

A Data-Driven Latent Variable Approach to Validating the Research Domain Criteria (RDoC) Framework

Quah, S.K.L.¹, Jo, B.¹, Geniesse C.¹, Uddin, L.Q.², Mumford, J.A.³, Barch, D.M.⁴, Fair, D.A.⁵, Gotlib, I.H.³, Poldrack, R.A.^{3,6}, Sagar, M.^{1,6}

¹Department of Psychiatry & Behavioral Sciences, Stanford University, Stanford, CA, USA

²Department of Psychiatry and Biobehavioral Sciences, University of California Los Angeles, CA USA

³Department of Psychology, Stanford University, Stanford, CA, USA

⁴Departments of Psychological & Brain Sciences, Psychiatry, and Radiology, Washington University in St. Louis, St Louis, MO, USA

⁵Department of Pediatrics, University of Minnesota Medical School, Minneapolis, MN, USA

⁶Wu Tsai Neurosciences Institute, Stanford University, Stanford, CA, USA

Abstract

Despite the widespread use of the Research Domain Criteria (RDoC) framework in psychiatry and neuroscience, recent studies suggest that the RDoC is insufficiently specific, or excessively broad, relative to the underlying brain circuitry it seeks to elucidate, leading to potential misrepresentation of circuit-function relations. We used a latent variable approach to address this issue, specifically utilizing bifactor analysis. We examined a total of 84 whole-brain task-based fMRI (tfMRI) activation maps from 19 studies with a total of 6,192 participants. Within this set of 84 maps, a curated subset of 37 maps with a balanced representation of RDoC domains constituted the training set of our analysis, and the remaining held-out maps formed the internal validation set. Furthermore, we externally validated the factor solutions from our curated training dataset using an independent set of 36 coordinate maps sourced through Neurosynth. We used RDoC constructs as seed terms for Neurosynth's topic meta-analysis. We hypothesized that if boundaries of RDoC domains warrant refinement, this would be indicated by the presence of overlapping domains or domains lacking specificity. Our findings suggest that a bifactor data-driven structure fits better with the current corpus of tfMRI data, with a general domain representing task-general patterns of brain activation. The data-driven model also proposes a different group of major domains, particularly splitting the RDoC cognitive systems domain into distinct domains. Data-driven models are useful for revising the posited circuit-function relations outlined in the current RDoC framework.

1. Introduction

The study of human neurobiology is a rapidly advancing field with significant implications for understanding brain function and, eventually, facilitating the development of valid biological markers and effective treatments for psychiatric disorders. Psychiatric disorders listed in the Diagnostic and Statistical Manual (DSM) have been considered to be discrete and unitary; recent research, however, suggests that they are both highly comorbid and heterogeneous across clinical samples^{1,2}. This heterogeneity may underlie the lack of well-established biomarkers to date for psychiatric disorders.

The Research Domain Criteria (RDoC) framework was developed by the National Institute of Mental Health (NIMH) to guide the development of a psychiatric nosology based on primary psychological functions and their associated biological features^{3,4}. The framework organizes core dimensions of behavior using a dimensional approach, viewing these aspects as varying along a continuum rather than in distinct categories. This approach spans multiple levels of analysis, from genes to behavior⁵. Within the RDoC framework, the fundamental neurobiological systems were defined and organized hierarchically into domains, with domain-specific constructs and sub-constructs. Now, over a decade since its inception, the framework's dimensional approach to psychopathology and its integration of multiple levels of analysis have contributed to a more nuanced and comprehensive understanding of brain function and mental disorders^{4,6}.

While the RDoC framework has helped guide research, a recent study using text-mining and machine learning found that a bottom-up data-driven ontological framework generated brain circuit-function links that were more reproducible than the RDoC or DSM frameworks⁷. They also showed that multiple RDoC domains shared underlying neural circuits, or some domains needed to be split. For example, Beam et al.⁷ showed that the RDoC domains of negative valence, positive valence, and arousal and regulation shared high mutual information across the

fronto-medial cortex and amygdala, indicating an overlap in the division of these domains.

Further, they also showed that the RDoC negative valence domain encompassed constructs that, from a data-driven framework, recombine elements of memory, reward, and cognitive systems. These findings prompt further investigation into potential refinements to RDoC's domain structure and mapping of brain function to neural circuits.

Researchers have made significant strides in attempting to develop a data-driven ontology that maps brain function to neural circuits through the meta-analysis of task-based fMRI (tfMRI) activation maps and topic modeling. Using data mining techniques, peak brain coordinate activation patterns during tasks have been categorized based on latent functional domains derived from study texts^{8,9} or task descriptions^{10,11}. While previous studies utilizing coordinate activation data have effectively harnessed the vast amounts of data available in databases like Neurosynth¹² and Brainmap¹³, they provide a very sparse representation of whole-brain activation. Image-based meta-analyses can provide a richer understanding of the intricate patterns of activation that occur during tasks¹⁴. It would be beneficial to compare RDoC directly with a data-driven model derived using image-based analyses to assess potential refinements to its framework.

To expand on the RDoC framework's hierarchical structure and address any potential overlap between domains or lack of specificity within a domain, we leveraged a latent variable approach with bifactor analysis to explore circuit-function relations. Bifactor models allow one to capture both shared variance across a number of latent constructs as well as variance unique to specific constructs. Assessing both general patterns of brain activity common across tasks^{15,16} and task-specific activation, Bolt et al. previously demonstrated that a bifactor model represents the relations between psychological constructs and underlying neural processes better than conventional non-hierarchical frameworks¹⁷. Using a bifactor model can help to identify shared and unique variance among the different constructs and provide more nuanced insight into the

organization of circuit-function relations. This approach can also help identify constructs that may be better conceptualized as part of a larger domain rather than as separate constructs. In this context, we used a bifactor analysis to examine the hierarchical structure of the RDoC framework across domains to provide data-driven evidence of complementary domain structures.

Specifically, we applied a latent variable approach with bifactor analysis to whole-brain task activation images from Neurovault and U.K. Biobank (n=84 select activation maps from 19 studies with a total of N=6,192 participants; adapted from Bolt et al¹⁷) to examine the organization of circuit-function relations. To ensure the robustness of our findings, we first derived our model solutions via a curated subset of the original dataset. Subsequently, we tested the model solution by applying it to the held-out maps, assessing its ability to generalize to previously unseen data. Moreover, we validated further using maps reconstructed from activation coordinates sourced from Neurosynth to assess the model's applicability to diverse data types. This comprehensive approach allows us to evaluate how well our model solution captures and represents brain activation patterns across various datasets and serves as a crucial step in advancing our understanding of circuit-function relations.

We posit that unclear boundaries between and within RDoC domains could lead some domains to lack sufficient specificity; for example, one domain may show strong connections to multiple latent factors. Conversely, other domains may be overly specified; for instance, multiple domains might share robust associations on the same latent factor. This examination advances our understanding of circuit-function relations in the brain and informs approaches to refine the RDoC framework.

2. Methods

2.1 Gathering and preparing activation maps: Our dataset comprises both whole-brain activation maps and maps reconstructed from activation coordinates (Fig. 1). These two sets of maps capture the primary published forms of neuroimaging data. Whole-brain activation maps underwent a rigorous selection process due to variations in contrast methodologies and acquisition parameters. The following subsections describe how we gathered and processed these maps.

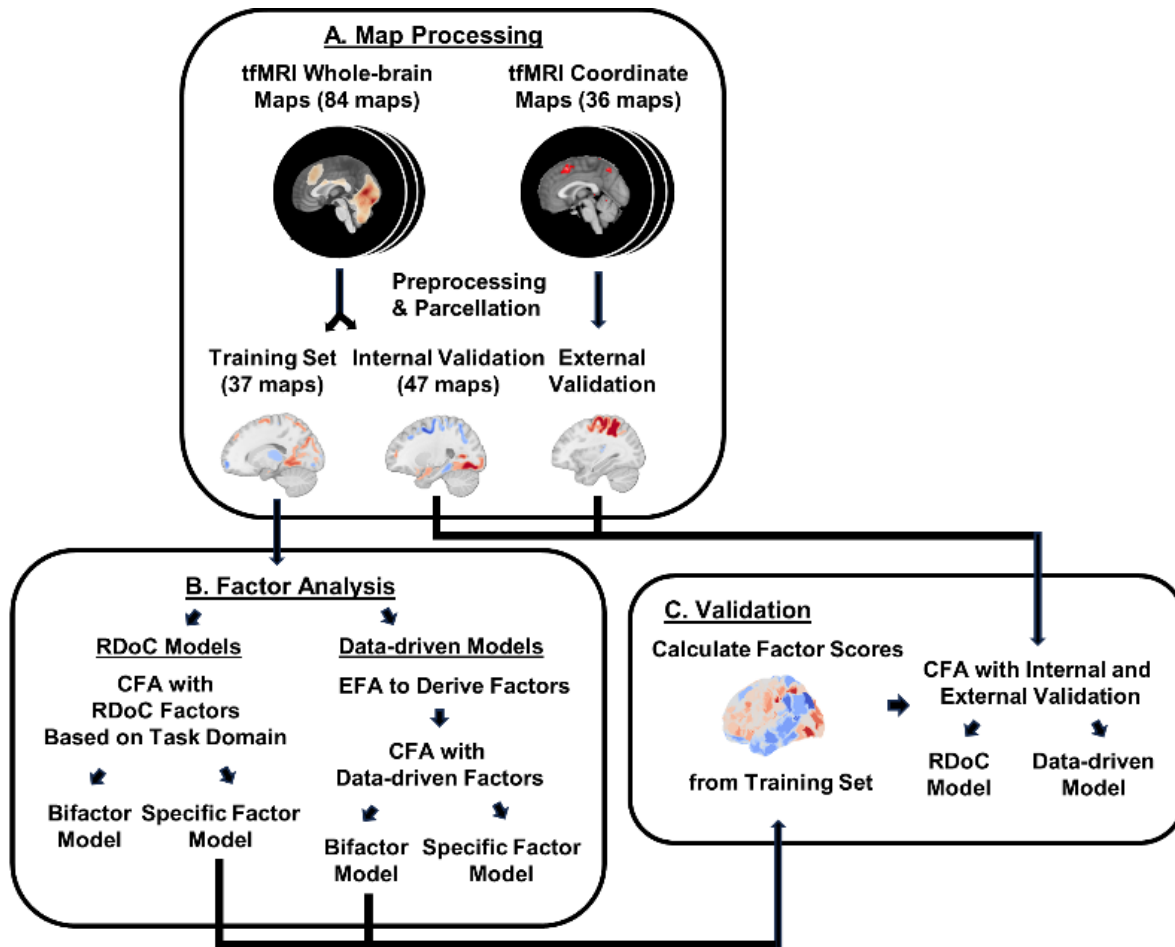


Figure 1: Approach to create and validate RDoC and data-driven factor models. (A) First, we divided tfMRI whole-brain maps into two subsets: the curated training dataset for building our factor models and the other as the internal validation set. Additionally, we processed tfMRI coordinate maps with peak activations to create an external validation set. (B) Regarding the confirmatory factor analysis (CFA) for the RDoC factor models, we assigned maps from the curated training dataset to specific factors corresponding to RDoC domains based on task

associations. Before conducting a CFA, we first performed an exploratory factor analysis (EFA) to determine factor assignments for data-driven factor models. (C) We employed a validation procedure to evaluate the model's performance on unseen data. We assigned maps to specific factors based on factor scores derived from the original data-driven and RDoC models. We then compared fit scores to assess the model's generalizability to new data.

2.1.1 Whole-brain activation maps. The collection of 84 whole-brain tfMRI maps was curated by Bolt et al.¹⁷ and sourced from two publicly accessible datasets: Neurovault¹⁸ and UK Biobank¹⁹. Although maps were also sourced from the Human Connectome Project by Bolt et al.¹⁷, these maps did not correlate sufficiently strongly with the other 84 maps within the dataset and, consequently, were not included for further analysis (Supplementary Figure 1). We used only unthresholded group-level BOLD contrasts for task conditions versus baseline. Contrast maps corresponding to the subtraction between two activation maps were not included because contrasts between events within the task would eliminate general activation patterns representing the task's domain.

We categorized contrast maps by matching the task descriptions extracted from the associated task contrasts (e.g., from <https://neurovault.org/> for NeuroVault) with descriptions of the RDoC domains and construct definitions from the RDoC matrix²⁰ (Supplementary Table 1). For instance, a contrast map created from a task where participants press a button as directed by visual instructions is categorized under the sensorimotor domain. We restricted our analysis to the following RDoC domains: cognitive systems, positive valence systems, negative valence systems, social processes, and sensorimotor systems, as no activation maps in the dataset fit within the domain of arousal and regulatory systems. Recognizing that a substantial proportion of the activation maps in the initial dataset originated from the cognitive systems domain (70%), we curated a sub-collection of maps. The curated training dataset was designed to achieve a more balanced representation of the constructs across all five domains and minimize study overlap. The curated training dataset is composed of 37 maps derived from a total of 6,119 participants, distributed as follows: cognitive systems (40.5%), negative valence systems

(13.5%), positive valence systems (18.9%), social processes (16.2%), and sensorimotor systems (10.8%). We also excluded maps representing tasks that strongly implicated multiple RDoC domains. Details of all 37 curated training maps and the 47 held-out maps (used for the internal validation set) are listed in (Supplementary Table 1).

Our initial collection of maps was composed of both *t*-stat and *z*-stat images. The unthresholded *t*-stat images were first converted to *z*-stat images before further processing. All maps were then resampled to the 2mm MNI-152 standard-space T1-weighted template (Nonlinear 6th generation).

2.1.3 Map Post-processing. All activation maps were parcellated into 333 cortical and 14 subcortical brain regions using the Gordon²¹ and Harvard-Oxford²² atlases, respectively.

2.2 Factor analysis. Latent variable models are designed to estimate latent constructs or classes that are not observed directly but are inferred from observed variables with measurement error²³. We conducted a comparative analysis of four distinct latent variable approaches, combining two methods of factor derivation (theory-driven RDoC factors or data-driven empirical factors) with two types of factor models (specific factor models or bifactor models). Specific factor models exclusively incorporate specific factors, while bifactor models have an additional general factor²⁴. To summarize, our study compared four models: (i) an RDoC specific factor model; (ii) an RDoC bifactor model; (iii) a data-driven specific factor model; and (iv) a data-driven bifactor model (Fig. 2). Data-driven models encompassed an exploratory factor analysis (EFA) step to first identify potential factor structures, followed by a confirmatory factor analysis (CFA) step to assess how well the factor model fits the observed data (Methods 2.2.1 and 2.2.3). In contrast, RDoC models involved only a CFA step, given that they incorporated pre-defined factors specific to RDoC (Methods 2.2.2 and 2.2.4).

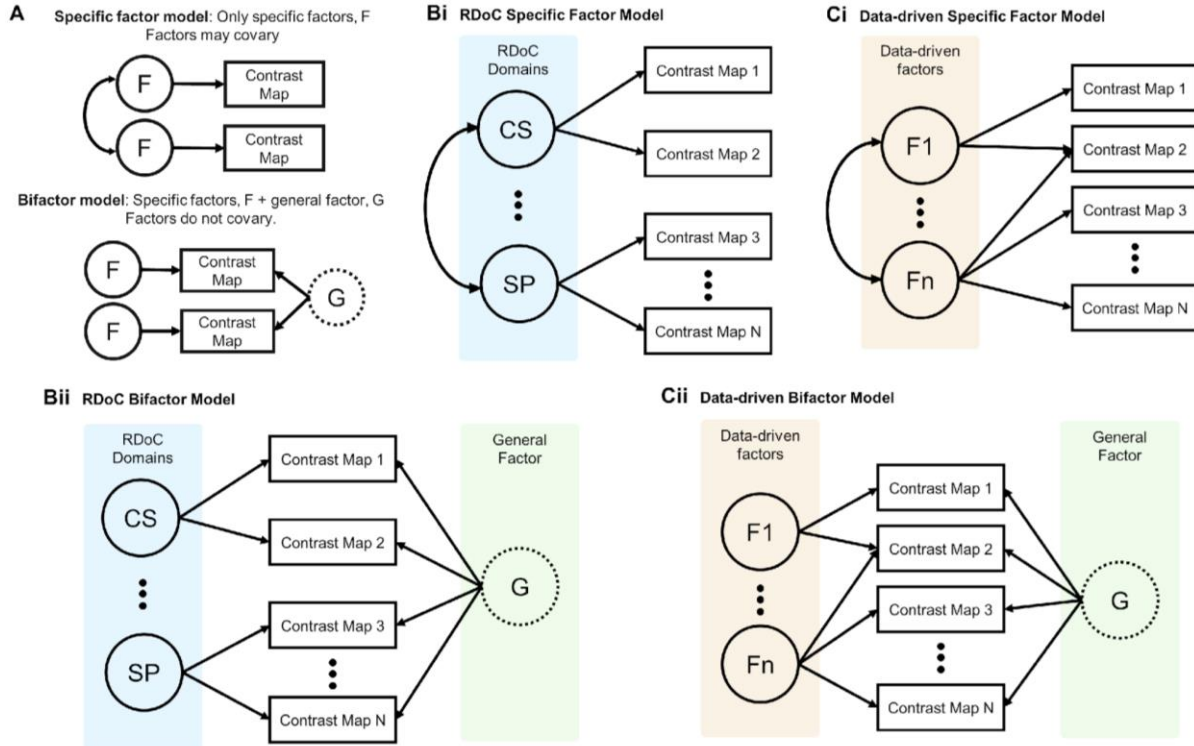


Figure 2: Factor model types. (A) Across all models, specific factors, F, denote brain activation patterns unique to a subset of tasks. In the bifactor models, the general factor, G, embodies brain activation patterns common across various tasks. (Bi-Bii) RDoC Models: These models are characterized by specific factors, each representing a distinct RDoC domain defined by task contrasts from whole-brain or coordinate maps. The Specific Factor Model (Bi): This model comprises specific factors. The Bifactor Model (Bii): This model is an extension of the specific factor model, with the addition of a general factor. (Ci-Cii) Data-Driven Models: These models (i & ii) are generated through EFA without predefined factors. CS and SP are representative RDoC domains, Cognitive Systems, and Social Processes, respectively.

Bootstrap distributions of fit indices were computed by resampling parcels over 5,000 iterations.

Factor scores were estimated using Bartlett's method to create brain maps (Fig. 5) reflecting each region's loading for each factor. This method is designed to yield factor scores that are strongly correlated with their respective factor, while maintaining minimal or no correlation with other factors.

2.2.1 Data-driven Factor Analysis with whole-brain maps. The factor analysis for our data-driven models (Fig. 2, Bi and Bii) was composed of three primary stages: (1) Horn's parallel analysis to determine the optimal number of factors to extract (see below); (2) EFA to extract specific factors; and (3) CFA with both the specific factors and a general factor for the bifactor model, and only specific factors for the specific factor model.

To determine the number of factors to extract, we conducted parallel analysis²⁵, which identifies the number of factors to extract based on where the calculated eigenvalues of the actual data intersect with the eigenvalues of random data generated²⁶. We then conducted an EFA using principal axis factoring and oblimin rotation to extract the identified number of specific factors in the subsequent confirmatory analysis. We also examined the scree plots to verify the suitability of the number of factors extracted (Supplementary Figure 3). To conduct the EFA, we used oblimin rotation to allow for correlated factors, but the correlation was constrained to be small. Based on previous work, each specific factor was defined by maps with a high absolute loading of 0.4 or higher²⁷. For the CFA, we used robust maximum likelihood estimation to account for non-normality in the data. Orthogonal rotation was used in the bifactor models to ensure that the general factor is not contaminated by the specific factors, making it difficult to interpret the factor structure. By constraining the general factor to be orthogonal to the specific factors, bifactor models can identify a general factor independent of the specific factors. The general factor captures the shared variance, while the specific factors capture the distinct variances that are unique to subsets of activation maps²⁸. We used the specific factors from the EFA and a general factor with all maps loaded onto it. For comparison, we also conducted an alternate CFA without the general factor (specific factor model). To account for interrelationships between factors within the specific factor models, which are not captured by a general factor, we maintained non-orthogonality and allowed all of our specific factor models to exhibit covariance.

2.2.2 RDoC Domain Factor Analysis with whole-brain maps. Our curated training set of whole-brain activation maps was grouped into RDoC domain-specific factors by matching the task description with the domain/construct definition. For our RDoC models (Fig. 2, Ai and Aii), we conducted a CFA utilizing robust maximum likelihood estimation and non-orthogonal factors. For comparison, we also conducted an alternate CFA using a bifactor model.

2.3 Statistical Analysis. Model fit was assessed using robust variants of fit indices, including the Root Mean Square Error of Approximation (RMSEA), Comparative Fit Index (CFI), and Tucker-Lewis Index (TLI). These fit indices were chosen to account for potential non-normality in the data. Additionally, information theoretical measures of model complexity, the Akaike Information Criterion (AIC) and Bayesian Information Criterion (BIC), were used for comparison. AIC and BIC consider the trade-off between model fit and complexity, with lower values indicating a more optimal balance²⁹. 95% confidence intervals of robust RMSEA were computed from the model using the observed robust RMSEA value, the degrees of freedom, and the sample size of the model. Bootstrap distributions of robust CFI, robust TLI, AIC, and BIC fit indices were computed by resampling over 5,000 iterations.

Pearson correlations were calculated between the factor scores of the RDoC specific factor model and the data-driven bifactor model. This analysis aimed to explore the extent to which the loadings of the base RDoC model align with the factors derived from the data-driven approach.

2.4 Validation using unseen data. We used a multi-prong validation strategy to assess the validity of the model solution derived from the curated training dataset. We compared the factor solutions from the RDoC specific factor model, representing the current RDoC framework, and the data-driven bifactor model, representing the best-performing data-driven model. For internal validation, we used the held-out maps from the original dataset, ensuring the model's reliability within the same type of dataset. Additionally, we used Neurosynth coordinate activation maps that were a different data type (compared to whole-brain) and had better coverage of the RDoC

domains (than the held-out maps) for external validation. This comprehensive validation strategy enabled us to evaluate the performance and generalizability of the factor structure we derived in varied contexts.

2.4.1 Internal Validation using held-out whole-brain maps. We systematically assigned individual maps to specific factors from the RDoC specific factor model (representing the RDoC framework) and the data-driven bifactor model. Factor assignment and loadings were determined using the factor scores derived from the original model using the curated training dataset. The factor assignment involved identifying, for each map, the factor from the original factor model that exhibited the highest product sum. After the map was assigned to a factor, the loading for each map was determined by dividing the map's product sum for that factor by the highest product sum of all other maps that were assigned to the same factor, providing an adjusted coefficient for its association with the respective factor (Supplementary Figure 2). Subsequently, we conducted a CFA with these factor assignments and loadings. We then compared the fit scores obtained from this validation analysis. This process allowed us to evaluate how well the training model solution generalized to unseen data, effectively probing the model's capability to capture brain activation patterns beyond the curated training dataset.

2.4.2 External validation using Neurosynth coordinate activation maps. In addition to using the held-out maps from our initial dataset to test the model solution derived using the curated dataset, we also utilized coordinate activation maps with topics matching RDoC construct seed terms for external validation. Seed terms adapted from Beam et. al.⁷ were compiled based on the name and synonyms of each RDoC domain construct, e.g., “acute threat” and “fear” for the “acute threat” construct under the negative valence system domain. These seed terms were then used to search for matching terms in a topic-based meta-analysis using Neurosynth. 200 topics were extracted using Latent Dirichlet Allocation (LDA) from the abstracts of all articles in the latest version of Neurosynth⁹ (ver. 5). Neurosynth's LDA topic-based meta-analysis is a

data-driven approach that uses natural language processing (NLP) techniques to uncover topics that share terms across a large set of studies. Each topic is associated with a probabilistic reverse inference map representing the likelihood that a given brain coordinate is activated during a study using these terms. Using this meta-analysis technique, we identified 36 coordinate activation maps with topics that matched RDoC construct seed terms. Seed terms with multiple topic maps had their activation averaged before further analysis. Spatial smoothing was applied using a 12-mm full-width half-maximum (FWHM) Gaussian kernel centered on each peak-activation coordinate in the maps, creating more realistic representations of brain activation patterns. Values were then thresholded ($z > 0.1$) to remove noise using a Gaussian kernel. A complete list of the seed terms and topics sourced from Neurosynth is presented in Supplementary Table 2. These maps were then used to validate the factor structure from the curated training dataset in the same way as the held-out maps (2.4.1). Bootstrap values for fit indices in the external validation models were determined over 10,000 iterations to address the high likelihood of convergence issues and to enhance the robustness of the computed fit indices.

3. Results

3.1 RDoC models with whole-brain maps

We conducted two CFAs with RDoC factors: one with only specific factors and another with an additional general factor (bifactor model). Based on the task description of each contrast map (Supplementary Table 1), maps were grouped into specific factors by matching respective RDoC domains' definitions.

In the specific factor model, most maps within each domain loaded significantly (i.e., $|\text{loading score}| \geq 0.4$) onto each factor representing their domains (cognitive systems: 11/15; negative

valence systems: 5/5; positive valence systems: 6/7; social processes: 6/6; sensorimotor systems: 4/4; Fig. 3A).

Comparing the RDoC specific factor model with the bifactor model to examine whether adding a general factor would improve the fit, we found that the bifactor model had a better fit according to all fit indices. This suggests that adding a general factor reflecting domain-general activation patterns improved the model fit of the conventional RDoC framework. This was also true after accounting for model complexity in the additional number of parameters estimated in the bifactor model, indicating that the addition of a general factor also provided a better balance between fit and complexity.

3.2 Data-driven models with whole-brain maps

In the data-driven approach, we also conducted two CFAs: one with only specific factors and another with an additional general factor (bifactor model). The specific factors for both models are latent variables derived using EFA that account for the unique variance among subsets of activation maps. They represent dimensions of task activation patterns that are not shared across all maps. Parallel analysis was first conducted to determine the appropriate number of factors to extract from the dataset. The parallel analysis indicated that models with eight factors or less had eigenvalues greater than what was expected by chance (Supplementary Figure 3). Thus, we extracted eight specific factors in the data-driven CFAs.

In the data-driven bifactor CFAs, all maps loaded significantly (i.e., $|\text{loading score}| \geq 0.4$) on the general, specific, or both factors. All but two maps across RDoC domains loaded on the general factor, indicating that maps across distinct studies and tasks showed overlap in activation patterns (Fig. 3B). Notably, the two maps that did not load on the general factor were associated

with contrasts related to button pressing in response to an auditory cue; in contrast, the tasks in the dataset primarily revolved around responses to visual cues.

Furthermore, maps labeled by RDoC domains showed divergent patterns in loadings across specific factors (Fig. 3B). Positive valence systems, social processes, and sensorimotor systems domain maps showed high loadings that were confined to relatively few specific factors. In contrast, cognitive and negative valence systems domain maps showed significant loadings spread across multiple specific factors.

The data-driven bifactor model also had a greater overall fit to the data compared with both RDoC models and the data-driven specific factor model (Fig. 3C). However, after accounting for the different number of parameters estimated in the models, the data-driven bifactor model had a better model fit than the RDoC specific factor model but not the data-driven specific factor model (Figure 3C). All model fit scores are shown in Table 1.

After deriving these models, we created a product matrix to study similarities in map loadings across factors from the RDoC specific factor model and the data-driven bifactor model (Fig. 4A). The values in the product matrix represent the average product of absolute non-zero value factor loadings in both models. The values range from 0-1, where 1 represents a complete 1-to-1 similarity in map loadings, and 0 represents no overlap. This matrix provides insight into the consistency of the boundaries within and without the RDoC domains. Maps of domains with cross-loading on many specific factors reflect heterogeneity within the domain's boundaries (low intra-domain consistency); maps of domains that share high loading with other domains on the same specific factor reflect overlap in the domains' boundaries (high inter-domain similarity).

The cognitive systems and negative valence systems domains load across multiple specific factors, indicating low intra-domain consistency. This suggests a degree of heterogeneity within the boundaries of these domains. In contrast, the sensorimotor systems domain shows notable

intra-domain consistency by loading heavily on only a single data-driven factor (Fig. 4A), indicating a relatively consistent pattern in the activation maps of this domain. The positive valence systems and social processes domains demonstrate loadings across various data-driven factors, with particularly high loadings for data-driven factors 8 and 1, respectively. This implies that the boundaries of these domains may benefit from some refinement, given the observed complexities in their activation patterns across different factors. RDoC specific factors that share high loadings with data-driven factors (Fig. 4A) also show high factor score correlations (Fig. 4B-C)

Brain maps of factor scores and map loading for the data-driven bifactor and RDoC specific factor model are shown in Fig. 5. All of the RDoC domains but the sensorimotor systems domain show positive factor scores across both visual and motor regions, implicating the frequent recruitment of these regions across tasks of different domains. The sensorimotor systems domain predictably showed notable positive factor scores across the motor cortex. Similarly, the factors score brain map of the data-driven bifactor model's general factor captured the predominant recruitment of visual and motor regions across most tasks. In contrast, the factor scores of the data-driven model's specific factors captured more specific and varied functional activation patterns.

3.3 Validation with held-out whole-brain maps and Neurosynth coordinate maps

To evaluate the robustness and generalizability of our model solutions, we conducted a validation procedure using both the held-out maps from the original dataset (internal validation) and the coordinate maps sourced from Neurosynth (external validation) (Fig. 6).

For internal validation using held-out whole-brain maps, we first compared the model fit of factors derived from the RDoC-specific factor model (conforming to the current RDoC

framework) with those from the data-driven bifactor model. Our analysis revealed that the data-driven model exhibited a better fit for the held-out maps. These findings highlight the data-driven model's superior fit and generalizability, even when tested on previously untrained data, compared with the RDoC model. All model fit scores are shown in Table 2.

External validation with Neurosynth coordinate maps was conducted to evaluate the model's generalizability to diverse data types. We did not include a general factor in our data-driven model. Here, coordinate maps are sparse and do not show substantial overlaps that a general factor would represent. Indeed, the general factor of a data-driven bifactor model from a CFA exhibits limited loading across all the coordinate maps, indicating a lack of substantial influence (Supplementary Figure 4). Similar to our findings with the held-out whole-brain activation maps, the data-driven specific factor model demonstrated a better fit for the Neurosynth coordinate maps compared to the RDoC model. These results indicate that the data-driven models demonstrated better fit and generalizability to both unseen whole-brain maps and coordinate maps compared to the RDoC models (Table 2).

Table 1: Model fit comparison. This table presents the model fit statistics for the factor models. The models are categorized into RDoC and data-driven, each further divided into specific factor and bifactor models. The data-driven models generally showed better model fit than all the RDoC models. The data-driven bifactor model showed the best model fit as measured by robust RMSEA, robust CFI, and robust TLI, but the data-driven specific factor model showed the best model fit as measured by AIC and BIC. Lower RMSEA, AIC, and BIC values and higher CFI and TLI values signify superior model fit. Each metric is accompanied by 95% CI. Bold values represent the best fit for each measure within each row (*significantly better-fit scores).

Map Type	Models	RDoC		Data-driven	
		Specific Factor	Bifactor	Specific Factor	Bifactor
Whole-brain Maps	Robust RMSEA	0.218	0.202	0.214	0.190*
	95% CI	0.214-0.222	0.198-0.207	0.209-0.220	0.185-0.194
	Robust CFI	0.495	0.582	0.590	0.633*
	95% CI	0.463-0.507	0.547-0.596	0.560-0.599	0.598-0.643
	Robust TLI	0.456	0.532	0.530	0.587*
	95% CI	0.423-0.469	0.492-0.547	0.497-0.541	0.548-0.599
	AIC	26,530	24,781	21,662*	23,770

	95% CI	25,725- 27,035	23,991- 25,238	20,952-22,041	22,948-24,226
	BIC	26,853	25,200	22,028*	24,193
	95% CI	26,048-27,359	24,410- 25,658	21,318-22,406	23,371-24,650

Table 2: Internal validation model fit using held-out whole-brain maps. The data-driven bifactor model showed better model fit across all model fit indices (robust RMSEA, robust CFI, robust TLI, AIC and BIC) compared to the RDoC specific factor model representing the RDoC framework. Lower RMSEA, AIC, and BIC values and higher CFI and TLI values signify superior model fit. Each metric is accompanied by 95% CIs. Bold values represent the best fit for each measure within each row (*significantly better fit scores).

Map Types	Models	RDoC Specific Factor	Data-driven Bifactor
Held-out whole-brain Maps	Robust RMSEA	0.205	0.198*
	95% CI	0.202- 0.209	0.194-0.201
	Robust CFI	0.426	0.484*
	95% CI	0.388-0.443	0.449-0.498
	Robust TLI	0.421	0.462*
	95% CI	0.383-0.438	0.426-0.477
	AIC	34,697	33,147*
	95% CI	33,725-35,434	32,296-33,705
	BIC	34,912	33,493*
	95% CI	33,941-35,649	32,643-34,051

Table 3: External validation model fit using Neurosynth coordinate maps. The data-driven specific factor model showed better model fit across all model fit indices (robust RMSEA, robust CFI, robust TLI, AIC and BIC) compared with the RDoC specific factor representing the RDoC framework. Lower RMSEA, AIC, and BIC values and higher CFI and TLI values signify superior model fit. Each metric is accompanied by 95% CIs. Bold values represent the best fit for each measure within each row (*significantly better fit scores). A significant proportion of bootstrapped models (approximately 30%) generated non-admissible solutions, which notably skewed the robust CFI and TLI confidence intervals (Supplementary Figure 5). Consequently, confidence intervals for these measures have not been included here.

Map Types	Models	RDoC Specific Factor	Data-driven Specific Factor
Coordinate Maps	Robust RMSEA	0.281	0.256*
	95% CI	0.277- 0.285	0.251-0.260
	Robust CFI	0.130	0.303
	95% CI	-	-
	Robust TLI	0.124	0.275
	95% CI	-	-
	AIC	33,060	29,897*
	95% CI	29,555-34,728	23,602-31,707
	BIC	33,214	30,128*
	95% CI	29,706-34,878	23,829-31,934

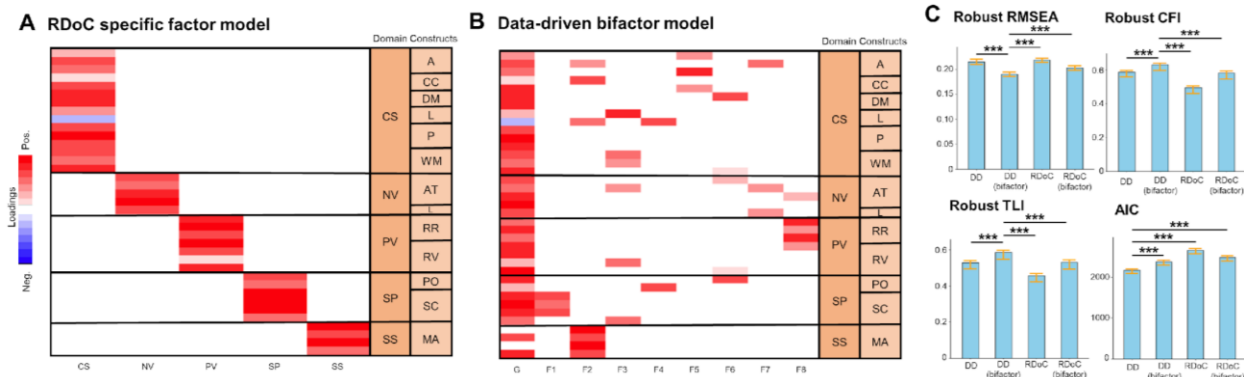


Figure 3: Comparison of RDoC and Data-driven models using whole-brain activation maps. (A-B) Heatmaps showing factor loadings of the RDoC and data-driven models across RDoC domain classified maps. Warmer colors = positive; Cooler colors = negative factor loadings shown. (C) Relative fit measures of different latent variable models from whole-brain activation maps. The data-driven bifactor model (DD (bifactor)) was the model with the best fit based on robust RMSEA, CFI and TLI, but the data-driven specific factor model (DD) was the model with the best fit based on AIC. 95% CI error bars. CS: Cognitive Systems; NV: Negative Valence systems; PV: Positive Valence systems; SS: Sensorimotor Systems; SP: Social Processes; G: General factor; F1-8: specific factors 1-8; and DD: Data-driven. Domain-Constructs: Cognitive Systems-A: Attention; CC: Cognitive Control; DM: Declarative Memory; L: Language; P: Perception; WM: Working Memory; Negative Valence systems-AT: Acute Threat; L: Loss; Positive Valence Systems-RR: Reward Response; RV: Reward Valuation; Social Processes-PO: Perception of Others; SC: Social Communication; Sensorimotor Systems-MA: Motor Action.

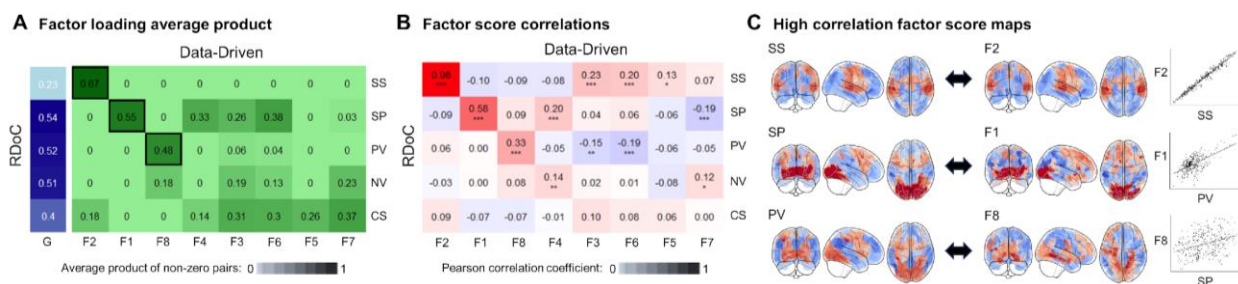


Figure 4: Factor convergence. (A) Heatmap showing the average product of factor loadings for maps in both the RDoC specific factor model and data-driven bifactor model. High values²⁷ above |.4| are highlighted with black borders. Maps of domains cross-loading on many specific factors reflect low intra-domain consistency; maps of domains sharing high loading with other domains on the same specific factor reflect high inter-domain similarity. (B) The heatmap shows the Pearson correlation of factor scores for the data-driven bifactor and RDoC specific factor model. (*p-value < .05; **p-value ≤ .01; ***p-value ≤ .001). Rows and columns are organized to

show the strongest correlations in the diagonal. (C) Glass brains of factors with the highest correlations are shown as illustrative examples of strong one-to-one convergence in factor scores on the brain (warmer colors = positive scores; cooler colors = negative scores). Specifically, the sensorimotor systems domain displayed a strong correspondence with data-driven factor 2, the positive valence systems domain aligned with data-driven factor 1, and the social processes domain strongly correlated with data-driven factor 8. Scatter plots of correlations are shown on the right. CS: Cognitive Systems; NV: Negative Valence systems; PV: Positive Valence systems; SS: Sensorimotor Systems; SP: Social Processes; F1-6: specific factors 1-8.

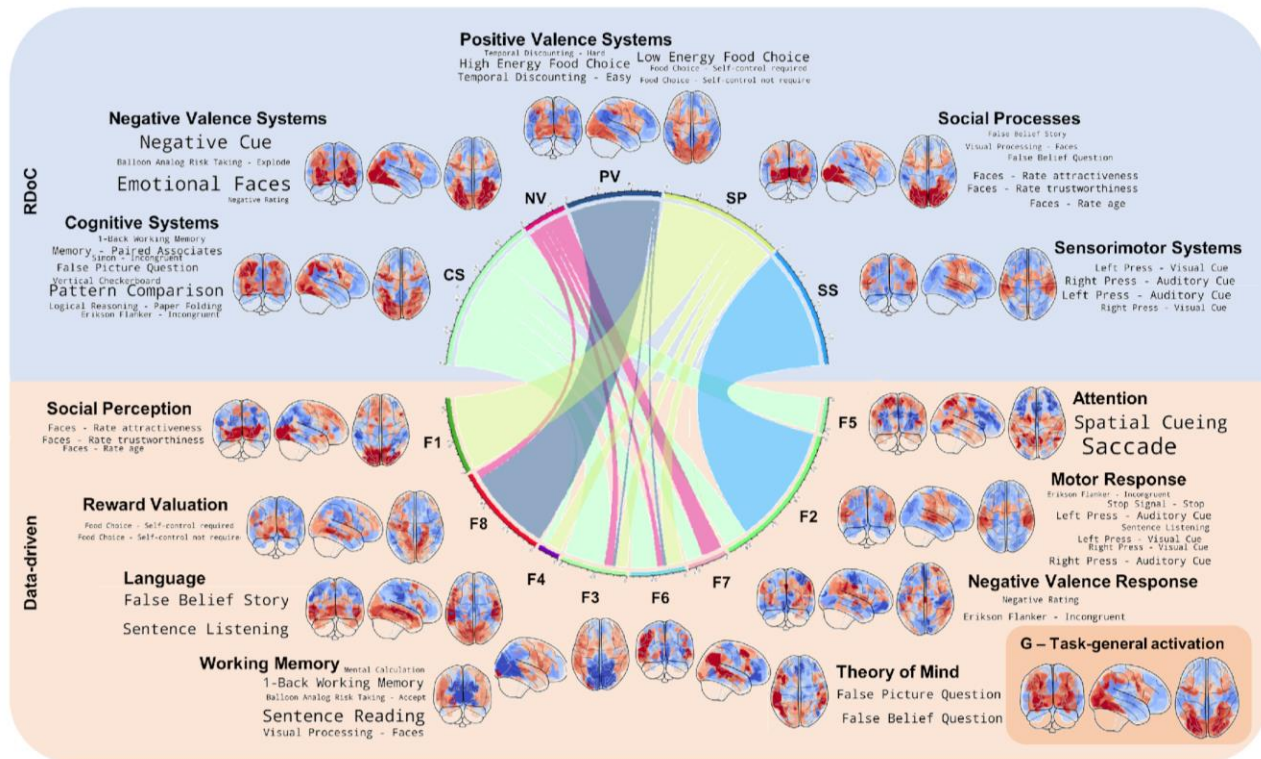
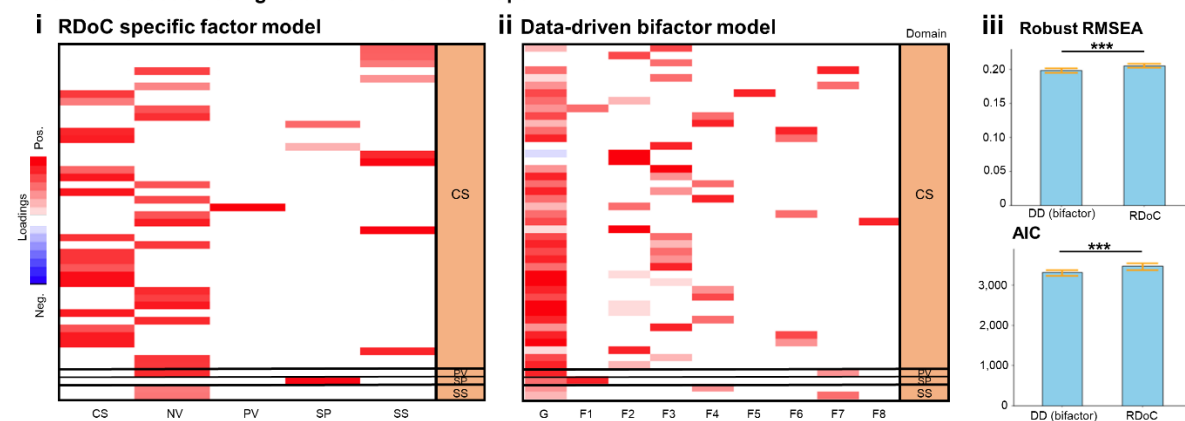


Figure 5: Mapping factors of the RDoC specific factor model and the data-driven bifactor model using data from whole-brain activation maps. Chord diagram showing the weighted links between maps loading on both models. The product of loadings in both models weighs links. Overlaps in maps showing loadings on both models illustrate the complex relations between RDoC factors and those derived using a data-driven bifactors modeling approach. Glass brain maps reflect factor scores (warmer colors = positive scores; cooler colors = negative scores). Word clouds of factors reflect the paradigm descriptors of the top eight maps loading on each factor. The size of words reflects the magnitude of the factor's loading.

A Internal validation using held-out whole-brain maps



B External validation using coordinate maps

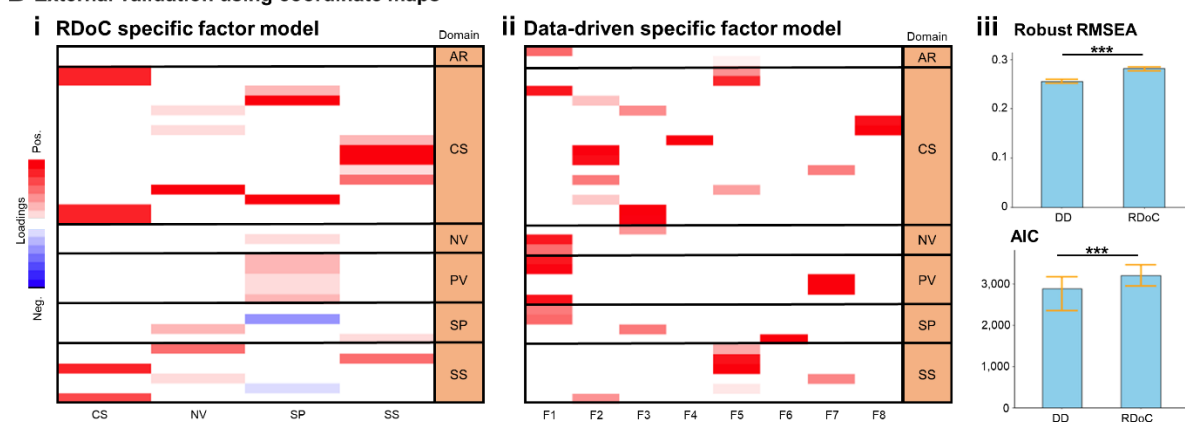


Figure 6: Validating RDoC and Data-driven models using held-out whole brain maps and

coordinate maps. Heatmaps showing factor loadings of the (i) RDoC and (ii) data-driven

models for (A) whole-brain and (B) coordinate activation maps. Whole-brain maps held-out from the curated training dataset used to create the original models formed the internal validation set.

The external validation set consists of coordinate activation maps obtained using Neurosynth's LDA topic-based meta-analysis. These coordinate maps represent RDoC construct seed terms.

The data-driven validation model with coordinate maps did not incorporate a general factor. This omission stemmed from the nature of sparse coordinate maps, which lacked significant overlaps that would warrant the representation of a general factor, as detailed in Supplementary Figure 4.

Warmer colors = positive; Cooler colors = negative factor loadings shown. (iii) Relative fit

measures of different latent variable models from whole-brain activation maps. These new maps were assigned to specific factors based on factor scores from data-driven and RDoC models.

The data-driven model outperformed the RDoC model in terms of fit scores when applied to whole-brain maps and coordinate maps, capturing brain activation patterns for different unseen datasets. 95% CI error bars. AR: Arousal and Regulatory Systems; CS: Cognitive Systems; NV: Negative Valence systems; PV: Positive Valence systems; SS: Sensorimotor Systems; SP: Social Processes; G: General factor; F1-8: specific factors 1-8; and DD: Data-driven.

4. Discussion

The current study aimed to advance the ontology of human brain functions by using a latent variable approach with bifactor analysis to examine the hierarchical structure of the RDoC framework. Our findings suggest that a data-driven approach provides a more accurate representation of the organization of the human brain's circuit-function relations than does the RDoC model.

The traditional RDoC model had most maps within each domain that loaded significantly onto each factor representing their domains; however, compared with data-driven models, the RDoC model also showed a relatively poor fit for both whole-brain and coordinate activation maps, indicating that the RDoC framework may not fully capture the complexity of brain-behavior relations. Adding a general factor to the conventional RDoC also improved the fit of the RDoC specific factor model, suggesting that the conventional RDoC framework may benefit by adding a superordinate domain representing task-general functioning. Incorporating a task-general functional domain into the RDoC model that extends beyond the existing task-specific functional domains would enhance the model's ability to represent brain functioning comprehensively. Impairments within this domain may reflect transdiagnostic alterations that cause changes to domain-general/task-nonspecific processing, including attention or awareness¹⁷.

Compared to the RDoC model, the data-driven model had a better fit to the data, indicating that it may provide a more accurate representation of the organization of circuit-function relations in the human brain. By differentiating general activation patterns common across different functional tasks from patterns specific to each construct, the data-driven bifactor model captured both shared and unique variance among different constructs, providing insight into the hierarchical organization of circuit-function relations. This is consistent with findings from recent studies that have advocated for a data-driven bifactor approach to understanding brain-behavior relations¹⁷. Notably, the data-driven specific factor model had better fitness scores after penalizing for model complexity as measured by AIC and BIC. This indicates that although the data-driven bifactor model had the best overall model fit, the improvement in fit from adding the general factor comes at a substantial cost in model complexity.

The product matrix (Fig. 4A) and factor score correlations (Fig. 4B) revealed divergent patterns in correspondence across data-driven factors for different RDoC domains. For instance, whereas the cognitive systems domain had low loadings and correlations spread across the data-driven factors, the maps labeled by the positive valence systems, social processes, and sensorimotor systems domains had significant loadings and correlations that were confined to relatively fewer specific factors. Finally, the negative valence systems domain did not have significant loadings on any data-driven factors (Fig. 3). Still, its factor scores correlated strongly with two data-driven factors (Fig. 4). This pattern suggests that activation patterns within some domains are more distinct and separable than others, supporting our hypothesis that the boundaries between RDoC domains need to be reconsidered. Specifically, constructs within the cognitive systems domain might be better defined by being divided into separate domains. For example, attention, working memory, semantic processing/perception, and theory of mind within the cognitive systems domain formed individual data-driven factors (Fig. 5), and may be better represented as a revised set of domains in a refined RDoC framework.

Visualization of the factor scores on the brain showed us that the RDoC factors, excluding the sensorimotor system's domain, consistently reveal activation patterns spanning visual and motor regions. This alignment with the general factor of the data-driven bifactor model suggests that there is shared task-general activation across tfMRI whole-brain maps. The utility of the general factor in the data-driven model lies in its ability to capture overarching patterns present across the entire dataset. This, in turn, allows the specific factors to focus on representing activation patterns that exhibit greater sensitivity to the nuances of specific task paradigms.

After constructing our factor models, we performed validation steps to assess how much our derived model, developed from the curated training dataset, could extend to unseen data. We used two distinct validation sets: whole-brain maps held out from the original dataset (internal validation) and coordinate maps sourced externally from Neurosynth (external validation). The internal validation using held-out whole-brain maps, while sharing the same data type as the original dataset, had a skewed distribution of maps (more cognitive maps) across the RDoC domains. To address this imbalance, we also conducted validation using Neurosynth coordinate maps, which provide a more balanced representation of the RDoC domains and constructs. This dual validation approach enhances the reliability of our findings and strengthens our model's applicability to diverse datasets and contexts. Our data-driven model consistently exhibited superior fit and generalizability in both cases compared to the RDoC specific factor model representing the conventional RDoC framework. These outcomes underscored the data-driven model's capability to capture brain activation patterns, extending beyond the initial dataset. Moreover, external validation using coordinate maps highlighted the model's adaptability to diverse data types, particularly in handling sparse coordinate activation maps commonly generated from extensive meta-analytic tools.

4.1 Limitations: Despite these advancements, it must be acknowledged that the overall model fit, even with the data-driven approaches, was not optimal. This limitation underscores the need

for continued refinement and development in this field, recognizing that the complexity of brain-behavior relations may pose challenges to modeling efforts. Further research verifying the validity of newly defined functional domains with different datasets and cohorts is also needed to consolidate the ontological advancements made in this study. It is also important to note that although task activation relative to baseline allowed us to capture general task activation in our models, tasks here often involve more than one functional domain. For example, even a simple button press-to-cue task involves perception (cognitive domain) and motor action (sensorimotor systems domain). Therefore, subtraction contrasts between tasks may reveal additional insights into the brain's function-circuit relations. Moreover, while the bifactor model offers a valuable framework for understanding the complexity of brain-behavior relations, future work is needed to explore other model structures.

4.2 Conclusion: In conclusion, our study indicates that a data-driven approach provides a more accurate representation of the organization of the human brain's circuit-function relations than the conventional RDoC model. Our findings support the use of data-driven approaches to inform revisions to the RDoC framework and to develop a more comprehensive ontology to guide further research. Integrating a task-general domain within the RDoC framework holds promise in broadening the capacity of the RDoC framework to capture brain functionality holistically. Furthermore, our research underscores the need to reassess the demarcations or boundaries within RDoC domains, especially in delineating constructs within the various domains. Future studies should continue to explore the utility of these approaches to refine the RDoC framework and unravel the intricate dynamics of circuit-function relations in the human brain.

References

1. Feczko, E. *et al.* The Heterogeneity Problem: Approaches to Identify Psychiatric Subtypes. *Trends Cogn. Sci.* **23**, 584–601 (2019).
2. Stephan, K. E. *et al.* Charting the landscape of priority problems in psychiatry, part 1: classification and diagnosis. *Lancet Psychiatry* **3**, 77–83 (2016).
3. Insel, T. *et al.* Research domain criteria (RDoC): toward a new classification framework for research on mental disorders. *Am. J. Psychiatry* **167**, 748–751 (2010).
4. Cuthbert, B. N. & Insel, T. R. Toward the future of psychiatric diagnosis: the seven pillars of RDoC. *BMC Med.* **11**, 126 (2013).
5. Cuthbert, B. N. & Insel, T. R. Toward new approaches to psychotic disorders: the NIMH Research Domain Criteria project. *Schizophr. Bull.* **36**, 1061–1062 (2010).
6. Cuthbert, B. N. Research Domain Criteria (RDoC): Progress and Potential. *Curr. Dir. Psychol. Sci.* **31**, 107–114 (2022).
7. Beam, E., Potts, C., Poldrack, R. A. & Etkin, A. A data-driven framework for mapping domains of human neurobiology. *Nat. Neurosci.* **24**, 1733–1744 (2021).
8. Rubin, T. N. *et al.* Decoding brain activity using a large-scale probabilistic functional-anatomical atlas of human cognition. *PLoS Comput. Biol.* **13**, e1005649 (2017).
9. Poldrack, R. A. *et al.* Discovering relations between mind, brain, and mental disorders using topic mapping. *PLoS Comput. Biol.* **8**, e1002707 (2012).
10. Yeo, B. T. T. *et al.* Functional Specialization and Flexibility in Human Association Cortex. *Cereb. Cortex* **25**, 3654–3672 (2015).
11. Bolt, T. *et al.* Ontological Dimensions of Cognitive-Neural Mappings. *Neuroinformatics* **18**, 451–463 (2020).

12. Yarkoni, T., Poldrack, R. A., Nichols, T. E., Van Essen, D. C. & Wager, T. D. Large-scale automated synthesis of human functional neuroimaging data. *Nat. Methods* **8**, 665–670 (2011).
13. Fox, P. T. & Lancaster, J. L. Mapping context and content: the BrainMap model. *Nat. Rev. Neurosci.* **3**, 319–321 (2002).
14. Salimi-Khorshidi, G., Smith, S. M., Keltner, J. R., Wager, T. D. & Nichols, T. E. Meta-analysis of neuroimaging data: a comparison of image-based and coordinate-based pooling of studies. *Neuroimage* **45**, 810–823 (2009).
15. Hugdahl, K., Raichle, M. E., Mitra, A. & Specht, K. On the existence of a generalized non-specific task-dependent network. *Front. Hum. Neurosci.* **9**, 430 (2015).
16. Fox, M. D. *et al.* The human brain is intrinsically organized into dynamic, anticorrelated functional networks. *Proceedings of the National Academy of Sciences* **102**, 9673–9678 (2005).
17. Bolt, T., Nomi, J. S., Yeo, B. T. T. & Uddin, L. Q. Data-Driven Extraction of a Nested Model of Human Brain Function. *J. Neurosci.* **37**, 7263 (2017).
18. Gorgolewski, K. J. *et al.* Neurovault.org: a web-based repository for collecting and sharing unthresholded statistical maps of the human brain. *Front. Neuroinform.* **9**, 8 (2015).
19. Miller, K. L. *et al.* Multimodal population brain imaging in the UK Biobank prospective epidemiological study. *Nat. Neurosci.* **19**, 1523–1536 (2016).
20. RDoC Matrix. *National Institute of Mental Health (NIMH)*
<https://www.nimh.nih.gov/research/research-funded-by-nimh/rdoc/constructs/rdoc-matrix>.
21. Gordon, E. M. *et al.* Generation and Evaluation of a Cortical Area Parcellation from Resting-State Correlations. *Cereb. Cortex* **26**, 288–303 (2016).
22. Desikan, R. S. *et al.* An automated labeling system for subdividing the human cerebral cortex on MRI scans into gyral based regions of interest. *Neuroimage* **31**, 968–980 (2006).

23. Bollen, K. A. *Structural Equations with Latent Variables*. (Wiley & Sons, Limited, John, 2017).
24. Fang, G., Guo, J., Xu, X., Ying, Z. & Zhang, S. IDENTIFIABILITY OF BIFACTOR MODELS. *Stat. Sin.* **31**, 2309–2330 (2021).
25. Horn, J. L. A RATIONALE AND TEST FOR THE NUMBER OF FACTORS IN FACTOR ANALYSIS. *Psychometrika* **30**, 179–185 (1965).
26. Hayton, J. C., Allen, D. G. & Scarpello, V. Factor Retention Decisions in Exploratory Factor Analysis: A Tutorial on Parallel Analysis. *Organizational Research Methods* **7**, 191–205 (2004).
27. Stevens, J. (James P. *Applied multivariate statistics for the social sciences*. 629 (L. Erlbaum Associates, 1992).
28. Carroll, J. B. Human cognitive abilities: A survey of factor-analytic studies. **819**, (1993).
29. Burnham, K. P. & Anderson, D. R. Multimodel Inference: Understanding AIC and BIC in Model Selection. *Sociol. Methods Res.* **33**, 261–304 (2004).

Acknowledgments

This work was supported by an NIH R01MH127608 and an MCHRI Faculty Scholar Award to M.S.

Data & code availability

The whole-brain task fMRI contrast maps used in this study are publicly available at the neurovault.org website. The coordinate maps used are available at neurosynth.org. The R code used for latent variable analysis and visualization will be available upon publication at this address: <https://github.com/braindynamicslab/>