical RSMs for every brain parcel separately in the Glasser et al. (2016) atlas. This was done by comparing the activation patterns of each of the 28 task conditions using the vertices within each parcel (12 task rule activations, 16 sensory stimulus activations). We then applied a Fisher's $z$-transform on both the empirical and ANN's RSMs, and then estimated the Spearman's rank correlation between the Fisher's $z$-transformed ANN and empirical RSMs (using the upper triangle values only). This procedure was performed on the RSM of every brain parcel, providing a similarity score between each brain parcel's and the ANN's representational geometry. For our main analysis, we selected the top 10 parcels with highest similarity to the ANN's hidden layer. However, we also performed additional analyses using the top 20, 30, and 40 parcels.

## Inter-layer FC weight estimation

We estimated the inter-layer resting-state FC to identify weights between regions and layers in our empirical model. This was similar to a previously published approach which identified FC weights between pairs of brain regions[8]. This involved identifying FC weight mappings between the task rule input layer to the hidden layer, sensory stimulus input layer to the hidden layer, and the hidden layer to the motor output layer. For each inter-layer FC mapping, we estimated the vertex-to-vertex FC weights using principal components linear

regression. We used principal components regression because most layers had more vertices (i.e., predictors) than samples in our resting-state data (resting-state fMRI data contained 1065 TRs). For all inter-layer FC estimations, we used principal components regression with 500 components. Specifically, inter-layer weights were estimated by fitting principal components to the regression equation

$$Y = \beta_0 + \sum_i^{500} X_i \beta_i + \varepsilon \qquad (3)$$

where $Y$ corresponds to the t x n matrix with t time points and n vertices (i.e., the target vertices to be predicted), $\beta_0$ corresponds to a constant term, $\beta_i$ corresponds to the 1 x n matrix reflecting the mapping from the component time series onto the n target vertices, $X_i$ corresponds to the t x 1 component time series for component i, and $\varepsilon$ corresponds to the error in the regression model. Note that $X$ corresponds to the t x 500 component matrix obtained from a PCA on the resting-state data from the source layer. Also note that these loadings onto these 500 components are saved for later, when task activation patterns from a source layer are projected onto a target layer. The loadings project the original vertex-wise task activation patterns in the source layer onto a lower-dimensional space enabling faster computations. A similar approach was used in a previous study[53]. FC weights were computed for each individual separately, but then averaged across subjects to obtain a group inter-layer weight FC matrix.

Note that in some cases, it was possible for overlap between the source and target vertices. (For example, some hidden area vertices may have coincided with the same vertices in the context layer.) In these cases, these overlapping vertices were excluded in the set of predictors (i.e., removed from the source layer) in the regression model.

## Simulating sensorimotor transformations with multi-step activity flow mapping

We generated predictions of motor response activations (in motor cortex) by assessing the correct motor response given a specific task context and sensory stimulus activation pattern (for additional details see Supplementary Fig. 1). For each subject, we simulated 960 trials. This consisted of the 64 unique task contexts paired with 15 randomly sampled stimulus combinations. For a trial, the task context input activation pattern was obtained by extracting the activation vector for the logic, sensory, and motor rule, and computing the mean rule vector (i.e., additive compositionality). The sensory stimulus input activation pattern was obtained by extracting the relevant sensory stimulus activation pattern. (Note that for a given trial, we only extracted the activation pattern for the sensory feature of interest. For example, if the rule was "Red", only color activation patterns would be extracted, and all other stimulus activations would be set to 0.) Thus, the context and sensory stimulus activation patterns could be defined as

$$X_{context} = (R_{logic} + R_{sensory} + R_{motor})/3 \qquad (4)$$

$$X_{stimulus} = X_{sensory} \qquad (5)$$

where $X_{context}$ corresponds to the input activation pattern for task context, $R_{logic}$ corresponds to extracted logic rule activation pattern (e.g., "Both", "Not Both", "Either", or "Neither") obtained

from the task GLM, $R_{sensory}$ corresponds to the extracted sensory rule activation pattern from the task GLM, $R_{motor}$ corresponds to the extracted motor rule activation pattern from the task GLM, and $X_{stimulus}$ corresponds to the extracted sensory stimulus activation pattern that is indicated by the task context.

$X_{context}$ and $X_{stimulus}$ reflect the input activation patterns that were used to predict motor response conditions. Importantly, these input activation patterns were both spatially *and* representationally distinct from the motor response activations (in motor cortex). They were representationally distinct because these input activation patterns contained no information about the motor response required for a correct response. (In addition, we also used cross-validation to predict the motor response of a held-out subject, described below).

We used the inter-layer FC weight maps to project $X_{context}$ and $X_{stimulus}$ onto the hidden layer vertices. The projections (or predicted activation patterns on the hidden layer) were then thresholded to remove any negative BOLD predictions. This thresholding is equivalent to a rectified linear unit (ReLU), a commonly used nonlinear function in artificial neural networks[34]. Thus, the hidden layer was defined by

$$X_{hidden} = f_r(X_{context} W_{context2hidden} + X_{stimulus} W_{stimulus2hidden}) \qquad (6)$$

where $X_{hidden}$ corresponds to the predicted hidden layer activation pattern, $f_r$ is a ReLU function (i.e., $f_r(x) = max(x, 0)$ ), $W_{context2hidden}$ corresponds to the inter-layer resting-state FC weights between the context and hidden layer, and $W_{stimulus2hidden}$ corresponds to the inter-layer resting-state FC weights between the stimulus and hidden layer. Note that all inter-layer FC weights ($W_x$) were computed using a principal component regression with 500 components. This requires that the vertex-wise activation space (e.g., $X_{context}$) be projected onto component space such that we define

$$W_x = U \hat{W}_{pc} \qquad (7)$$

where $U$ corresponds a *m* x 500 matrix which maps the source layer's *m* vertices into component space, and $\hat{W}_{pc}$ is a 500 x *n* matrix that maps the components onto the target layer's *n* vertices. (Note that $\hat{W}_{pc}$ corresponds to the regression coefficients from equation 3., and that both $U$ and $\hat{W}_{pc}$ are estimated from resting-state data.) Thus, $W_x$ is an *m* x *n* transformation from a source layer's spatial pattern to a target layer's spatial pattern that is achieved through principal component regression on resting-state fMRI data.

Finally, we generated a predicted motor output response by computing

$$X_{output} = X_{hidden} W_{hidden2output} \qquad (8)$$

where $X_{output}$ corresponds to the predicted motor response (in motor cortex), and $W_{hidden2output}$ corresponds to the inter-layer resting-state FC weights between the hidden and output layer. The full model computation can thus be formalized as

$$X_{output} = f_r(X_{context} W_{context2hidden} + X_{stimulus} W_{stimulus2hidden}) W_{hidden2output} \qquad (9)$$

$X_{output}$ only yields a predicted activation pattern for the motor cortex for a given context and stimulus input activation pattern. To evaluate whether $X_{output}$ could successfully predict the *correct* motor response for a given trial, we constructed an ideal 'task solver' that would indicate the correct motor response on a given trial (Supplementary Fig. 1). This solver would then be

used to extract the correct motor response activation pattern, and compare the predicted motor cortex activation with the actual motor cortex activation pattern.

We simulated 960 trials per subject, randomly sampling context and stimulus input activation patterns. Because we sampled across the 64 task contexts equally (15 trials per context), the correct motor responses were equally balanced across 960 trials. Thus, of the 960 simulated trials for each subject, 240 trials yielded a left middle, left index, right middle, and right index response each. Each of these 240 predicted motor response patterns were subsequently averaged across trials such that we only obtained 4 predicted motor response patterns for each subject. Averaging was performed to remove any potential biases that a trial may have (e.g., a task context with the 'left middle' motor rule might be more biased towards a 'left middle' motor response).

## Statistical and permutation testing of predicted motor response activations

The simulated empirical model generated predicted activations of motor activations in motor cortex. However, the predictions would only be interesting if they resembled actual motor response activations directly estimated the response period via task GLM. In other words, without a ground truth reference to the actual motor response activation pattern, the predicted activation patterns would hold little meaning. The simulated empirical model generated four predicted activation patterns corresponding to predicted motor responses for each subject. We also had four *actual* activation patterns corresponding to motor responses that were extracted from the motor response period using a standard task GLM for each subject. To test whether the predicted activation patterns actually conformed to the actual motor response activation patterns, we trained a decoder on the predicted motor response activations and tested on the actual motor response activations of held-out subjects. We used the same cross-validation decoding scheme as before, with the exception that training was exclusively performed on predicted activation patterns of 88 subjects, while testing was exclusively performed on the actual activation patterns of 8 held-out subjects. Training a decoder on the predicted activations and decoding the actual activations made this analysis consistent with a prediction perspective – we could test if, in the absence of any motor task activation information, the ENN could predict actual motor response activation patterns that correspond to behavior. All other details (e.g., minimum-distance classifier, leave-8-subjects out cross-validation) remained the same.

Statistical significance was assessed using permutation tests. We permuted the labels of the predicted motor responses while testing on the actual motor responses. Null distributions are visualized in gray (Fig. 7h). Statistical significance was assessed by comparing the mean of the bootstrapped predicted-to-actual accuracy scores, and comparing them against a non-parametric p-value that was estimated from the null distribution. Statistical significance was defined by a $p < 0.05$ threshold.

## Code and data availability

All code and data related to this study will be made available on a public repository upon (or before) publication. In the interim, code and data are available on request.

# References

1. Cole, M. W., Braver, T. S. & Meiran, N. The task novelty paradox: Flexible control of inflexible neural pathways during rapid instructed task learning. *Neurosci. Biobehav. Rev.* **81**, 4–15 (2017).

2. Schneider, W. & Chein, J. M. Controlled & automatic processing: behavior, theory, and biological mechanisms. *Cogn. Sci.* **27**, 525–559 (2003).

3. Miller, E. K. & Cohen, J. D. An integrative theory of prefrontal cortex function. *Annu. Rev. Neurosci.* **24**, 167–202 (2001).

4. Kanwisher, N. Functional specificity in the human brain: A window into the functional architecture of the mind. *Proc. Natl. Acad. Sci.* **107**, 11163–11170 (2010).

5. Nishimoto, S. *et al.* Reconstructing visual experiences from brain activity evoked by natural movies. *Curr. Biol.* **21**, 1641–1646 (2011).

6. Yokoi, A. & Diedrichsen, J. Parcellation of motor sequence representations in the human neocortex. *bioRxiv* 419754 (2018) doi:10.1101/419754.

7. Cole, M. W., Ito, T. & Braver, T. S. The Behavioral Relevance of Task Information in Human Prefrontal Cortex. *Cereb. Cortex N. Y. N 1991* bhv072– (2015) doi:10.1093/cercor/bhv072.

8. Ito, T. *et al.* Cognitive task information is transferred between brain regions via resting-state network topology. *Nat. Commun.* (2017) doi:10.1038/s41467-017-01000-w.

9. Rigotti, M. *et al.* The importance of mixed selectivity in complex cognitive tasks. *Nature* **497**, 585–90 (2013).

10. Reverberi, C., Görgen, K. & Haynes, J.-D. Compositionality of rule representations in human prefrontal cortex. *Cereb. Cortex* **22**, 1237–1246 (2012).

11. Kikumoto, A. & Mayr, U. Conjunctive representations that integrate stimuli, responses, and

rules are critical for action selection. *Proc. Natl. Acad. Sci.* (2020)

doi:10.1073/pnas.1922166117.

12. Yang, G. R., Joglekar, M. R., Song, H. F., Newsome, W. T. & Wang, X.-J. Task

representations in neural networks trained to perform many cognitive tasks. *Nat. Neurosci.* 1

(2019) doi:10.1038/s41593-018-0310-2.

13. Mante, V., Sussillo, D., Shenoy, K. V. & Newsome, W. T. Context-dependent computation

by recurrent dynamics in prefrontal cortex. *Nature* **503**, 78–84 (2013).

14. Cole, M. W. *et al.* Multi-task connectivity reveals flexible hubs for adaptive task control. *Nat.*

*Neurosci.* **16**, 1348–1355 (2013).

15. Cocuzza, C. V., Ito, T., Schultz, D., Bassett, D. S. & Cole, M. W. Flexible coordinator and

switcher hubs for adaptive task control. *J. Neurosci.* (2020)

doi:10.1523/JNEUROSCI.2559-19.2020.

16. Cole, M. W., Ito, T., Bassett, D. S. & Schultz, D. H. Activity flow over resting-state networks

shapes cognitive task activations. *Nat. Neurosci.* (2016) doi:10.1038/nn.4406.

17. Song, H. F., Yang, G. R. & Wang, X.-J. Training Excitatory-Inhibitory Recurrent Neural

Networks for Cognitive Tasks: A Simple and Flexible Framework. *PLOS Comput. Biol.* **12**,

e1004792 (2016).

18. Yamins, D. L. K. *et al.* Performance-optimized hierarchical models predict neural responses

in higher visual cortex. *Proc. Natl. Acad. Sci.* **111**, 8619–8624 (2014).

19. Khaligh-Razavi, S. M. & Kriegeskorte, N. Deep Supervised, but Not Unsupervised, Models

May Explain IT Cortical Representation. *PLoS Comput. Biol.* **10**, (2014).

20. Bashivan, P., Kar, K. & DiCarlo, J. J. Neural population control via deep image synthesis.

*Science* **364**, eaav9436 (2019).

21. Norman, K. A., Polyn, S. M., Detre, G. J. & Haxby, J. V. Beyond mind-reading: multi-voxel

pattern analysis of fMRI data. *Trends Cogn. Sci.* **10**, 424–30 (2006).

22. Mur, M., Bandettini, P. A. & Kriegeskorte, N. Revealing representational content with pattern-information fMRI—an introductory guide. *Soc. Cogn. Affect. Neurosci.* **4**, 101–109 (2009).

23. Glasser, M. F. *et al.* A multi-modal parcellation of human cerebral cortex. *Nature* 1–11 (2016) doi:10.1038/nature18933.

24. Ji, J. L. *et al.* Mapping the human brain's cortical-subcortical functional network organization. *NeuroImage* **185**, 35–57 (2019).

25. Power, J. D. & Petersen, S. E. Control-related systems in the human brain. *Curr. Opin. Neurobiol.* **23**, 223–228 (2013).

26. Cohen, J. D., Dunbar, K. & McClelland, J. L. On the control of automatic processes: a parallel distributed processing account of the Stroop effect. *Psychol. Rev.* **97**, 332–61 (1990).

27. Cohen, J. D., Aston-Jones, G. & Gilzenrat, M. S. A Systems-Level Perspective on Attention and Cognitive Control: Guided Activation, Adaptive Gating, Conflict Monitoring, and Exploitation versus Exploration. in *Cognitive neuroscience of attention* 71–90 (The Guilford Press, 2004).

28. Ito, T., Hearne, L., Mill, R., Cocuzza, C. & Cole, M. W. Discovering the Computational Relevance of Brain Network Organization. *Trends Cogn. Sci.* (2019) doi:10.1016/j.tics.2019.10.005.

29. De-Wit, L., Alexander, D., Ekroll, V. & Wagemans, J. Is neuroimaging measuring information in the brain? *Psychon. Bull. Rev.* 1–14 (2016) doi:10.3758/s13423-016-1002-0.

30. Brette, R. Is coding a relevant metaphor for the brain? *Behav. Brain Sci.* 1–44 (2019).

31. Dosenbach, N. U. F. *et al.* Distinct brain networks for adaptive and stable task control in

humans. *Proc. Natl. Acad. Sci. U. S. A.* **104**, 11073–11078 (2007).

32. Waskom, M. L., Kumaran, D., Gordon, A. M., Rissman, J. & Wagner, A. D. Frontoparietal representations of task context support the flexible control of goal-directed cognition. *J. Neurosci. Off. J. Soc. Neurosci.* **34**, 10743–55 (2014).

33. Silver, D. *et al.* Mastering the game of Go without human knowledge. *Nature* **550**, 354–359 (2017).

34. Yang, G. R. & Wang, X.-J. Artificial neural networks for neuroscientists: A primer. *ArXiv200601001 Cs Q-Bio* (2020).

35. Hagmann, P. *et al.* Mapping the structural core of human cerebral cortex. *PLoS Biol.* **6**, 1479–1493 (2008).

36. Deco, G. *et al.* Resting-state functional connectivity emerges from structurally and dynamically shaped slow linear fluctuations. *J. Neurosci. Off. J. Soc. Neurosci.* **33**, 11239–52 (2013).

37. Wang, P. *et al.* Inversion of a large-scale circuit model reveals a cortical hierarchy in the dynamic resting human brain. *Sci. Adv.* **5**, eaat7854 (2019).

38. Demirtaş, M. *et al.* Hierarchical Heterogeneity across Human Cortex Shapes Large-Scale Neural Dynamics. *Neuron* **101**, 1181-1194.e13 (2019).

39. Tschopp, F. D., Reiser, M. B. & Turaga, S. C. A Connectome Based Hexagonal Lattice Convolutional Network Model of the Drosophila Visual System. *ArXiv180604793 Cs Q-Bio* (2018).

40. Ito, T., Hearne, L., Mill, R., Cocuzza, C. & Cole, M. W. Discovering the Computational Relevance of Brain Network Organization. *Trends Cogn. Sci.* **24**, 25–38 (2020).

41. Litwin-Kumar, A. & Turaga, S. C. Constraining computational models using electron microscopy wiring diagrams. *Curr. Opin. Neurobiol.* **58**, 94–100 (2019).

42. Wu, Y., Zhang, S., Zhang, Y., Bengio, Y. & Salakhutdinov, R. R. On Multiplicative Integration with Recurrent Neural Networks. in *Advances in Neural Information Processing Systems 29* (eds. Lee, D. D., Sugiyama, M., Luxburg, U. V., Guyon, I. & Garnett, R.) 2856–2864 (Curran Associates, Inc., 2016).

43. Schultz, D. H. *et al.* Global connectivity of the fronto-parietal cognitive control network is related to depression symptoms in the general population. *Netw. Neurosci.* **3**, 107–123 (2019).

44. Cole, M. W., Bagic, A., Kass, R. & Schneider, W. Prefrontal Dynamics Underlying Rapid Instructed Task Learning Reverse with Practice. *J. Neurosci.* **30**, 14245–14254 (2010).

45. Schneider, W., Eschman, A. & Zuccolotto, A. *E-Prime: User's guide*. (Psychology Software Incorporated, 2002).

46. Advani, M. S. & Saxe, A. M. High-dimensional dynamics of generalization error in neural networks. *ArXiv171003667 Phys. Q-Bio Stat* (2017).

47. Kingma, D. P. & Ba, J. Adam: A Method for Stochastic Optimization. *ArXiv14126980 Cs* (2017).

48. Kriegeskorte, N., Mur, M. & Bandettini, P. Representational similarity analysis - connecting the branches of systems neuroscience. *Front. Syst. Neurosci.* **2**, 4 (2008).

49. Glasser, M. F. *et al.* The Human Connectome Project's neuroimaging approach. *Nat. Neurosci.* **19**, 1175–87 (2016).

50. Ito, T. *et al.* Task-evoked activity quenches neural correlations and variability across cortical areas. *PLOS Comput. Biol.* **16**, e1007983 (2020).

51. Ciric, R. *et al.* Benchmarking of participant-level confound regression strategies for the control of motion artifact in studies of functional connectivity. *NeuroImage* **154**, 174–187 (2017).

52. Behzadi, Y., Restom, K., Liau, J. & Liu, T. T. A component based noise correction method (CompCor) for BOLD and perfusion based fMRI. *NeuroImage* **37**, 90–101 (2007).

53. Anzellotti, S., Fedorenko, E., Caramazza, A. & Saxe, R. Measuring and Modeling Transformations of Information Between Brain Regions with fMRI. *bioRxiv* 1–13 (2016) doi:10.1101/074856.

# Supplementary Figures



**Supplementary Figure 1. Flow chart describing neural network simulations with empirical data via activity flow mapping.** We generate a subject's predicted motor response activations using only task rule and sensory stimulus activation patterns as inputs. We then test these predictions against the actual motor response activations of held-out subjects.

**Supplementary Figure 2. Network affiliations of conjunction hubs and the task rule input layer using a previously defined multimodal atlas and network partition**[23,24]**. a)** The network affiliations of the 10 conjunction hub brain areas. **b)** Network affiliations of the 228 brain regions that contained decodable task rule information.

**Supplementary Figure 3. Example of task GLM approach to obtain task activation estimates. a)** An example miniblock containing one encoding block (task rule set) and three trials. Note that while stimulus presentation and response periods overlap, they are not collinear. **b)** The regressors for the relevant task conditions in the example miniblock. We obtain regressors (estimated across all 128 miniblocks) for all task rule, sensory stimuli, and motor response conditions. Altogether there are 32 different task conditions (12 task rules, 16 sensory stimuli pairs, and four motor response periods). Note that task rule regressors (logic, sensory, and motor rule examples) appear collinear in this example, but that across all 128 miniblocks task rule conditions are properly counterbalanced to avoid collinearity. Regressors shown here are illustrated without convolution with SPM's canonical HRF.

**a**



Supplementary Figure 4. A task GLM design matrix for an example subject.

**Supplementary Table 1. Regions containing decodable color (red/blue) stimulus information.**

| Label | GlasserID | Network Affiliation | Hemisphere |
|---|---|---|---|
| Visual2-54_L-Ctx | L_VVC_ROI | Visual2 | L |
| Visual2-05_R-Ctx | R_V4_ROI | Visual2 | R |
| Visual1-03_R-Ctx | R_DVT_ROI | Visual1 | R |
| Visual2-22_R-Ctx | R_V4t_ROI | Visual2 | R |

**Supplementary Table 2. Regions containing decodable orientation (vertical/horizontal) stimulus information.**

| Label | GlasserID | Network Affiliation | Hemisphere |
|---|---|---|---|
| Visual1-04_L-Ctx | L_V1_ROI | Visual1 | L |
| Visual2-28_L-Ctx | L_MST_ROI | Visual2 | L |
| Visual2-30_L-Ctx | L_V2_ROI | Visual2 | L |
| Visual2-31_L-Ctx | L_V3_ROI | Visual2 | L |
| Visual2-32_L-Ctx | L_V4_ROI | Visual2 | L |
| Visual2-33_L-Ctx | L_V8_ROI | Visual2 | L |
| Visual2-35_L-Ctx | L_V7_ROI | Visual2 | L |
| Visual2-40_L-Ctx | L_LO2_ROI | Visual2 | L |
| Visual2-41_L-Ctx | L_PIT_ROI | Visual2 | L |
| Visual2-42_L-Ctx | L_MT_ROI | Visual2 | L |
| Visual2-51_L-Ctx | L_V3CD_ROI | Visual2 | L |
| Visual1-01_R-Ctx | R_V1_ROI | Visual1 | R |
| Visual2-03_R-Ctx | R_V2_ROI | Visual2 | R |
| Visual2-04_R-Ctx | R_V3_ROI | Visual2 | R |
| Visual2-05_R-Ctx | R_V4_ROI | Visual2 | R |
| Visual2-06_R-Ctx | R_V8_ROI | Visual2 | R |
| Visual2-09_R-Ctx | R_IPS1_ROI | Visual2 | R |
| Visual2-12_R-Ctx | R_LO1_ROI | Visual2 | R |

**Supplementary Table 3. Regions containing decodable pitch (high/low) stimulus information.**

| Label | GlasserID | Network Affiliation | Hemisphere |
|---|---|---|---|
| Auditory-08_L-Ctx | L_A1_ROI | Auditory | L |
| Auditory-09_L-Ctx | L_52_ROI | Auditory | L |
| Auditory-10_L-Ctx | L_RI_ROI | Auditory | L |
| Auditory-11_L-Ctx | L_TA2_ROI | Auditory | L |
| Auditory-12_L-Ctx | L_PBelt_ROI | Auditory | L |
| Auditory-13_L-Ctx | L_MBelt_ROI | Auditory | L |
| Auditory-14_L-Ctx | L_LBelt_ROI | Auditory | L |
| Auditory-15_L-Ctx | L_A4_ROI | Auditory | L |
| Auditory-01_R-Ctx | R_A1_ROI | Auditory | R |
| Auditory-02_R-Ctx | R_52_ROI | Auditory | R |
| Auditory-03_R-Ctx | R_TA2_ROI | Auditory | R |
| Auditory-04_R-Ctx | R_PBelt_ROI | Auditory | R |
| Auditory-05_R-Ctx | R_MBelt_ROI | Auditory | R |
| Auditory-06_R-Ctx | R_LBelt_ROI | Auditory | R |
| Auditory-07_R-Ctx | R_A4_ROI | Auditory | R |

**Supplementary Table 4. Regions containing decodable constant/beeping stimulus information.**

| Label | GlasserID | Network Affiliation | Hemisphere |
|---|---|---|---|
| Auditory-08_L-Ctx | L_A1_ROI | Auditory | L |
| Auditory-09_L-Ctx | L_52_ROI | Auditory | L |
| Auditory-10_L-Ctx | L_RI_ROI | Auditory | L |
| Auditory-11_L-Ctx | L_TA2_ROI | Auditory | L |
| Auditory-12_L-Ctx | L_PBelt_ROI | Auditory | L |
| Auditory-13_L-Ctx | L_MBelt_ROI | Auditory | L |
| Auditory-14_L-Ctx | L_LBelt_ROI | Auditory | L |
| Auditory-15_L-Ctx | L_A4_ROI | Auditory | L |
| Auditory-01_R-Ctx | R_A1_ROI | Auditory | R |
| Auditory-02_R-Ctx | R_52_ROI | Auditory | R |
| Auditory-03_R-Ctx | R_TA2_ROI | Auditory | R |
| Auditory-04_R-Ctx | R_PBelt_ROI | Auditory | R |
| Auditory-05_R-Ctx | R_MBelt_ROI | Auditory | R |
| Auditory-06_R-Ctx | R_LBelt_ROI | Auditory | R |
| Auditory-07_R-Ctx | R_A4_ROI | Auditory | R |

**Supplementary Table 5. Conjunction hubs.**

| Label | GlasserID | Network Affiliation | Hemisphere |
|---|---|---|---|
| Cingulo-Opercular-49_L-Ctx | L_FOP3_ROI | Cingulo-Opercular | L |
| Somatomotor-15_R-Ctx | R_OP4_ROI | Somatomotor | R |
| Dorsal-Attention-18_L-Ctx | L_AIP_ROI | Dorsal-Attention | L |
| Cingulo-Opercular-04_R-Ctx | R_5mv_ROI | Cingulo-Opercular | R |
| Cingulo-Opercular-22_R-Ctx | R_FOP3_ROI | Cingulo-Opercular | R |
| Somatomotor-08_R-Ctx | R_7PC_ROI | Somatomotor | R |
| Cingulo-Opercular-44_L-Ctx | L_PFcm_ROI | Cingulo-Opercular | L |
| Frontoparietal-34_L-Ctx | L_a47r_ROI | Frontoparietal | L |
| Frontoparietal-03_R-Ctx | R_7Pm_ROI | Frontoparietal | R |
| Somatomotor-07_R-Ctx | R_7AL_ROI | Somatomotor | R |

**Supplementary Table 6. Task rule (input) regions.**

| Label | GlasserID | Network Affiliation | Hemisphere |
|---|---|---|---|
| Visual1-04_L-Ctx | L_V1_ROI | Visual1 | L |
| Visual2-30_L-Ctx | L_V2_ROI | Visual2 | L |
| Visual2-31_L-Ctx | L_V3_ROI | Visual2 | L |
| Visual2-32_L-Ctx | L_V4_ROI | Visual2 | L |
| Somatomotor-21_L-Ctx | L_4_ROI | Somatomotor | L |
| Somatomotor-22_L-Ctx | L_3b_ROI | Somatomotor | L |
| Cingulo-Opercular-30_L-Ctx | L_FEF_ROI | Cingulo-Opercular | L |
| Dorsal-Attention-12_L-Ctx | L_PEF_ROI | Dorsal-Attention | L |
| Language-10_L-Ctx | L_55b_ROI | Language | L |
| Frontoparietal-29_L-Ctx | L_RSC_ROI | Frontoparietal | L |
| Frontoparietal-30_L-Ctx | L_POS2_ROI | Frontoparietal | L |
| Visual2-35_L-Ctx | L_V7_ROI | Visual2 | L |
| Visual2-36_L-Ctx | L_IPS1_ROI | Visual2 | L |
| Visual2-38_L-Ctx | L_V3B_ROI | Visual2 | L |
| Visual2-42_L-Ctx | L_MT_ROI | Visual2 | L |
| Auditory-08_L-Ctx | L_A1_ROI | Auditory | L |
| Language-11_L-Ctx | L_PSL_ROI | Language | L |
| Language-12_L-Ctx | L_SFL_ROI | Language | L |
| Posterior-Multimodal-05_L-Ctx | L_PCV_ROI | Posterior-Multimodal | L |
| Default-38_L-Ctx | L_7m_ROI | Default | L |
| Default-39_L-Ctx | L_POS1_ROI | Default | L |
| Default-40_L-Ctx | L_23d_ROI | Default | L |
| Default-41_L-Ctx | L_v23ab_ROI | Default | L |
| Default-42_L-Ctx | L_d23ab_ROI | Default | L |
| Cingulo-Opercular-32_L-Ctx | L_23c_ROI | Cingulo-Opercular | L |
| Somatomotor-25_L-Ctx | L_24dd_ROI | Somatomotor | L |
| Somatomotor-26_L-Ctx | L_24dv_ROI | Somatomotor | L |
| Somatomotor-27_L-Ctx | L_7AL_ROI | Somatomotor | L |
| Cingulo-Opercular-33_L-Ctx | L_SCEF_ROI | Cingulo-Opercular | L |
| Cingulo-Opercular-34_L-Ctx | L_6ma_ROI | Cingulo-Opercular | L |
| Cingulo-Opercular-35_L-Ctx | L_7Am_ROI | Cingulo-Opercular | L |
| Somatomotor-28_L-Ctx | L_7PC_ROI | Somatomotor | L |
| Visual2-43_L-Ctx | L_LIPv_ROI | Visual2 | L |
| Visual2-44_L-Ctx | L_VIP_ROI | Visual2 | L |
| Somatomotor-29_L-Ctx | L_1_ROI | Somatomotor | L |
| Somatomotor-30_L-Ctx | L_2_ROI | Somatomotor | L |
| Somatomotor-31_L-Ctx | L_3a_ROI | Somatomotor | L |
| Somatomotor-32_L-Ctx | L_6d_ROI | Somatomotor | L |
| Somatomotor-33_L-Ctx | L_6mp_ROI | Somatomotor | L |
| Somatomotor-34_L-Ctx | L_6v_ROI | Somatomotor | L |
| Cingulo-Opercular-36_L-Ctx | L_p24pr_ROI | Cingulo-Opercular | L |

| | | | |
|---|---|---|---|
| Cingulo-Opercular-37_L-Ctx | L_33pr_ROI | Cingulo-Opercular | L |
| Cingulo-Opercular-38_L-Ctx | L_a24pr_ROI | Cingulo-Opercular | L |
| Cingulo-Opercular-39_L-Ctx | L_p32pr_ROI | Cingulo-Opercular | L |
| Default-44_L-Ctx | L_a24_ROI | Default | L |
| Default-45_L-Ctx | L_d32_ROI | Default | L |
| Frontoparietal-32_L-Ctx | L_8BM_ROI | Frontoparietal | L |
| Default-47_L-Ctx | L_10r_ROI | Default | L |
| Default-49_L-Ctx | L_8Av_ROI | Default | L |
| Default-51_L-Ctx | L_9m_ROI | Default | L |
| Default-52_L-Ctx | L_8BL_ROI | Default | L |
| Default-54_L-Ctx | L_10d_ROI | Default | L |
| Frontoparietal-33_L-Ctx | L_8C_ROI | Frontoparietal | L |
| Language-14_L-Ctx | L_44_ROI | Language | L |
| Language-15_L-Ctx | L_45_ROI | Language | L |
| Default-55_L-Ctx | L_47l_ROI | Default | L |
| Frontoparietal-34_L-Ctx | L_a47r_ROI | Frontoparietal | L |
| Cingulo-Opercular-40_L-Ctx | L_6r_ROI | Cingulo-Opercular | L |
| Language-16_L-Ctx | L_IFJa_ROI | Language | L |
| Frontoparietal-35_L-Ctx | L_IFJp_ROI | Frontoparietal | L |
| Language-17_L-Ctx | L_IFSp_ROI | Language | L |
| Frontoparietal-36_L-Ctx | L_IFSa_ROI | Frontoparietal | L |
| Frontoparietal-37_L-Ctx | L_p9-46v_ROI | Frontoparietal | L |
| Cingulo-Opercular-42_L-Ctx | L_9-46d_ROI | Cingulo-Opercular | L |
| Default-56_L-Ctx | L_9a_ROI | Default | L |
| Frontoparietal-39_L-Ctx | L_a10p_ROI | Frontoparietal | L |
| Frontoparietal-40_L-Ctx | L_11l_ROI | Frontoparietal | L |
| Dorsal-Attention-15_L-Ctx | L_LIPd_ROI | Dorsal-Attention | L |
| Dorsal-Attention-16_L-Ctx | L_6a_ROI | Dorsal-Attention | L |
| Frontoparietal-42_L-Ctx | L_i6-8_ROI | Frontoparietal | L |
| Cingulo-Opercular-43_L-Ctx | L_43_ROI | Cingulo-Opercular | L |
| Somatomotor-35_L-Ctx | L_OP4_ROI | Somatomotor | L |
| Somatomotor-36_L-Ctx | L_OP1_ROI | Somatomotor | L |
| Somatomotor-37_L-Ctx | L_OP2-3_ROI | Somatomotor | L |
| Auditory-09_L-Ctx | L_52_ROI | Auditory | L |
| Auditory-10_L-Ctx | L_RI_ROI | Auditory | L |
| Cingulo-Opercular-44_L-Ctx | L_PFcm_ROI | Cingulo-Opercular | L |
| Cingulo-Opercular-45_L-Ctx | L_PoI2_ROI | Cingulo-Opercular | L |
| Auditory-11_L-Ctx | L_TA2_ROI | Auditory | L |
| Cingulo-Opercular-46_L-Ctx | L_FOP4_ROI | Cingulo-Opercular | L |
| Cingulo-Opercular-47_L-Ctx | L_MI_ROI | Cingulo-Opercular | L |
| Frontoparietal-44_L-Ctx | L_AVI_ROI | Frontoparietal | L |
| Orbito-Affective-05_L-Ctx | L_AAIC_ROI | Orbito-Affective | L |
| Cingulo-Opercular-48_L-Ctx | L_FOP1_ROI | Cingulo-Opercular | L |

| Cingulo-Opercular-49_L-Ctx | L_FOP3_ROI | Cingulo-Opercular | L |
|---|---|---|---|
| Somatomotor-38_L-Ctx | L_FOP2_ROI | Somatomotor | L |
| Dorsal-Attention-17_L-Ctx | L_PFt_ROI | Dorsal-Attention | L |
| Dorsal-Attention-18_L-Ctx | L_AIP_ROI | Dorsal-Attention | L |
| Default-62_L-Ctx | L_PreS_ROI | Default | L |
| Language-18_L-Ctx | L_STGa_ROI | Language | L |
| Language-19_L-Ctx | L_A5_ROI | Language | L |
| Dorsal-Attention-19_L-Ctx | L_PHA3_ROI | Dorsal-Attention | L |
| Language-21_L-Ctx | L_STSdp_ROI | Language | L |
| Default-65_L-Ctx | L_STSvp_ROI | Default | L |
| Frontoparietal-45_L-Ctx | L_TE1p_ROI | Frontoparietal | L |
| Dorsal-Attention-20_L-Ctx | L_TE2p_ROI | Dorsal-Attention | L |
| Dorsal-Attention-21_L-Ctx | L_PHT_ROI | Dorsal-Attention | L |
| Visual2-45_L-Ctx | L_PH_ROI | Visual2 | L |
| Language-22_L-Ctx | L_TPOJ1_ROI | Language | L |
| Posterior-Multimodal-06_L-Ctx | L_TPOJ2_ROI | Posterior-Multimodal | L |
| Visual1-06_L-Ctx | L_DVT_ROI | Visual1 | L |
| Dorsal-Attention-22_L-Ctx | L_PGp_ROI | Dorsal-Attention | L |
| Frontoparietal-47_L-Ctx | L_IP1_ROI | Frontoparietal | L |
| Dorsal-Attention-23_L-Ctx | L_IP0_ROI | Dorsal-Attention | L |
| Cingulo-Opercular-50_L-Ctx | L_PFop_ROI | Cingulo-Opercular | L |
| Cingulo-Opercular-51_L-Ctx | L_PF_ROI | Cingulo-Opercular | L |
| Frontoparietal-48_L-Ctx | L_PFm_ROI | Frontoparietal | L |
| Default-69_L-Ctx | L_PGi_ROI | Default | L |
| Default-70_L-Ctx | L_PGs_ROI | Default | L |
| Visual2-46_L-Ctx | L_V6A_ROI | Visual2 | L |
| Default-71_L-Ctx | L_PHA2_ROI | Default | L |
| Default-73_L-Ctx | L_31a_ROI | Default | L |
| Visual2-54_L-Ctx | L_VVC_ROI | Visual2 | L |
| Cingulo-Opercular-52_L-Ctx | L_Pol1_ROI | Cingulo-Opercular | L |
| Somatomotor-39_L-Ctx | L_Ig_ROI | Somatomotor | L |
| Cingulo-Opercular-53_L-Ctx | L_FOP5_ROI | Cingulo-Opercular | L |
| Frontoparietal-50_L-Ctx | L_p47r_ROI | Frontoparietal | L |
| Auditory-14_L-Ctx | L_LBelt_ROI | Auditory | L |
| Auditory-15_L-Ctx | L_A4_ROI | Auditory | L |
| Default-77_L-Ctx | L_TE1m_ROI | Default | L |
| Cingulo-Opercular-55_L-Ctx | L_a32pr_ROI | Cingulo-Opercular | L |
| Cingulo-Opercular-56_L-Ctx | L_p24_ROI | Cingulo-Opercular | L |
| Visual1-01_R-Ctx | R_V1_ROI | Visual1 | R |
| Visual2-02_R-Ctx | R_V6_ROI | Visual2 | R |
| Visual2-03_R-Ctx | R_V2_ROI | Visual2 | R |
| Visual2-04_R-Ctx | R_V3_ROI | Visual2 | R |
| Visual2-05_R-Ctx | R_V4_ROI | Visual2 | R |

| Somatomotor-01_R-Ctx | R_4_ROI | Somatomotor | R |
|---|---|---|---|
| Somatomotor-02_R-Ctx | R_3b_ROI | Somatomotor | R |
| Cingulo-Opercular-01_R-Ctx | R_FEF_ROI | Cingulo-Opercular | R |
| Cingulo-Opercular-02_R-Ctx | R_PEF_ROI | Cingulo-Opercular | R |
| Frontoparietal-01_R-Ctx | R_RSC_ROI | Frontoparietal | R |
| Frontoparietal-02_R-Ctx | R_POS2_ROI | Frontoparietal | R |
| Visual2-08_R-Ctx | R_V7_ROI | Visual2 | R |
| Visual2-09_R-Ctx | R_IPS1_ROI | Visual2 | R |
| Visual2-11_R-Ctx | R_V3B_ROI | Visual2 | R |
| Visual2-12_R-Ctx | R_LO1_ROI | Visual2 | R |
| Cingulo-Opercular-03_R-Ctx | R_PSL_ROI | Cingulo-Opercular | R |
| Posterior-Multimodal-01_R-Ctx | R_PCV_ROI | Posterior-Multimodal | R |
| Default-01_R-Ctx | R_7m_ROI | Default | R |
| Default-02_R-Ctx | R_POS1_ROI | Default | R |
| Default-03_R-Ctx | R_23d_ROI | Default | R |
| Default-04_R-Ctx | R_v23ab_ROI | Default | R |
| Default-05_R-Ctx | R_d23ab_ROI | Default | R |
| Somatomotor-03_R-Ctx | R_5m_ROI | Somatomotor | R |
| Cingulo-Opercular-04_R-Ctx | R_5mv_ROI | Cingulo-Opercular | R |
| Cingulo-Opercular-05_R-Ctx | R_23c_ROI | Cingulo-Opercular | R |
| Somatomotor-04_R-Ctx | R_5L_ROI | Somatomotor | R |
| Somatomotor-05_R-Ctx | R_24dd_ROI | Somatomotor | R |
| Somatomotor-06_R-Ctx | R_24dv_ROI | Somatomotor | R |
| Somatomotor-07_R-Ctx | R_7AL_ROI | Somatomotor | R |
| Cingulo-Opercular-06_R-Ctx | R_SCEF_ROI | Cingulo-Opercular | R |
| Cingulo-Opercular-07_R-Ctx | R_6ma_ROI | Cingulo-Opercular | R |
| Cingulo-Opercular-08_R-Ctx | R_7Am_ROI | Cingulo-Opercular | R |
| Somatomotor-08_R-Ctx | R_7PC_ROI | Somatomotor | R |
| Visual2-16_R-Ctx | R_LIPv_ROI | Visual2 | R |
| Visual2-17_R-Ctx | R_VIP_ROI | Visual2 | R |
| Dorsal-Attention-02_R-Ctx | R_MIP_ROI | Dorsal-Attention | R |
| Somatomotor-09_R-Ctx | R_1_ROI | Somatomotor | R |
| Somatomotor-10_R-Ctx | R_2_ROI | Somatomotor | R |
| Somatomotor-11_R-Ctx | R_3a_ROI | Somatomotor | R |
| Somatomotor-12_R-Ctx | R_6d_ROI | Somatomotor | R |
| Somatomotor-13_R-Ctx | R_6mp_ROI | Somatomotor | R |
| Somatomotor-14_R-Ctx | R_6v_ROI | Somatomotor | R |
| Cingulo-Opercular-09_R-Ctx | R_p24pr_ROI | Cingulo-Opercular | R |
| Cingulo-Opercular-10_R-Ctx | R_a24pr_ROI | Cingulo-Opercular | R |
| Cingulo-Opercular-11_R-Ctx | R_p32pr_ROI | Cingulo-Opercular | R |
| Frontoparietal-05_R-Ctx | R_d32_ROI | Frontoparietal | R |
| Frontoparietal-06_R-Ctx | R_8BM_ROI | Frontoparietal | R |
| Default-11_R-Ctx | R_8Av_ROI | Default | R |

| Default-12_R-Ctx | R_8Ad_ROI | Default | R |
|---|---|---|---|
| Default-13_R-Ctx | R_9m_ROI | Default | R |
| Default-14_R-Ctx | R_8BL_ROI | Default | R |
| Frontoparietal-07_R-Ctx | R_8C_ROI | Frontoparietal | R |
| Frontoparietal-08_R-Ctx | R_44_ROI | Frontoparietal | R |
| Default-17_R-Ctx | R_47l_ROI | Default | R |
| Cingulo-Opercular-12_R-Ctx | R_6r_ROI | Cingulo-Opercular | R |
| Language-04_R-Ctx | R_IFJa_ROI | Language | R |
| Frontoparietal-11_R-Ctx | R_IFSp_ROI | Frontoparietal | R |
| Cingulo-Opercular-13_R-Ctx | R_IFSa_ROI | Cingulo-Opercular | R |
| Frontoparietal-12_R-Ctx | R_p9-46v_ROI | Frontoparietal | R |
| Cingulo-Opercular-15_R-Ctx | R_9-46d_ROI | Cingulo-Opercular | R |
| Dorsal-Attention-04_R-Ctx | R_6a_ROI | Dorsal-Attention | R |
| Cingulo-Opercular-16_R-Ctx | R_43_ROI | Cingulo-Opercular | R |
| Somatomotor-15_R-Ctx | R_OP4_ROI | Somatomotor | R |
| Somatomotor-16_R-Ctx | R_OP1_ROI | Somatomotor | R |
| Auditory-02_R-Ctx | R_52_ROI | Auditory | R |
| Cingulo-Opercular-17_R-Ctx | R_PFcm_ROI | Cingulo-Opercular | R |
| Cingulo-Opercular-18_R-Ctx | R_PoI2_ROI | Cingulo-Opercular | R |
| Cingulo-Opercular-19_R-Ctx | R_FOP4_ROI | Cingulo-Opercular | R |
| Cingulo-Opercular-20_R-Ctx | R_MI_ROI | Cingulo-Opercular | R |
| Frontoparietal-20_R-Ctx | R_AVI_ROI | Frontoparietal | R |
| Orbito-Affective-02_R-Ctx | R_AAIC_ROI | Orbito-Affective | R |
| Cingulo-Opercular-21_R-Ctx | R_FOP1_ROI | Cingulo-Opercular | R |
| Cingulo-Opercular-22_R-Ctx | R_FOP3_ROI | Cingulo-Opercular | R |
| Somatomotor-19_R-Ctx | R_FOP2_ROI | Somatomotor | R |
| Dorsal-Attention-05_R-Ctx | R_PFt_ROI | Dorsal-Attention | R |
| Dorsal-Attention-06_R-Ctx | R_AIP_ROI | Dorsal-Attention | R |
| Default-23_R-Ctx | R_PreS_ROI | Default | R |
| Default-24_R-Ctx | R_H_ROI | Default | R |
| Language-06_R-Ctx | R_A5_ROI | Language | R |
| Language-07_R-Ctx | R_STSdp_ROI | Language | R |
| Default-27_R-Ctx | R_STSvp_ROI | Default | R |
| Frontoparietal-21_R-Ctx | R_TE1p_ROI | Frontoparietal | R |
| Dorsal-Attention-09_R-Ctx | R_PHT_ROI | Dorsal-Attention | R |
| Visual2-18_R-Ctx | R_PH_ROI | Visual2 | R |
| Language-08_R-Ctx | R_TPOJ1_ROI | Language | R |
| Posterior-Multimodal-03_R-Ctx | R_TPOJ2_ROI | Posterior-Multimodal | R |
| Posterior-Multimodal-04_R-Ctx | R_TPOJ3_ROI | Posterior-Multimodal | R |
| Visual1-03_R-Ctx | R_DVT_ROI | Visual1 | R |
| Dorsal-Attention-10_R-Ctx | R_PGp_ROI | Dorsal-Attention | R |
| Frontoparietal-23_R-Ctx | R_IP1_ROI | Frontoparietal | R |
| Dorsal-Attention-11_R-Ctx | R_IP0_ROI | Dorsal-Attention | R |

| Cingulo-Opercular-23_R-Ctx | R_PFop_ROI | Cingulo-Opercular | R |
| Cingulo-Opercular-24_R-Ctx | R_PF_ROI | Cingulo-Opercular | R |
| Frontoparietal-24_R-Ctx | R_PFm_ROI | Frontoparietal | R |
| Default-31_R-Ctx | R_PGi_ROI | Default | R |
| Default-32_R-Ctx | R_PGs_ROI | Default | R |
| Visual2-23_R-Ctx | R_FST_ROI | Visual2 | R |
| Visual2-26_R-Ctx | R_VMV2_ROI | Visual2 | R |
| Default-34_R-Ctx | R_31pd_ROI | Default | R |
| Cingulo-Opercular-25_R-Ctx | R_PoI1_ROI | Cingulo-Opercular | R |
| Somatomotor-20_R-Ctx | R_Ig_ROI | Somatomotor | R |
| Cingulo-Opercular-26_R-Ctx | R_FOP5_ROI | Cingulo-Opercular | R |
| Frontoparietal-27_R-Ctx | R_p47r_ROI | Frontoparietal | R |
| Auditory-07_R-Ctx | R_A4_ROI | Auditory | R |
| Frontoparietal-28_R-Ctx | R_TE1m_ROI | Frontoparietal | R |
| Cingulo-Opercular-28_R-Ctx | R_a32pr_ROI | Cingulo-Opercular | R |
| Cingulo-Opercular-29_R-Ctx | R_p24_ROI | Cingulo-Opercular | R |