

1 **Title:** Representational learning of brain responses in executive
2 function and higher-order cognition using deep graph convolutions

3 **Running Title:** Decoding high-order cognition of human brain

4 **Authors:** Yu Zhang^{1,2,3,*}, Nicolas Farrugia⁴, Alain Dagher⁵ and Pierre Bellec^{2,3,*}

5

6 ¹ Research Center for Healthcare Data Science, Zhejiang Lab, Hangzhou, China

7 ² Centre de recherche de l'Institut universitaire de gériatrie de Montréal, Montréal, QC H3W 1W6,
8 Canada

9 ³ Department of Psychology, Université de Montréal, Montréal, QC H3C 3J7, Canada

10 ⁴ Department of Mathematical and Electrical Engineering, IMT Atlantique, Brest, France

11 ⁵ McConnell Brain Imaging Centre, Montreal Neurological Institute, McGill University, Montreal,
12 QC, Canada

13 *** Corresponding Author:**

14 Yu Zhang

15 Research Center for Healthcare Data Science, Zhejiang Lab

16 Zhongtai Street, Yuhang District, Hangzhou, Zhejiang, China

17 yuzhang2bic@gmail.com

18

19 Pierre Bellec

20 Département de Psychologie, Université de Montréal

21 4565, Chemin Queen-Mary, Montréal (Québec) H3W 1W5

22 pierre.bellec@gmail.com

23 **Abstract**

24 Brain decoding aims to infer human cognition from recordings of neural activity using modern
25 neuroimaging techniques. Studies so far often concentrated on a limited number of cognitive states
26 and aimed to classifying patterns of brain activity within a local area. This procedure demonstrated a
27 great success on classifying motor and sensory processes but showed limited power over higher
28 cognitive functions. In this work, we investigate a high-order graph convolution model, named
29 ChebNet, to model the segregation and integration organizational principles in neural dynamics, and
30 to decode brain activity across a large number of cognitive domains. By leveraging our prior
31 knowledge on brain organization using a graph-based model, ChebNet graph convolution learns a
32 new representation from task-evoked neural activity, which demonstrates a highly predictive signature
33 of cognitive states and task performance. Our results reveal that between-network integration
34 significantly boosts the decoding of high-order cognition such as visual working memory tasks, while
35 the segregation of localized brain activity is sufficient to classify motor and sensory processes. Using
36 twin and family data from the Human Connectome Project ($n = 1,070$), we provide evidence that
37 individual variability in the graph representations of working-memory tasks are under genetic control
38 and strongly associated with participants in-scanner behaviors. These findings uncover the essential
39 role of functional integration in brain decoding, especially when decoding high-order cognition other
40 than sensory and motor functions.

41

42 **Keywords:** fMRI, brain decoding, functional integration, graph convolutional network, brain atlas,
43 brain connectome, saliency map, representation similarity

44

45 **Teaser**

- 46 • Modelling functional integration through graph convolution is a necessary step towards
47 decoding high-order human cognition.

48

49 **Significance statement**

50 Over the past two decades, many studies have applied multivariate pattern analysis to decode what
51 task a human participant is performing, based on a scan of her brain. The vast majority of these
52 studies have however concentrated on select regions and a specific domain, because of the
53 computational complexity of handling full brain data in a multivariate model. With the fast progress
54 in the field of deep learning, it is now possible to decode a variety of cognitive domains
55 simultaneously using a full-brain model. By leveraging our prior knowledge on brain organization
56 using a graph-based model, we uncovered different organizational principles in brain decoding for
57 motor execution and high-order cognition by modelling functional integration through graph
58 convolution.

59 Introduction

60 Understanding the neural substrates of human cognition is a main goal of neuroscience research.
61 Modern imaging techniques, such as functional magnetic resonance imaging (fMRI), provide an
62 opportunity to map cognitive function in-vivo, and to decode the dynamics of cognitive processes
63 from neural activity. Brain decoding has been an active topic since Haxby and colleagues first
64 proposed the idea of using fMRI brain responses to predict the category of visual stimuli presented to
65 a subject (1). Nowadays, a variety of computational models are used in the field, including multi-
66 voxel pattern recognition, linear regression models, as well as nonlinear models such as deep artificial
67 neural networks (DNN). Among which, DNN showed promising advantages over other linear models
68 by providing an end-to-end solution to a direct mapping from recorded brain activity to brain
69 cognition, for instance, using convolutional (2) and recurrent neural networks (3). However, most
70 previous decoding studies aimed to segregate the spatial patterns of brain activation under different
71 task conditions, but largely ignored the integration of brain dynamics during cognitive processes.
72 Functional segregation into highly localized brain areas, and functional integration at the levels of
73 distributed brain regions, modules and networks, are fundamental principles of brain organization and
74 have been widely observed in different populations and among a variety of cognitive tasks (4–7). So
75 far, the majority of brain decoding studies only utilized the functional specialization hypothesis that
76 aims to distinguish the localized brain activation patterns under a small number of experiment tasks,
77 for instance, the involvement of the motor and sensory cortex during the movement of different body
78 parts (8), or the engagement of different regions in the visual cortex for the recognition of various
79 types of visual stimuli (9). As a result, such brain decoders were restricted to mostly motor and
80 sensory processes (e.g. recognition of visual stimuli) and highly relied on domain knowledge (e.g.
81 activating different parts of the visual cortex). This assumption of functional segregation also limited
82 the generalizability of brain decoding towards high-order cognitive functions that were known to
83 engage multiple brain systems. One typical example is the visual working memory task (VWM), for
84 which multiple brain networks were involved through intense interactions among memory
85 representations and other basic attention and sensory processes (10). For instance, early visual cortex

86 played an important role in the detection of visual features including orientation (11), motion (12) and
87 content (13), while the parietal and prefrontal cortex contributed to maintenance of visual information
88 over a delayed interval (14). In these cases, both local and global information of brain activity may
89 contribute to the decoding of cognitive processes (15,16).

90 We started to tackle this problem in our previous paper (17) by generalizing the convolutional
91 operations from DNN onto brain organization. This approach can effectively capture both segregated
92 brain activity of task-related brain regions, and information integration of neural dynamics within
93 brain networks. Compared to previous linear and nonlinear decoding models, the proposed decoding
94 model showed high generalizability over a variety of cognitive domains without relying on any prior
95 information on the tested domain. However, this model relied on a simplified version of graph
96 convolution which only took into account information integration within the same brain network at
97 each layer. It showed limited power of representational learning on high-order cognitions that may
98 involve complex forms of functional interactions across multiple brain systems.

99 To address this issue, we investigated a more sophisticated form of graph convolution in this study,
100 namely ChebNet, which approximates the calculation of graph convolution using high-order
101 Chebyshev polynomials. It has been proved that the ChebNet graph convolution is K -localized in
102 space (on the graph) by taking up to K th order Chebychev polynomials (18). In other words, ChebNet
103 integrates information within a relatively larger neighborhood by taking multiple steps of random
104 walks on the brain graph. As a result, ChebNet graph convolution is capable of characterizing the
105 complex forms of information processing during cognitive processes, i.e. segregating task-evoked
106 activity from localized brain regions ($K=0$), integrating neural activity within the same brain networks
107 ($K=1$), as well as information integration between different networks and among multiple brain
108 systems ($K>1$). ChebNet provides a generalized form to encode this multiscale hierarchical
109 organization of brain cognition in a single graph convolutional layer. The decoding model started with
110 a parcellation that divides the whole brain into hundreds of brain regions and a brain graph that
111 captures hierarchical and modular structures in brain organization. The brain graph as well as the
112 dynamic information flow (i.e. task-evoked brain response at each brain region) on the graph was then
113 mapped onto a new representational space through multilayer spatiotemporal graph convolutions.

114 These embedded graph representations naturally disassociate different cognitive tasks with large
115 distances between task conditions and small distances within the same condition, and can improve the
116 prediction of cognitive states by achieving better functional alignment between multiple trials and
117 across different subjects.

118 In order to verify this hypothesis, we constructed the decoding model based on ChebNet graph
119 convolution at different orders, ranging from local brain regions ($K=0$), to the same brain network
120 ($K=1$), to multiple brain systems ($K>1$). All decoding models were evaluated on the task-fMRI
121 database from the Human Connectome Project (HCP)(19) and simultaneously distinguished 6
122 cognitive domains or 21 task conditions by using a short time window of fMRI scans (e.g. 10
123 seconds). Under this framework, we systematically investigated how large-scale functional integration
124 impact on brain decoding especially for multidomain decoding and decoding of high-order cognitive
125 functions. Taking Motor and Working-memory tasks as examples, we further investigated the
126 organizational structures among ChebNet layers within and across decoding models, and explored
127 their relations to the two principles of brain organization, i.e. functional segregation and integration.
128 Moreover, we investigated whether the representations learned through ChebNet graph convolutions
129 were able to improve inter-subject alignment in brain responses and preserve individual variability in
130 brain organization at the same time.

131

132 **Results**

133 Decoding cognitive functions with fine cognitive granularity and high accuracy

134 We proposed a decoding pipeline based on ChebNet graph convolution which automatically learns the
135 spatiotemporal dynamics of brain activity from a short series of fMRI responses and predicts brain
136 states based on learned feature representations (as shown in Figure 1-S1). The model starts with a
137 brain graph with nodes representing brain parcels and edges representing brain connectivity, maps
138 task-evoked fMRI responses onto the predefined brain graph, and learns high-level graph
139 representations of neural activity by using stacked graph convolutions, taking into account both
140 segregated neural activity within localized brain regions and functional interactions among between
141 brain networks. For a more detailed description of the decoding model, please see the “Methods”
142 section and Supplemental Information.

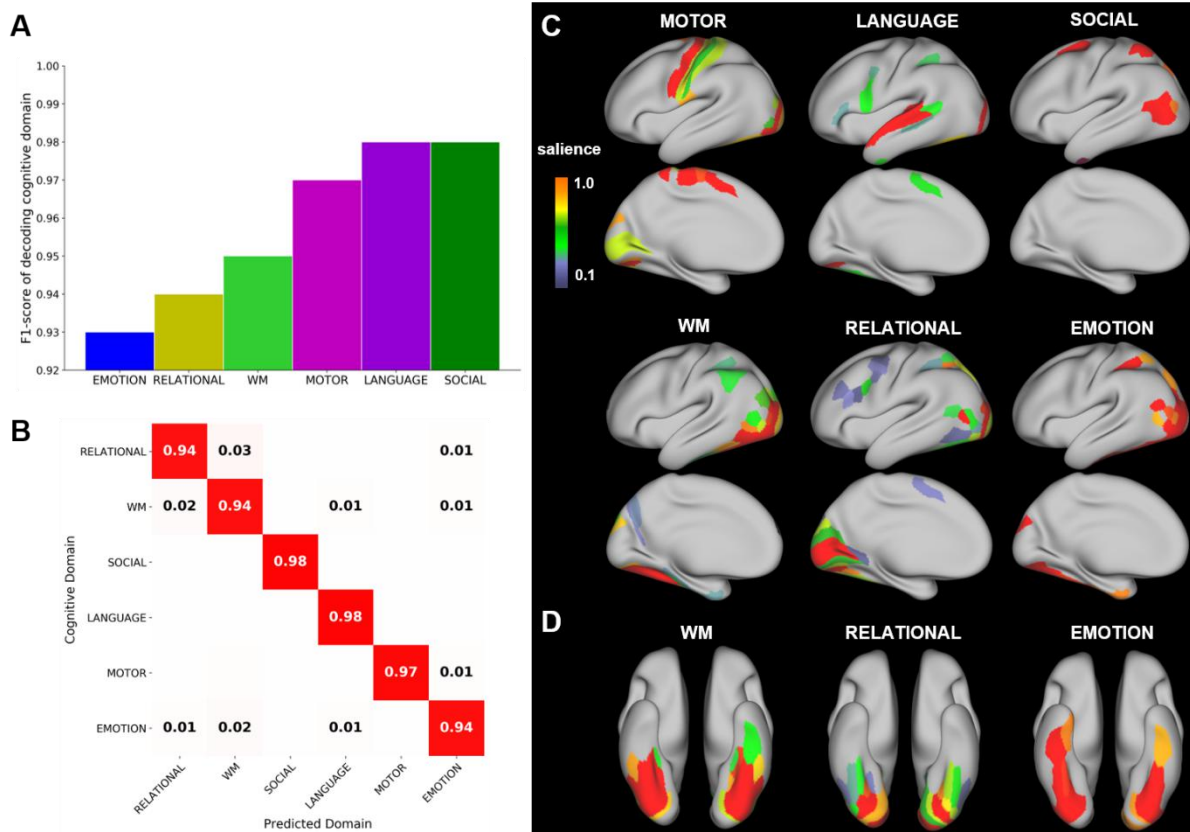
143 The ChebNet decoding model was evaluated using the cognitive battery of HCP task-fMRI dataset
144 acquired from 1200 healthy subjects. Using the ChebNet-*K*5 model (i.e. ChebNet graph convolution
145 with $K=5$), the six cognitive domains were nearly perfectly differentiated from each other by only
146 using 10s of brain recordings (approximately the shortest duration of task conditions in HCP), with an
147 average test accuracy of 96% (mean=95.81%, STD =0.15% by using 10 fold cross-validation with
148 shuffle splits). Moreover, the pipeline was capable of distinguishing experimental conditions with fine
149 cognitive granularity and fine temporal resolution, either across multiple domains (see Table S3) or
150 within each cognitive domain (see Table S2), and achieved high decoding accuracy on both tasks.
151 Among the six cognitive domains (as shown in Figure 1-S2 and Table S2), the language tasks (2
152 conditions, story vs math), and motor tasks (5 conditions, left/right hand, left/right foot and tongue)
153 were the most easily recognizable conditions, and showed the highest precision and recall scores (F1-
154 score = 98.45% and 99.38%, respectively for classifying two language conditions and five motor
155 conditions). The model achieved high decoding performance on other high-order cognitive functions
156 when longer duration of task blocks was available, for instance working-memory (94.51%, classifying
157 8 conditions using 25s) and social cognition (96.58%, classifying 2 conditions using 23s). Our

158 decoding model outperformed existing linear and nonlinear models including other deep learning
159 architectures, which neglect the hierarchical brain organization during cognitive processes, for either
160 classifying between cognitive domains (e.g. 93.7% when using 27 TRs reported in (2)) or decoding
161 task states within specific domain (e.g. 92.6% and 92.9% for working-memory and social cognition
162 respectively when decoding on 30s of fMRI data (3)).

163 Brain decoding captured reliable and task-specific salient features

164 In order to validate that the decoding model used biological meaningful features, we generated the
165 saliency maps on the trained decoding model by propagating the non-negative gradients backwards to
166 the input layer (20). An input feature is *salient or important* only if its little variation causes big
167 changes in the decoding output. The saliency scores were evaluated for each task trial independently
168 and then averaged within each subject and for each condition (cognitive domain or task state). First,
169 different sets of salient brain regions were detected for each cognitive domain (as shown in Figure 1C
170 and D), for instance the involvement of the somatosensory cortex for motor execution (MOTOR) and
171 the engagement of perisylvian language areas for language comprehension (LANGUAGE). Second,
172 the salient features were not only highly selective to specific cognitive tasks but also very stable
173 across trials and subjects. We took the Motor and Working-memory tasks as examples. The reliability
174 of saliency values was evaluated by using repeated-measure ANOVA, controlling for the random
175 effect of subjects and experimental trials. Only the salient brain regions that having high saliency
176 values (>0.3) and showing a significant effect of task ($p < 0.001$) were reported in the following
177 analysis. As shown in Figure 2, salient brain regions in the sensorimotor cortex were identified in the
178 Motor tasks, including region “a” (labelled as “area 5m” in the Glasser’s atlas) selectively activated
179 during foot movements, region “b” (labelled as “area 2”) selectively activated during hand movements,
180 regions c (labelled as “area OP4”) selectively activated during tongue movements. This distinctive
181 pattern in the saliency maps were highly consistent across trials, sessions, and even subjects. For
182 Working memory tasks, which involved both the differentiation between 0back vs 2back tasks and the
183 recognition of different image categories, the decoding model learned reliable features related to both
184 aspects, i.e. memory-load and image category. Here, we plotted the salient features for 0back and

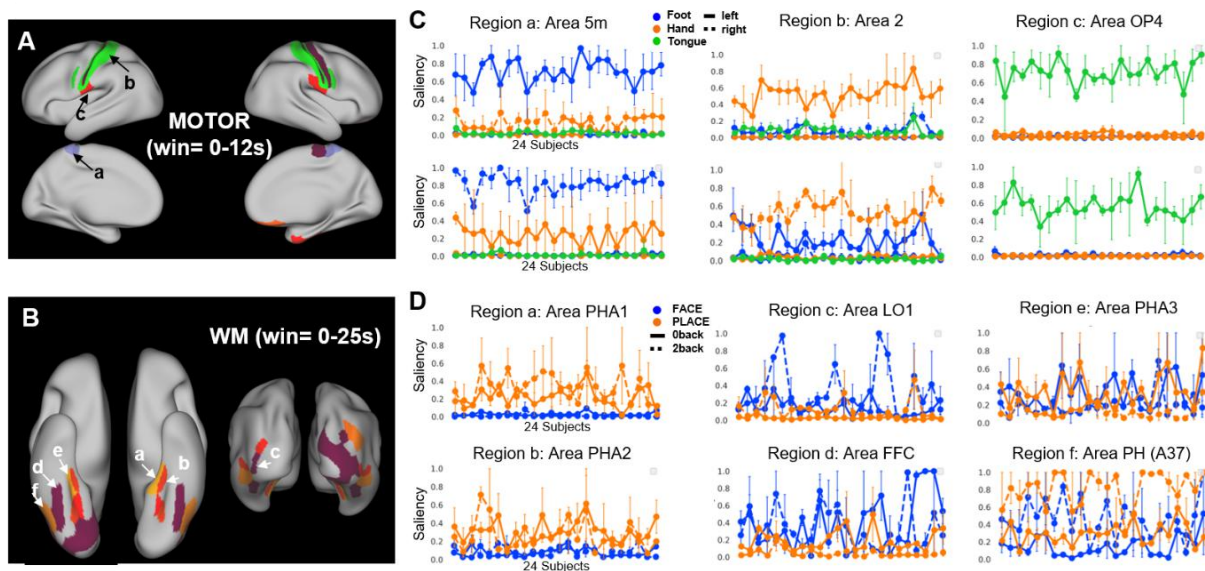
185 2back tasks on face and place images. As shown in Figure 2D, ParaHippocampal Area 1 (PHA1) and
 186 ParaHippocampal Area 2 (PHA2) were selectively involved for the recognition of place images
 187 (repeated measure ANOVA, F-score=70.96 and 38.12, p-value=1.74e-8 and 3.0e-6 respectively for
 188 PHA1 and PHA2), while Fusiform Face Complex (FFC) and Lateral Occipital Area 1 (LO1) were
 189 selectively engaged for the recognition of faces (F-score=57.75 and 91.47, p-value=1.02e-7 and
 190 1.75e-9 respectively for FFC and LO1). On other hand, for both place and face images,
 191 ParaHippocampal Area 3 (PHA3) was more involved in 0back tasks than 2back tasks (F-score=26.38,
 192 p-value=3.3e-5) while area PH was selectively engaged in the 2back tasks (F-score=102.56, p-
 193 value=6.01e-10) when fixing the stimuli category. Our results revealed that the decoding model
 194 captured reliable salient features from task-evoked brain activities, in order to distinguish among
 195 cognitive domains and task states. These salient features were derived from task-related brain regions
 196 and showed selective responses to different task conditions with high consistency not only with the
 197 same subject but also between different subjects (as shown in Figure 2), possibly revealing the
 198 biological basis of the decoding model.
 199



201
202
203
204
205
206
207
208
209
210
211
212

Figure 1. Decoding on six cognitive domains and the corresponding saliency maps.

The decoding model predicted the cognitive domain from each 10s of fMRI responses and achieved an average test accuracy of 96%. The F1-score on each domain was shown in A. The corresponding cross-domain confusion matrix was shown in B. The saliency maps were evaluated for each task trial independently and then averaged within each subject and each domain. Different sets of salient brain regions were detected for each cognitive domain (C). Due to the similarity in task stimuli, the salient features in the ventral visual stream were identified for image recognition in three cognitive tasks, i.e. Working-memory (WM), relational processing (RELATIONAL) and emotional processing (EMOTION). Still, the decoding model captured different sets of visual areas for the three cognitive domains (D).



213
214
215
216
217
218
219
220

Figure 2. Salient features for the Motor (A) and Working-memory (B) tasks.

Saliency value of each individual trial was estimated by using the guided backpropagation approach. The stability of saliency was evaluated by plotting the saliency values across randomly selected HCP subjects. The effect of task condition in the saliency values was then evaluated by using repeated-measure ANOVA, with the 'subject' as the random effect and 'task condition' as the within-subject effect. Only salient brain regions (saliency values > 0.3, the full range of saliency is (0,1)) with a significant 'task condition' effect ($p < 0.001$) were shown in the final saliency maps (A and B). For

221 Motor task (C), three salient brain regions were selected that showed selective responses to the
222 movement of foot (region “a”), hand (region “b”) and tongue (region “c”). The task trials
223 corresponding to the movements of the left body parts were plotted in solid lines and in dashed lines
224 for the right body parts. Brain regions in the left hemisphere were shown in the 1st row and the right
225 hemisphere shown in the 2nd row. For Working-memory task (D), three sets of salient brain regions
226 were selected that showed selective responses to the image category, e.g. place (1st column, in orange)
227 and face image (2nd column, in blue), or to memory load, e.g. 0back (solid line) and 2back (3rd column,
228 dashed line).

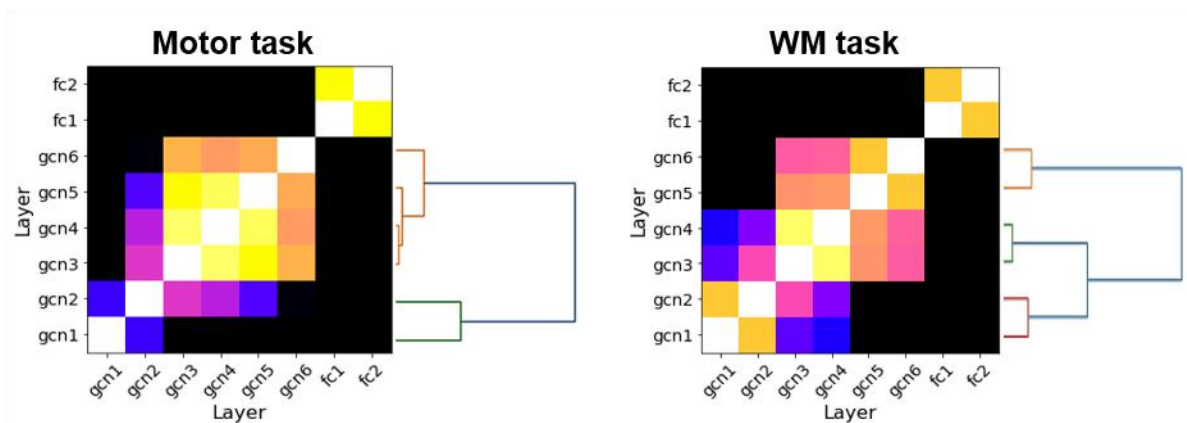
229

230 Decoding model learned hierarchical representations among ChebNet layers

231 The decoding model not only extracted biologically meaningful features associated with task-related
232 brain regions, as illustrated by the saliency map analysis, but also learned hierarchical representations
233 of brain response in each ChebNet layer. For instance, in the first graph convolutional layer (gcn1),
234 the model learned various shapes of temporal kernels (accounting for the hemodynamic response in
235 BOLD signals). Using these kernels, the model extracted a collection of spatial “*activation maps*”,
236 which resembled the actual brain activation maps detected by the canonical GLM approach (Figure 3-
237 S1 and Figure 3-S2). More sophisticated and task-specific feature representations were captured in
238 deeper layers. In order to verify the hierarchy among ChebNet layers, we evaluated the similarity of
239 feature representations using centered kernel alignment (CKA) with a linear kernel (21), with $0 <$
240 $CKA < 1$. As shown in Figure 3, a block-diagonal structure was detected in the CKA matrix of
241 Working-memory (WM) tasks, indicating a hierarchical organization of representations among
242 stacked ChebNet layers such that each layer inherited some information from previous layers, learned
243 new representations in the current layer and passed these features onto the next layer.

244 The hierarchical clustering was applied to the CKA matrix and revealed a strong disassociation
245 between the low-level features (gcn1 to gcn2), hidden representations (gcn3 to gcn4), and high-level
246 representations (gcn5 to gcn6). Weak associations were detected across different levels (CKA=0.94
247 and 0.76 for within- and between-level similarity), with a stepwise progression towards the last

248 ChebNet layer (CKA=0.54, 0.83, 0.92 for low, middle, high-level features as compared to gcn6),
249 where category-specific information was present (i.e. different representations between task
250 conditions). A similar hierarchical organizational structure was detected on the Motor task (as shown
251 in Figure 3) but with fewer levels in the hierarchy, i.e. low- and high-level features, and with high
252 redundancy in the middle layers (gcn3 to gcn5, average similarity with gcn6 is CKA=0.92). Still,
253 distinct features were learned in the low- and high-level representations (CKA=0.58 for gcn1-gcn2 as
254 compared to gcn6). Besides, the model already captured category-specific information starting in
255 early ChebNet layers (Figure 6-S2). Our results indicated that the ChebNet decoding model learned
256 hierarchical representations across graph convolutional layers in order to capture the underlying
257 neural dynamics during cognitive processes. The hierarchy in the representations of the decoding
258 model resembled the hierarchy in brain organization which has been reported in a variety of cognition
259 functions especially for high-order cognition (22,23). Moreover, the different organizational patterns
260 in ChebNet representations between cognitive tasks to some extent reflects different scales of
261 information integration in cognitive processes, such that a deep ChebNet architecture was required to
262 encode the complex forms of functional integration in WM, while a shallow ChebNet was sufficient
263 to encode the segregation of localized brain activity during motor execution.



265 **Figure 3. Hierarchical organization of layer representations learned through ChebNet graph**
266 **convolutions.**

267 The similarity of representations between ChebNet layers was first calculated using CKA with a
268 linear kernel. A distance metric was then generated from the CKA matrix ($dis = 1 - cka$). After that, the
269 hierarchical clustering was applied to the distance matrix using Ward's linkage. The resulting

270 dendrogram illustrated the hierarchical organization among ChebNet layers, for instance two-level
271 organization in the Motor task and a tripartite organization in the Working-memory task.

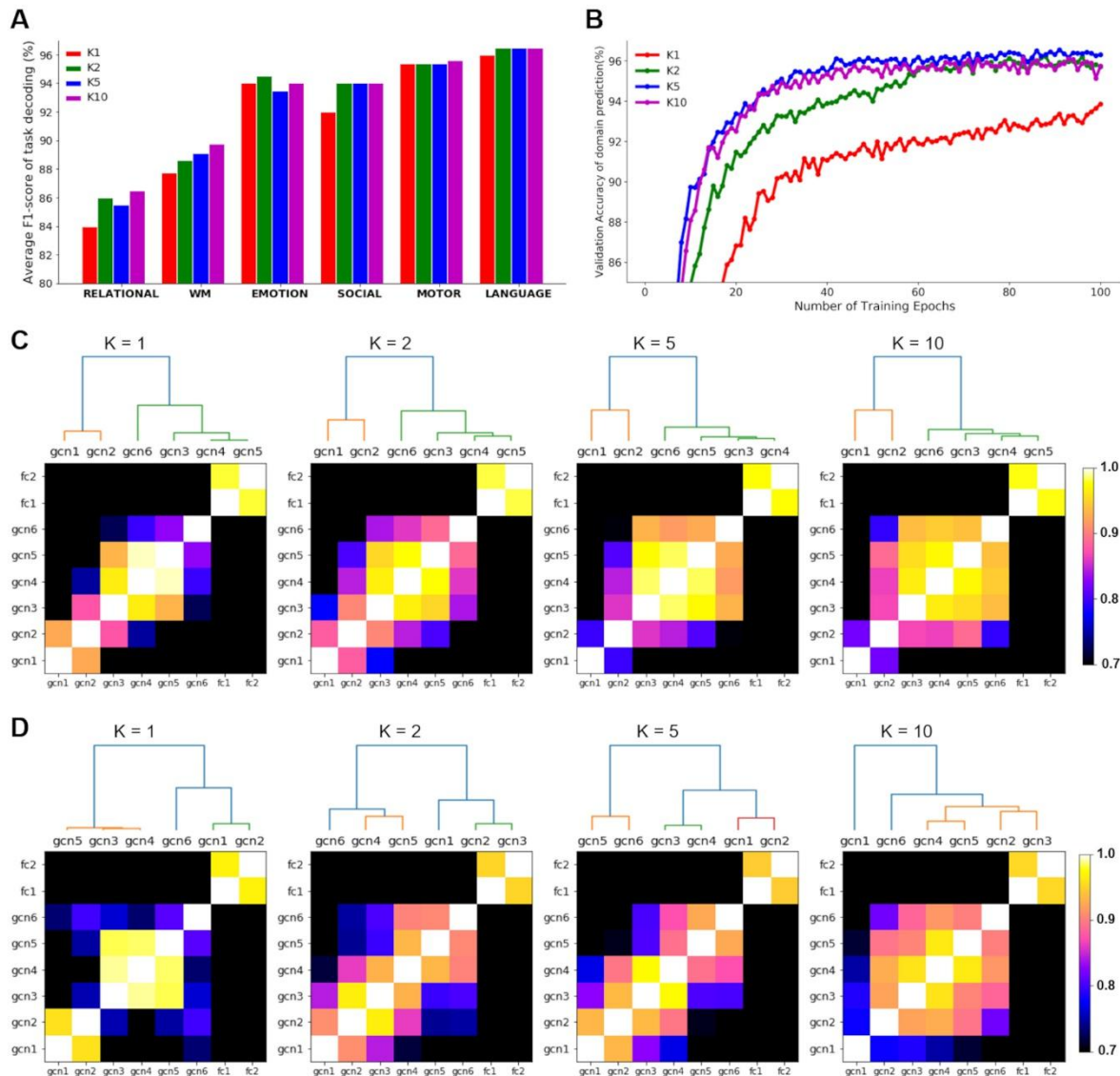
272

273 Variable sensitivity to the K -order uncovers different organizational principles in cognitive
274 processes

275 Another factor that impacts information integration in brain decoding is the K -order of ChebNet, by
276 taking into account multi-level integration of neural dynamics at each graph convolutional layer,
277 ranging from localized brain areas ($K=0$) to spatially distributed regions within the same network
278 ($K=1$) and towards inter-connected brain networks ($K>1$). The choice of K -order not only showed an
279 impact on the decoding performance, but also changed the hierarchy of feature representations learned
280 in each ChebNet layer.

281 First of all, the decoding of six cognitive domains significantly impacted by the choice of K -order in
282 ChebNet, indicating a faster convergence speed as well as higher decoding accuracy when using high-
283 order models (Figure 4B). Significant improvements in decoding were detected between $K=1$
284 (integration of brain activity within the same network) and $K>1$ (between-network communication)
285 (test accuracy = 93% vs 96% respectively for $K=1$ and $K>1$), significantly boosted compared to the
286 localized decoding model (test accuracy = 83% for $K=0$). Second, variable sensitivity to the K -order
287 was detected among different cognitive domains (Figure 4A). Specifically, for the Motor task, the
288 decoding performance showed no improvement when increasing K , which means no gain from
289 between-network communication during motor execution. Coinciding with this, the hierarchical
290 organization of layer representations in the Motor task showed a very stable bipartition pattern when
291 increasing K , i.e. low- and high-level features (as shown in Figure 4C). By contrast, the decoding of
292 WM tasks gradually improved as increasing K and reaching the plateau after $K > 5$, which means that
293 between-network communication and high-order integration plays an important role in WM,
294 especially for distinguishing between 0back and 2back tasks (as shown in Figure 4-S1). Interestingly,
295 the hierarchical organizational structure in WM (as shown in Figure 4D) started with three isolated

296 clusters at $K=1$, gradually fused the representations by filling the gaps between neighboring layers,
297 and converged to a stable tripartite organization at $K=5$ (i.e. low-, middle- and high-level
298 representations). Further increase in the K -order did not change this organization but instead
299 expanded the middle-level through encoding redundant hidden representations. Our results indicated
300 that the variable sensitivity to the choice of K -order may uncover distinct organizational principles in
301 cognitive processes, for instance, localized information processing within the motor and sensory
302 cortex for motor execution, while complex forms of functional interaction and information integration
303 across multiple brain systems/networks during WM tasks. Our findings coincided with the notion of
304 functional segregation and integration in brain cognition (24), for instance, within-network
305 communication is essential for motor execution, whereas integrative, between-network
306 communication is critical for visual working memory (4).



307

308 **Figure 4. The effect of K-order on brain decoding and hierarchical organization of ChebNet.**

309 The effect of K -order on brain decoding was investigated by spanning over the list of $[0, 1, 2, 5, 10]$. The

310 decoding performance on $K=0$ was not shown in this figure due to its low overall performance

311 (decoding accuracy = 83.76%, 84.21%, 83.51% on training, validation and test sets). (B) High-order

312 decoding models showed a faster convergence speed during model training and also achieved better

313 decoding accuracy. Significant improvements were detected between $K=1$ (information integration

314 within the same network) and $K>1$ (transmission of brain activity among inter-connected brain

315 networks). (A) Variable sensitivity to the K -order was detected among different cognitive domains.

316 The effect of K -order on each cognitive domain was estimated by averaging the F1-score on the test

317 set. Among which, the Motor tasks showed stable decoding performance when increasing K while the

318 decoding of WM tasks gradually improved as increasing K . (C) A stable two-level organization
319 among ChebNet layers was revealed for the Motor tasks when increasing K . (D) For the Working-
320 memory task, it started with an unstable bipartition and gradually evolved into a tripartite organization
321 among ChebNet layers. The similarity of representations among ChebNet layers was calculated using
322 CKA with a linear kernel. The hierarchical clustering was then applied to the distance matrix ($dis = 1 -$
323 cka) using Ward's linkage and revealed the organizational principles among ChebNet layers.

324

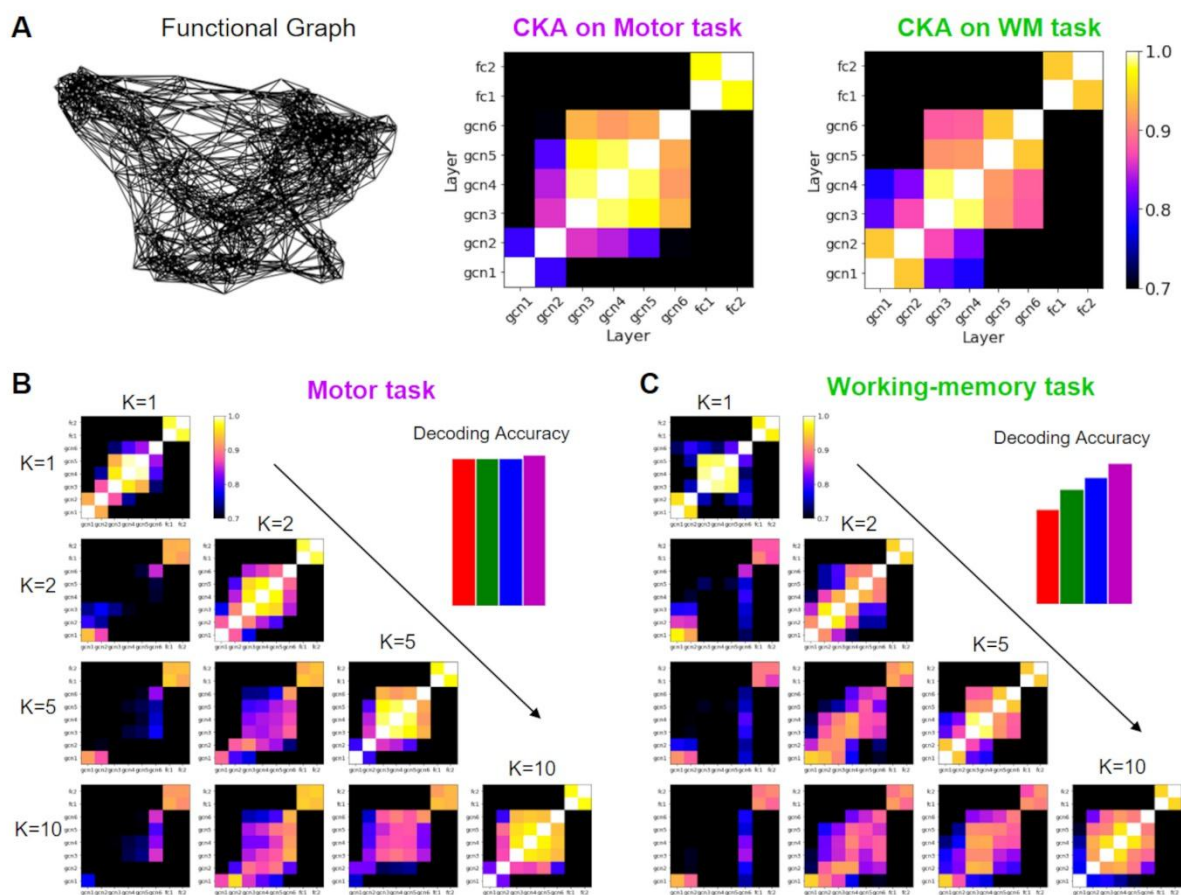
325 Functional integration in Working-memory tasks and segregation in Motor tasks

326 To further validate the functional segregation and integration hypothesis in brain decoding, we
327 conducted a systematic analysis on the decoding models at different K -orders by calculating the
328 similarity of representations between ChebNet models. We used the ChebNet- $K5$ model as the
329 reference model for the similarity analysis.

330 For Motor tasks, the ChebNet- $K1$ model already captured the low-to-high-level organization in graph
331 representations. Further increase in K did not change this organization but only caused redundant
332 representations in the high-level features (average similarity with $gcn6$ in $gcn3$ - $gcn5$ is $CKA=0.78$
333 and 0.92 for ChebNet- $K1$ and ChebNet- $K5$). A direct comparison between the two models (3rd row
334 and 1st column in Figure 5B) revealed that, compared to the ChebNet- $K5$ model, the ChebNet- $K1$
335 model captured highly similar low-level features ($CKA=0.92$ for $gcn1$ when comparing between
336 ChebNet- $K1$ and ChebNet- $K5$) and learned closely related high-level representations ($CKA=0.84$ for
337 $gcn6$ between the two models). However, very different hidden representations were learned in the
338 middle layers between the two models (averaged $CKA=0.70$ for $gcn2$ to $gcn5$). These results
339 indicated that the highly segregated brain function, such as the sensory and motor tasks, did not
340 involve high-level of information integration, but rather relied on neural transmission of brain activity
341 within a local area or segregated networks.

342 On the other hand, for the Working Memory tasks, the ChebNet- $K5$ model captured a nice
343 disassociation between low-level features ($gcn1$ to $gcn2$), hidden representations ($gcn3$ to $gcn4$), and
344 high-level features ($gcn5$ to $gcn6$). Such hierarchical organization was broken in the ChebNet- $K1$

345 model due to poor between-layer communication (i.e. big gaps in the representations between
 346 neighboring layers, CKA=0.98 and 0.69 for within- and between-level similarity in ChebNet-K1).
 347 Moreover, the ChebNet-K1 model successfully captured the low-level features by showing high
 348 similarity to ChebNet-K5 in the first two ChebNet layers, but it was not capable of encoding high-
 349 level representations in the last ChebNet layer (3rd row and 1st column in Figure 5C, compared
 350 between ChebNet-K1 and ChebNet-K5, CKA=0.93 for gcn1, 0.88 for gcn2, 0.74 for gcn6). By
 351 contrast, the ChebNet-K10 model learned very similar representations in the low, middle and high
 352 ChebNet layers as in ChebNet-K5 (4th row and 3rd column in Figure 5C, compared between
 353 ChebNet-K5 and ChebNet-K10, CKA=0.94 for gcn1, 0.90 for gcn6, average CKA=0.90 for gcn3 to
 354 gcn5). These results indicated that the high-order cognitive functions required a large scale of
 355 information propagation and integration on the brain graph, not only involving the local connections
 356 within a specific brain network ($K = 1$) but also engaging the long-range connections across multiple
 357 networks ($K \geq 5$).



359 **Figure 5. Similarity analysis of the decoding model with different K orders for the Motor and**
360 **Working-memory tasks.**

361 The similarity analysis of layer representations demonstrated a hierarchical organization among
362 stacked ChebNet layers in both Motor and Working-memory tasks (A). The decoding models were
363 built on the same functional graph derived from the resting-state functional connectivity. The
364 similarity of ChebNet representations was estimated not only between different layers in the same
365 model but also between different models. For the Motor task (B), the decoding model already
366 achieved the best performance at $K = 1$, and no further improvement on the decoding performance
367 when increasing the K -order. In terms of layer representations, the ChebNet- $K5$ model captured
368 similar low-level and high-level representations as ChebNet- $K1$ (3rd row and 1st column in B), but
369 learned different representations in the hidden layers. Besides, higher redundancy was captured in the
370 ChebNet- $K5$ model. This analysis indicated that ChebNet- $K1$ was enough to capture the functional
371 segregation in brain activity during Motor tasks. For the Working-memory task (C), the decoding
372 model showed high sensitivity to the choice of K -order and achieved the best decoding accuracy
373 when $K = 10$. In terms of layer representations, the ChebNet- $K1$ model captured similar low-level
374 representations as ChebNet- $K5$ (3rd row and 1st column in C), but learned very different hidden and
375 high-level representations. On the other hand, ChebNet- $K10$ was highly similar to ChebNet- $K5$ (4th
376 row and 3rd column in C), not only in the low-level representations (gcn1-gcn2), hidden
377 representations (gcn3-gcn5), as well as high-level representations (gcn6). This analysis indicated that
378 a high-order model was required in order to capture the complex forms of functional integration
379 during Working-memory tasks.

380

381 Improved functional alignment of cognitive states using graph convolution

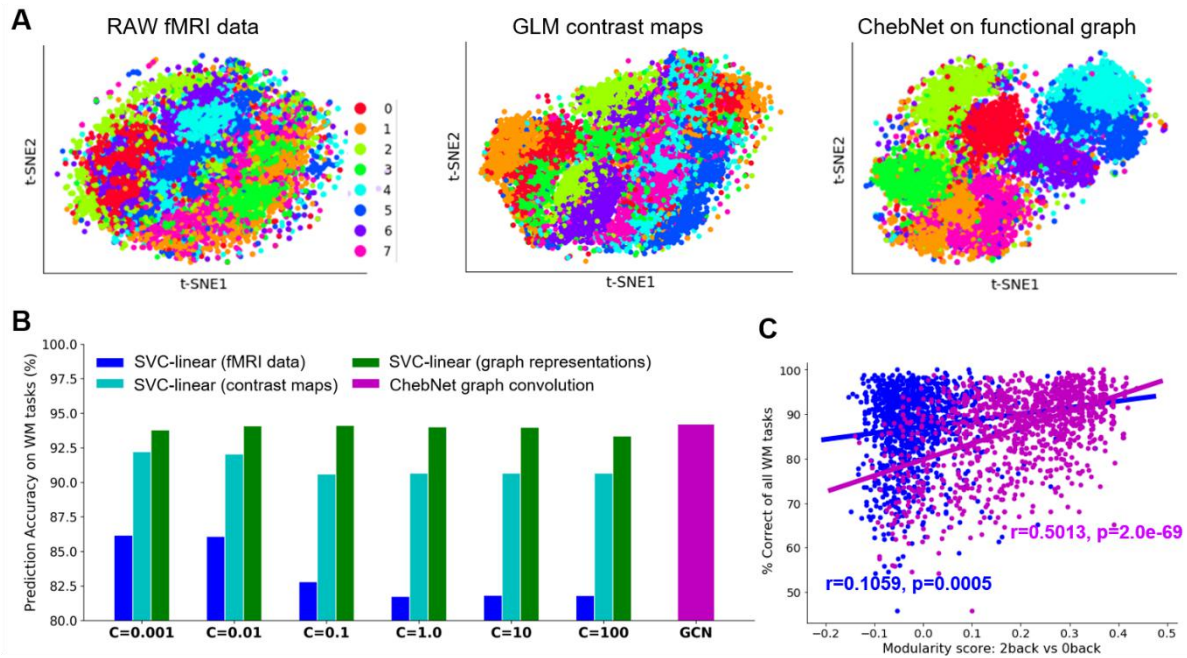
382 The layer representations in ChebNet improved inter-subject alignment of task-evoked brain
383 responses. For visualization purposes, we projected the feature representations of each ChebNet layer
384 onto a 2-dimensional space using t-SNE (25). Compared to raw fMRI data or the activation maps
385 derived from GLM analysis, the ChebNet representations of different task conditions were highly

386 clustered and easily separated from each other, demonstrating a strong effect of task segregation
387 (Figure 6A). The projections using other dimension reduction techniques were shown in Figure 6-S1,
388 including PCA, UMAP (McInnes et al., 2018), and PHATE (Moon et al., 2019). The segregation
389 effect was evaluated by calculating the modularity score (Q) of the state-transition graph on the
390 projections of layer representations. As we went deeper along ChebNet layers, the segregation effect
391 gradually strengthened and reached the peak in the last ChebNet layer. As shown in Figure 6-S2, for
392 the Motor task, a low segregation was detected in the raw fMRI data ($Q = 0.25$), with slightly higher
393 values in early ChebNet layers (e.g. $Q = 0.41$ in gcn1) and reaching the peak in the last ChebNet
394 layer ($Q = 0.60$ in gcn6). A similar level of task segregation was observed when using a high-order
395 ChebNet model, except for a faster convergence speed among ChebNet layers (Figure 6-S2B). For the
396 WM tasks, the segregation effect was evaluated separately for the memory-load and image category.
397 Interestingly, stronger segregation effect was detected among different image categories, e.g. place,
398 face, body and tool images, than between different levels of memory loads, e.g. 0back and 2back (e.g.
399 $Q = 0.55$ vs 0.38 in gcn6 respectively for the image category and memory load), but both higher than
400 the effects in the raw fMRI data ($Q = 0.06$ vs 0.03 respectively).

401 Association between ChebNet representations and behavioral performance

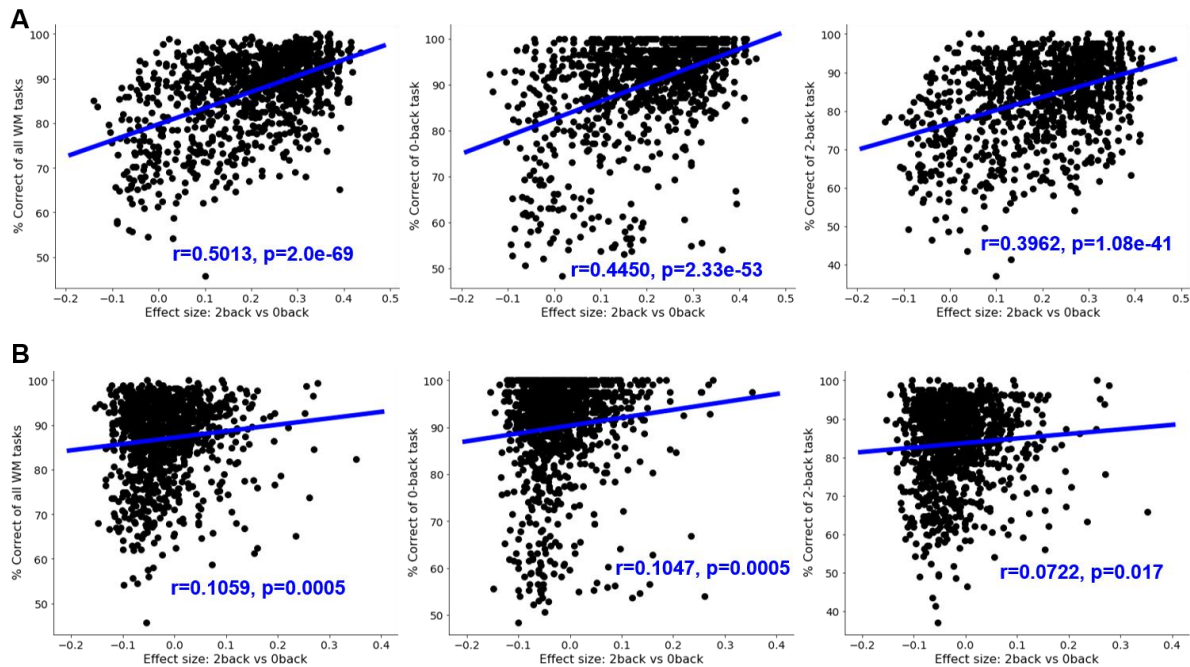
402 The segregation effect in the ChebNet representations not only boosted the decoding of cognitive states,
403 through better alignment of brain response across trials and subjects, but also improved the
404 association between behavior and brain organization, largely preserving individual variability in
405 cognitive processes. Specifically, the decoding model achieved much higher decoding accuracy when
406 using the ChebNet representations as features, regardless of the choices of classifiers and parameters,
407 for instance either using a multi-class support vector machine classification (SVC) or deep neural
408 networks (as shown in Figure 6B). In addition, the segregation in the representations of brain response
409 was significantly associated with behavioral performance during visual working memory task (as
410 shown in Figure 6C). It has been shown in previous studies that the modularity of state-transition on
411 individual fMRI data was significantly associated with participants' in-scanner task performance (26).
412 This association was also observed in our analysis. Moreover, we found a much stronger brain-

413 behavior association when constructing the state-transition graph based on ChebNet representations
414 rather than using raw fMRI data (Figure 6C). Specifically, the segregation of memory load in graph
415 representations was highly associated with subjects' in-scanner performance (as shown in Figure 7
416 and Figure 7-S1), including the average accuracy on all WM tasks ($r = 0.5013$, $p = 2.0e-69$), on 0back
417 tasks ($r = 0.4450$, $p = 2.33e-53$) and on 2back tasks ($r = 0.3962$, $p = 1.08e-13$), as well as the reaction
418 time on all WM tasks ($r = -0.2601$, $p = 5.74e-18$), on 0back tasks ($r = -0.3592$, $p = 1.13e-33$) and on
419 2back tasks ($r = -0.1173$, $p = 0.0001$). This analysis was done by using all subjects from the *HCP*
420 *SI200* database ($N = 1074$ of all subjects with available behavioral and imaging data for WM tasks).
421 The significant correlations were sustained after controlling for the effect of confounds including age,
422 gender, handedness and head motion ($r = 0.4659$, $p = 5.74e-59$ for the average accuracy; $r = -0.2552$,
423 $p = 2.0e-16$ for the reaction time). Moreover, the segregation of ChebNet representations as well as in-
424 scanner behavioral performance during WM tasks were significantly heritable in HCP population (h^2
425 $= 0.2882$ for ChebNet representations, $h^2 = 0.5624$ and 0.4118 for average accuracy and reaction time
426 in WM tasks, see Table S4 for all heritability estimates) and demonstrated significant shared genetic
427 variance in ChebNet representations and behavioral scores ($\rho_g = 0.80$ and -0.39 respectively for the
428 average accuracy and reaction time, see Table 1 for shared genetic influences in brain-behavioral
429 associations).



430
431 **Figure 6. ChebNet graph representation improved the functional alignment of cognitive states,**
432 **and induced higher decoding accuracy (B) and better prediction of task performance (C).**

433 Graph representations were extracted from the last ChebNet layer of the decoding model. Both fMRI
434 data and contrast maps were mapped onto the same brain atlas, i.e. Glasser's atlas (27) in the example,
435 by averaging the fMRI time-series or z-scores of task activation within each brain region. (A) All
436 three types of features, i.e. fMRI data, contrast maps and graph representations, were projected onto a
437 2-dimensional space using t-SNE (25) for the visualization purpose. Among them, graph
438 representations showed high distinctions among different task conditions. (B) The decoding of eight
439 working-memory tasks was re-evaluated by using multi-class support vector machine classification
440 (SVC) on these features. Among them, graph representations showed the highest decoding accuracy,
441 regardless of the chosen classifiers and parameters, e.g. linear classifier like SVC or nonlinear
442 classifier such as ChebNet. (C) The effect of task segregation was evaluated by the modularity score
443 on individual state-transition graph, as proposed by (26). We found a strong association between the
444 task segregation of ChebNet representations and participants' in-scanner task performance. The purple
445 line indicated the association of participants' in-scanner task performances with graph representations
446 derived from ChebNet, while the blue line indicated the association with raw fMRI data.



447

448

449 **Figure 7. Modularity scores in the state-transition graph significantly correlated with correct**
450 **responses during Working-memory tasks.**

451 The modularity score was calculated based on the state-transition graph of each subject, as proposed
452 by (26). Specifically, we first constructed a kNN graph from the t-SNE projections of fMRI signals (B)
453 or learned graph representations (A) of each subject. The modularity scores were then evaluated based
454 on the kNN graph with the partition provided by task conditions (e.g. 0back vs 2back). We found
455 significant correlations between the modularity scores of graph representations (A) and correct
456 responses during working-memory tasks (1st panel), 0back task conditions (2nd panel) and 2back task
457 conditions (3rd panel). Much weaker associations were detected in the raw fMRI data (B). The blue
458 lines indicated the linear regression models between the modularity score of the state-transition graph
459 and the average accuracy during task performance. The analysis was done among all subjects from
460 HCP S1200 release, with complete records of behavioral and imaging data for working memory tasks
461 (N=1074).

462

463

464 **Table 1: Shared genetic influences in ChebNet representations and behavioral scores.**

465 Bivariate genetic analyses were applied to quantify the shared genetic variance between ChebNet
466 representations of brain responses and behavioral measures. Strong associations of ChebNet
467 representations with the average accuracy (Acc) and reaction time (RT) during WM tasks were
468 observed, mainly due to shared genetic effects in brain response and behaviors. Both genetic and
469 phenotypic correlations reached a high-level of significance (FDR corrected). ***: $p < 0.001$, **: $p < 0.01$, *: $p < 0.05$
470

| | Phenotypic correlation (ρ_p) | Genetic correlation (ρ_g) |
|-----------------|---|--|
| WM_Task_Acc | 0.4659 *** | 0.7992 *** |
| WM_Task_2bk_Acc | 0.3716 *** | 0.7731 *** |
| WM_Task_0bk_Acc | 0.4189 *** | 0.8650 *** |
| WM_Task_RT | -0.2552 *** | -0.3895 ** |
| WM_Task_2bk_RT | -0.1173 *** | -0.2455 |
| WM_Task_0bk_RT | -0.3408 *** | -0.4967 ** |

471

472 **Discussion**

473 We proposed a generalized framework for brain decoding based on ChebNet graph convolutions. The
474 model takes in a short window of fMRI time series and a brain graph (with nodes representing brain
475 parcels and edges representing brain connectivity), and then annotates brain activity with fine
476 temporal resolution, and fine cognitive granularity. Using a 10s window of fMRI signals, our model
477 identified 6 cognitive domains with a test accuracy of 96%, and distinguished fine-grained cognitive
478 states on a trial basis with an accuracy above 93%, outperforming existing linear and nonlinear
479 decoding models. This gain in brain decoding was mainly contributed by high-level integration of
480 brain dynamics not only within limited subsets of brain regions but also between multiple brain
481 networks. Specifically, we used high-order ChebNet graph convolution to encode the complex forms
482 of functional integration during cognitive processes and captured hierarchical representations of brain
483 activities at different levels. This hierarchical organizational pattern as well as the decoding accuracy
484 was selectively impacted by the K-order in graph convolution, due to different organizational
485 principles of the cognitive tasks. For segregated brain function like motor execution, the K=1 model
486 achieved the best performance and revealed a stable 2-level hierarchy in neural representations. By
487 contrast, for high-order cognition such as visual working-memory tasks, the model plateaued at K=5
488 and uncovered a tripartite organization in neural activity. Our findings revealed the essential role of
489 functional integration in brain decoding, especially when decoding high-order cognition other than
490 sensory and motor functions.

491 **Functional segregation and integration in brain decoding**

492 Brain decoding has been a popular topic in neuroscience literature for decades since Haxby first
493 proposed the idea of recognition of different visual stimuli using brain activity from the visual cortex
494 (1,28). In the last decades, a variety of decoding models have been proposed with the aim to learn a
495 linear discriminative function on the spatial patterns of brain activations associated with different task
496 conditions. For instance, researchers have successfully attempted to use brain activity to reconstruct

497 the frames of movies (29), or to decode the semantic context from words (30) and visual scenes (31)
498 by using linear regression models. Recently, the fast development of deep artificial neural networks
499 (DNNs) has also drawn a lot of attention in neuroscience research. Several different DNN
500 architectures have been proposed to map human cognition from recorded brain activity, for instance
501 using classical convolutional (2) and recurrent neural networks (3), or a generalized form of
502 convolutions in the graph domain (17). However, the majority of brain decoding studies only utilized
503 the functional specialization hypothesis that aims to distinguish the localized brain activation patterns
504 by either training a linear classifier (16,32) or a nonlinear model through DNNs (2).
505 However, the majority of brain decoding studies so far only utilized the functional specialization
506 hypothesis that aims to distinguish the localized brain activities from a single brain region or a small
507 set of areas. This approach has shown promising results in the recognition of visual stimuli (28) and
508 decoding the direction of finger movements (33). It may suffer from limited decoding power when
509 dealing with large populations and a variety of cognitive states that involve not only motor and
510 sensory perception but also high-order cognition (16). Such large-scale decoding is still challenging
511 and may require a large collection of brain imaging data and incorporating brain responses from the
512 whole brain in the decoding model, including both local and global information (34). So far, the
513 functional integration at the whole-brain has been largely ignored in the brain decoding literature, but
514 started to draw the attention of neuroscientists. Cole and colleagues (35) first showed that the
515 information flow within functional networks was able to predict brain activation during cognitive
516 tasks, specifically to predict activation patterns of unseen brain regions from regions of the same
517 network. A similar idea was recently used in (3) by first extracting an integrated signal from each of
518 90 resting-state networks and then inferring brain states based on the temporal dependencies of these
519 brain signals. Following this line of work, we recently proposed a graph convolutional network to
520 decode brain states by propagating temporal dynamics of brain activity based on functional networks
521 (17). In the present study we further extended this framework by exploring more variants in the graph
522 convolutional network architecture. By using high-order graph convolutions, the model projected the
523 spatiotemporal dynamics of cognitive processes onto a new representational space and integrated the

524 context of brain activity in both local and global extent, ranging from brain region to functional
525 networks and towards the whole brain.

526 Compared to previous linear and nonlinear decoding models, our proposed decoding model provided
527 a generalized solution over a large population and a variety of cognitive domains. Besides, our model
528 outperformed other approaches on the same dataset (as shown in Table S3), most of which followed
529 the functional segregation assumption by predicting cognitive states from localized features of each
530 brain parcel covering the entire cerebral cortex. After incorporating the network architecture of the
531 human brain and integrating information flow within functional networks (e.g. first-order GCN), the
532 classification accuracy was largely improved (90%, as stated in (17)). The decoding accuracy was
533 further improved after taking into account the high-order interactions on the graph, not only within
534 functional networks but also across multiple brain systems (93% using the 5-order ChebNet model).

535 Our results suggest that not only the segregated brain activations played an important part in
536 distinguishing between cognitive processes as illustrated in previous brain decoding literature, but the
537 functional integration within and between brain networks can also contribute to the classification of
538 cognitive states to some degree. The tradeoff between functional segregation and integration largely
539 depends on the nature of cognitive processes, for instance, localized brain signatures from motor and
540 sensory cortex for Motor tasks (as shown in Figure 2A), while complex forms of functional
541 interactions among multiple brain systems during WM tasks (as shown in Figure 2B). Their relations
542 were automatically captured during the training of deep neural networks. Coinciding with this
543 hypothesis, the decoding model achieved excellent performance in distinguishing different types of
544 body movements when only considered the local context of brain activity either from a local area
545 (94.7% in (2)) or within a functional network (96.6% in (3)). A similar level of performance was
546 achieved when using either first-order or high-order graph convolutions (95.6% when using 10s fMRI
547 signals). By contrast, when classifying 0-back and 2-back WM tasks, much higher classification
548 errors were detected only using local brain activity (14% in (2)) compared to functional integration
549 within the functional networks (10% in (3)). The classification errors among WM tasks were highly
550 reduced when applying graph convolutions ($\leq 9\%$ when using ChebNet-K1, $\leq 4\%$ when using the
551 ChebNet-K5, as shown in Figure 4-S1).

552 To conclude, our results demonstrated that an efficient brain decoder not only involved functional
553 segregation, e.g. distinguishing localized brain activation during body movements, but also engaged
554 functional integration, e.g. integrating brain activity among multiple brain networks during visual
555 working memory tasks.

556 Saliency of brain decoding goes beyond brain activation

557 Both saliency maps and brain activations aimed to reveal the neural substrates of cognitive processes.
558 They also shared some common features, for instance, both relying on task-evoked brain responses
559 and showing selective responses to different task conditions. However, their relations need to be
560 addressed with caution. Brain activation was commonly used in neuroscience research to study the
561 neural basis of cognitive processes by convolving neural activity with a canonical hemodynamic
562 response function and to find the localization of each cognitive function using a generalized linear
563 model (GLM) approach. However, as stated in Poldark's paper (34), not all brain activations were
564 "diagnostic" in terms of brain state prediction, i.e. distinguishing among different cognitive processes.
565 To address this issue, we used saliency maps to detect the important features that show high
566 contributions to the prediction.

567 First, common brain regions were revealed in both approaches, i.e. areas that not only strongly
568 activated during task performance (in GLM analysis) but also largely contributed to the classification
569 of different tasks (in saliency maps). For example, salient features in the sensorimotor cortex were
570 detected for motor execution and in the ventral visual stream for image recognition during WM tasks
571 (Figure 2 A and B). These brain regions have been well validated in the literature, that the primary
572 motor cortex was engaged during movements of the human body (36) and ventral temporal cortex was
573 responsible for the recognition of face and place images (37). Most previous decoding studies were
574 based on this set of brain regions, for instance, using brain activity of the visual cortex to decode the
575 category of seen images including faces vs objects as well as different animal species (28).

576 Second, some inconsistent findings were reported by the two techniques. On one hand, low saliency
577 values did not mean no brain activation. Instead, the regions might be activated for all task conditions
578 (i.e. low saliency but showing high activation on all tasks). For example, the visual cortex was

579 consistently activated during the Motor task (38), in response to the presentation of the cue images
580 and visual instructions (e.g. images of hand, foot and tongue in the cue phase). But this activation
581 pattern was not related to actual movements and not informative to distinguish between hand and
582 tongue movements, i.e. much lower decoding accuracy in the cue phase compared to the movement
583 blocks, as indicated in the Figure 6A in (17). Consequently, the visual cortex was not detected in the
584 saliency maps of Motor tasks (Figure 2A). On the other hand, some brain regions with absence of
585 strong activations might show high saliency to the prediction (i.e. high saliency but with low
586 activation scores). For instance, high saliency values were detected in the bilateral area OP4 for
587 tongue movements and in the left area PHA3 for recognition of both place and face images. These
588 patterns were not detected by random but instead highly consistent across a number of subjects
589 (Figure 2C and D). By contrast, when using the GLM approach, these regions were not detected in
590 either the group activation maps from HCP database (Figure 5-S1 in (17)) or meta-analysis from a
591 collection of previous studies (Figure 5-S2 in (17)). Specifically, weak activation was detected in the
592 left but not right OP4 for the tongue movement (T-score=4.6 and -1.66 respectively for the left and
593 right OP4 in the contrast maps of tongue vs rest, using group activation maps from HCPS500 release),
594 while area PHA3 was not even activated for the recognition of face and place images (T-score= -1.78
595 and -7.02 respectively for the contrast maps of face and place images vs rest). One possible
596 explanation of this is that task-evoked activities in these brain regions did not follow the shape of the
597 canonical hemodynamic response function and thus were not detected by the GLM approach and not
598 shown in the contrast maps. However, these brain regions still showed distinctive patterns of response
599 to different cognitive tasks, and thus were detected in the saliency maps when such constraint of
600 temporal dynamics was not applied in the decoding model. Our results suggest that not only brain
601 activations but also deactivations or even brain responses not shown in the contrast maps might highly
602 contribute to the classification of cognitive states. In addition, the detected salient features were
603 highly consistent across subjects and showed selective responses to different cognitive states, may
604 uncover the biological basis of the decoding model and shed a light on the anatomical and functional
605 substrates of cognitive processes.

606

607 To conclude, our results suggested that a more generalized framework was required for brain
608 decoding which does not rely on localized brain activations or spatial patterns of contrast maps, but
609 instead decoding information from task-evoked brain responses of the whole brain. Consequently, the
610 model was able to distinguish the neural dynamics of various cognitive functions in both spatial and
611 temporal domains and learn new representations of brain organization during cognitive processes (39).
612 The spatiotemporal graph convolution provided a promising solution for this problem by leveraging
613 our prior knowledge on brain organization using a graph-based model. Instead of classifying patterns
614 of brain activity within a local area, graph convolution takes into account the functional interactions
615 of neural dynamics across multiple networks and projects the spatiotemporal dynamics of cognitive
616 processes onto a new representational space. Moreover, ChebNet graph convolution naturally
617 incorporates both functional segregation and integration for brain decoding, i.e. distinguishing
618 localized brain activities from a subset of brain regions or within a specific brain network (first-order
619 convolution) and encoding complex forms of functional interactions among multiple brain systems
620 (high-order convolution), in line with different organizational principles of cognitive processes. As a
621 result, the learned representations improved the functional alignment among trials and subjects and
622 therefore increased the decoding of cognitive states. More importantly, by largely preserving the
623 individual variability in brain organization, the ChebNet representations achieved better associations
624 with human behaviors during task, demonstrating shared genetic influences in brain responses and
625 behaviors. The present work suggests the feasibility of large-scale multi-domain decoding with full-
626 brain models, opening new avenues for modelling of naturalistic tasks such as movies or video games.

627

628 **Materials and Methods**

629 fMRI Datasets and Preprocessing

630 We used the block-design task-fMRI dataset from the Human Connectome Project S1200 release
631 (https://db.humanconnectome.org/data/projects/HCP_1200). The minimal preprocessed fMRI data in
632 CIFTI formats were selected. The preprocessing pipelines includes two steps (40): 1) fMRIVolume
633 pipeline generates “minimally preprocessed” 4D time-series (i.e. “.nii.gz” file) that includes gradient
634 unwarping, motion correction, fieldmap-based EPI distortion correction, brain-boundary-based
635 registration of EPI to structural T1-weighted scan, non-linear (FNIRT) registration into MNI152 space,
636 and grand-mean intensity normalization. 2) fMRISurface pipeline projects fMRI data from the cortical
637 gray matter ribbon onto the individual brain surface and then onto template surface meshes (i.e.
638 “dtseries.nii” file), followed by surface-based smoothing using a geodesic Gaussian algorithm.
639 Further details on fMRI data acquisition, task design and preprocessing can be found in (38,40). The
640 task fMRI database includes six cognitive domains, which are emotion, language, motor, relational,
641 social, and working memory. In total, there are 21 different experimental conditions. We excluded the
642 two gambling conditions in our analysis due to the short event design of the gambling trials (1.5s for
643 button press, 1s for feedback and 1s for ITI). The detailed description of the task paradigms as well as
644 the selected cognitive domains can be found in (17,38)

645 Decoding brain activity using graph convolution

646 A brain graph provides a network representation of brain organization by associating nodes with brain
647 regions and defining edges via anatomical or functional connections (41). We recently found that
648 convolutional operations on brain graph can be used to encode the within-network interactions of
649 task-evoked brain responses and to decode a large number of cognitive tasks (17). Here, we proposed
650 a more generalized form of graph convolution by using Chebyshev polynomials and explored how
651 functional segregation and high-order functional interactions affects brain decoding.

652 Step 1: Construction of brain graph

653 The decoding pipeline started with a weighted graph $\mathcal{G} = (\mathcal{V}, \mathcal{E}, W)$, where \mathcal{V} is a parcellation of
654 cerebral cortex into N regions, \mathcal{E} is a set of connections between each pair of brain regions, with its
655 weights defined as $W_{i,j}$. Many alternative approaches can be used to build such brain graph \mathcal{G} , for
656 instance using different brain parcellation schemes and constructing various types of brain
657 connectomes (for a review, see (41)). Here, we used the Glasser's multi-modal parcellation, consisting
658 of 360 areas in the cerebral cortex, bounded by sharp changes in cortical architecture, function,
659 connectivity, and topography (27). The edges between each pair of nodes were estimated by
660 calculating the group averaged resting-state functional connectivity (RSFC) based on minimal
661 preprocessed resting-state fMRI data from $N = 1080$ HCP subjects (Glasser et al., 2013). Additional
662 preprocessing steps were applied before the calculation of RSFC, including regressing out the signals
663 from white matter and csf, and bandpass temporal filtering on frequencies between 0.01 to 0.1 HZ.
664 Functional connectivity was calculated on individual brains using Pearson correlation and then
665 normalized using Fisher z-transform before averaging among the entire group of subjects. After that, a
666 k-nearest-neighbor (k-NN) graph was built by only connecting each node to its 8 neighbors with the
667 highest connectivity strength.

668 Step 2: Mapping of brain activity onto the graph

669 After the construction of brain graph (i.e. defining brain parcels and edges), for each functional run
670 and each subject, the preprocessed task-fMRI data was then mapped onto the set of brain parcels,
671 resulting in a 2-dimensional time-series matrix. This time-series matrix was first split into multiple
672 task blocks according to fMRI paradigms and then cut into sets of time-series of the chosen window
673 size (e.g. 10 second). Shorter time windows were discarded in the process. The remaining time-series
674 were treated as independent data samples during model training. As a result, we generated a large
675 number of fMRI time-series matrices from all cognitive domains, i.e. a short time-series with duration
676 of T for each of N brain parcels, $x \in \mathbb{R}^{N \times T}$. The entire dataset consists of over 1000 subjects for
677 each cognitive domain (see Table S2 for detailed information), in total of 14,895 functional runs

678 across the six cognitive domains, and 138,662 data samples of fMRI signals $x \in \mathbb{R}^{N \times T}$ when using a
 679 10s time window (i.e. 15 functional volumes at TR=0.72s).

680 Step 3: Spatiotemporal graph convolutions using ChebNet

681 Graph convolution relies on the graph Laplacian, which is a smooth operator characterizing the
 682 magnitude of signal changes between adjacent nodes. The normalized graph Laplacian is defined as:

$$683 \quad L = I - D^{-1/2} W D^{-1/2} \quad (\text{Eq. 1})$$

684 where D is a diagonal matrix of node degrees, I is the identity matrix, and W is the weight matrix. The
 685 eigendecomposition of Laplacian matrix is defined as $L = U \Delta U^T$, where $U = (u_0, u_1, \dots, u_{N-1})$ is the
 686 matrix of Laplacian eigenvectors and is also called graph Fourier modes, and $\Delta =$
 687 $\text{diag}(\lambda_0, \lambda_1, \dots, \lambda_{N-1})$ is a diagonal matrix of the corresponding eigenvalues, specifying the frequency
 688 of the graph modes. In other words, the eigenvalues quantify the smoothness of signal changes on the
 689 graph, while the eigenvectors indicate the patterns of signal distribution on the graph.

690 For a signal x defined on graph, i.e. assigning a feature vector to each brain region, the convolution
 691 between the graph signal $x \in \mathbb{R}^{N \times T}$ and a graph filter $g_\theta \in \mathbb{R}^{N \times T}$ based on graph \mathcal{G} , is defined as
 692 their element-wise Hadamard product in the spectral domain, i.e.:

$$693 \quad x * g_\theta = U(U^T g_\theta) \odot (U^T x) = U G_\theta U^T x \quad (\text{Eq. 2})$$

694 where $G_\theta = \text{diag}(U^T g_\theta)$ and θ indicate a parametric model for graph convolution g_θ , $U =$
 695 $(u_0, u_1, \dots, u_{N-1})$ is the matrix of Laplacian eigenvectors and $U^T x$ is actually projecting the graph
 696 signal onto the full spectrum of graph modes. To avoid calculating the spectral decomposition of the
 697 graph Laplacian, ChebNet convolution (Defferrard et al., 2016) uses a truncated expansion of the
 698 Chebychev polynomials, which are defined recursively by:

$$699 \quad T_k(x) = 2xT_{k-1}(x) - T_{k-2}(x), \quad T_0(x) = 1, T_1(x) = x \quad (\text{Eq. 3})$$

700 Consequently, the ChebNet graph convolution is defined as:

$$701 \quad x * g_\theta = \sum_{k=0}^K \theta_k T_k(\tilde{L})x \quad (\text{Eq. 4})$$

702 where $\tilde{L} = 2L/\lambda_{max} - I$ is a normalized version of graph Laplacian with λ_{max} being the largest
 703 eigenvalue, θ_k is the model parameter to be learned at each order of the Chebychev polynomials. It
 704 has been proved that the ChebNet graph convolution was naturally K -localized in space by taking up

705 to Kth order Chebychev polynomials (18), which means that each ChebNet convolutional layer
706 integrates the context of brain activity within a K-step neighborhood.

707 Brain-decoding pipeline using ChebNet graph convolution

708 We used a similar decoding pipeline as proposed in (17), consisting of 6 graph convolutional layers
709 with 32 graph filters at each layer, followed by a flatten layer and 2 fully connected layers (256, 64
710 units). The model takes in a short series of fMRI volumes as input, maps the fMRI data onto the
711 predefined brain graph and results in a 2-dimensional time-series matrix $X^1 \in \mathbb{R}^{N \times T}$, i.e. a short
712 time-series with duration of T for each of N brain parcels at the first ChebNet layer. The first ChebNet
713 layer learns various shapes of temporal convolution kernels by treating multiple time steps as input
714 channels ($C = T$) and propagates such temporal dynamics within ($K=1$) and between ($K>1$) brain
715 networks. As a result, a set of “brain activation” maps are generated (see Figure 3-S1) and passed onto
716 the next ChebNet layer for higher-order information integration (see Figure 3-S2), and so on. The
717 learned graph representations in the last ChebNet layer were then imported to a 2-layer multilayer
718 perceptron (MLP) to predict the cognitive state.

719 The entire dataset was split into training (60%), validation (20%), test (20%) sets using a subject-
720 specific split scheme, which ensures that all fMRI data from the same subject was assigned to only
721 one of the three sets. Approximately, the training set includes fMRI data from 700 unique subjects
722 (depending on data availability for different cognitive tasks ranging from 1043 to 1085 subjects, see
723 Table S2), with 176 subjects for validation set and 219 subjects for test set. Specifically, we used the
724 training set to train model parameters, the validation set to evaluate the model at the end of each
725 training epoch, and saved the best model with the highest prediction score on the validation set after
726 100 training epochs. The saved model was evaluated on the test set and reported the final decoding
727 performance. We used Adam as the optimizer with the initial learning rate as 0.0001 on all cognitive
728 domains. Additional l2 regularization of 0.0005 on weights was used to control model overfitting and
729 the noise effect of fMRI signals. Dropout of 0.5 was additionally applied to the fully connected layers.
730 The implementation of the ChebNet graph convolution was based on PyTorch 1.1.0, and was made
731 publicly available in the following repository: https://github.com/zhangyu2ustc/gcn_tutorial_test.git.

732 Effects of K -order in ChebNet brain decoding

733 ChebNet graph convolution used truncated expansion of Chebychev polynomials of order K for the
734 approximation of graph convolution in the spectral space (see Eq. 4). The choice of K -order controls
735 the scale of the information integration on the graph. When $K = 0$, $x * g_\theta = \theta_0 x$, which indicates a
736 global scaling factor on the input signal x by treating each node independently, similar to the classical
737 mass univariate analysis for brain activation detection. When $K = 1$, $x * g_\theta = \theta_0 x + \theta_1 \tilde{L}x$, which
738 indicates information integration between the direct neighbors and the current node on the graph
739 (integrating signals within the same network). When $K = 2$, $x * g_\theta = \theta_0 x + \theta_1 \tilde{L}x + \theta_2 \tilde{L}^2 x$, which
740 indicates information integration within a two-step neighborhood on the graph (consisting of
741 information from local area, within network and between networks). Generally speaking, when $K > 1$,
742 the graph convolution integrates the information within a K -step neighbourhood by propagating graph
743 signals not only within the same network but also among inter-connected brain networks. Thus, the K -
744 order controls the propagation rate of information flow on the brain graph. We explored different
745 choices of K -order in ChebNet spanning over the list of $[0,1,2,5,10]$ and found a significant boost in
746 both brain decoding and representational learning by using high-order graph convolutions.

747 Saliency map of the decoding model: contribution of brain regions

748 The saliency map analysis aims to locate which part of the brain (or input features) helps to
749 differentiate between cognitive tasks. We used a gradient approach named Guided BackProp (20) to
750 visualize the contribution of inputs. This approach has been commonly used to visualize a deep neural
751 network and easily generalized to graph convolutions. The basic idea is that if an input is relevant, a
752 little variation on it will cause high change in the layer activation. This can be characterized by the
753 gradient of the output given the input, with the positive gradients indicating that a small change to the
754 input signals increases the output value. Specifically, for the graph signal X^l of layer l and its gradient
755 R^l , the overwritten gradient $\nabla_{X^l} R^l$ can be calculated as follows:

$$756 \quad \nabla_{X^l} R^l = (X^l > 0) \odot (\nabla_{X^{l+1}} R^{l+1} > 0) \odot \nabla_{X^{l+1}} R^{l+1} \quad (\text{Eq. 5})$$

757 In order to generate the saliency map, we started from the output layer of a pre-trained model and
758 used the above chain rule to propagate the gradients at each layer until reaching the input layer. This
759 guided-backpropagation approach can provide a high-resolution saliency map which has the same
760 dimension as the input data. We further calculated a heatmap of saliency by taking the variance across
761 all time steps for each parcel, considering that the variance of the saliency curve provides a simplified
762 way to evaluate the contribution of task-evoked hemodynamic response. To make it comparable
763 across subjects, the saliency value was additionally normalized to the range [0,1], with its highest
764 value at 1 (a dominant effect for task prediction) and lowest at 0 (no contribution to task prediction).

765 Similarity analysis of layer representations in graph convolutions

766 ChebNet graph convolution maps the spatiotemporal dynamics of fMRI brain activity onto a new
767 representational space in the spectral domain. Different representations were learned at each ChebNet
768 graph convolutional layer by integrating the information flow within ($K = 1$) and between networks
769 ($K > 1$). Besides, by using a multi-layer architecture, the scale of information integration was
770 gradually enhanced, ranging from a local area (first ChebNet layer) to the whole-brain (last ChebNet
771 layer). For a better understanding of the ChebNet models, we analyzed the similarity of learned
772 representations between ChebNet layers as well as across different decoding models (e.g. using
773 different K-orders). Considering the high-dimensional nature of learned representations (360*32 in
774 our case), we evaluated the cross-layer and cross-model similarity using centered kernel alignment
775 (CKA) with a linear kernel, which was proposed to compare layer representations of deep neural
776 networks, not only in the same network trained from different initializations, but also across different
777 models (21). Linear CKA is closely related to CCA and linear regression. Studies showed that CKA
778 was invariant to orthogonal transformation and isotropic scaling, and consistently identified the
779 correspondences of representations between layers, and thus can reveal pathology in neural networks
780 representations (21). Here, we used CKA to evaluate the effects of K-orders in ChebNet brain
781 decoding for both Motor and Working-memory tasks. First, we extracted the layer representations
782 from each ChebNet layer by passing all data samples from the test set into the pre-trained decoding
783 model and reshaped the representations (samples * nodes * channels) into a matrix

784 $X \in \mathbb{R}^{samples \times features}$. Then, the linear CKA of two representation matrices X and Y, either from
785 different layers or different models, was defined as:

$$786 \quad CKA(X, Y) = \|Y^T X\|_F^2 / (\|X^T X\|_F \|Y^T Y\|_F) \quad (\text{Eq. 5})$$

787 where $\|A\|_F = \sqrt{\sum_{i,j} a_{i,j}^2}$ indicates the Frobenius norm of the cross-correlation matrix A. The CKA
788 value was within the range [0,1], with its highest value at 1 (the same layer representation) and lowest
789 at 0 (totally different layer representations).

790 Projections of layer representations using t-SNE

791 For visualization purposes, we projected the high-dimensional layer representations to a 2-
792 dimensional (2D) space by using t-SNE (25). Based on the t-SNE projections, we calculated the
793 modularity score among task conditions as a measure of task segregation, representing the cost of
794 brain state transition between task conditions. It has been shown that the modularity score on the
795 shape graph constructed from individual fMRI data was significantly associated with participants' in-
796 scanner task performance (26). Here, we estimated the modularity score on t-SNE projections derived
797 from not only fMRI signals but also layer representations of graph convolutional networks.
798 Specifically, fMRI signals and layer representations were first mapped onto a 2D space by using t-
799 SNE. Then, a k-NN graph (k=5) was constructed based on the coordinates of t-SNE projections by
800 connecting each data sample with its five nearest neighbors in the 2D space. After that, a segregation
801 index was defined by calculating the modularity score (Q) based on the partition of communities
802 using task conditions, with a high separation value indicating more edges within the same task
803 condition that expected by chance (42).

$$804 \quad Q = \frac{1}{4m} \sum_{i,j} \sum_t (A_{ij} - \frac{k_i k_j}{2m}) \delta(c_i, c_j) \quad (\text{Eq. 6})$$

805 where k_i and k_j are the degrees of the nodes on the kNN graph, $m = \frac{1}{2} \sum_i k_i$ is the total number of
806 edges on the kNN graph, and $\delta(c_i, c_j)$ indicates whether node i and node j belong to the same
807 community (or task condition). The value of the task segregation index was within the range [-0.5,1],

808 with the value close to 1 indicating a strong community structure and a positive value indicating the
809 number of edges within the same task conditions exceeds the number expected by chance (e.g. on a
810 random graph). The same procedure could also be applied to individual subject, i.e. constructing a
811 kNN graph using t-SNE projections of fMRI signals or layer representations from the same subject.
812 Their association with participants' in-scanner task performance were also investigated by calculating
813 the Pearson correlation coefficient of individual segregation index (i.e. modularity score on fMRI
814 signals or layer representations) with the correct response and reaction time during working-memory
815 tasks.

816 Heritability analysis of brain representations and behavioral performances

817 For the heritability estimates of brain response and behavioral performance during WM tasks, we used
818 the Sequential Oligogenic Linkage Analysis Routines (SOLAR) Eclipse software package
819 (http://www.nitrc.org/projects/se_linux). SOLAR relies on the maximum variance decomposition
820 of the covariance matrix Ω for a pedigree:

$$821 \quad \Omega = 2\Phi\sigma_g^2 + I\sigma_e^2 \quad (\text{Eq. 7})$$

822 where σ_g^2 is the genetic variance due to the additive genetic factors, Φ is the kinship matrix
823 representing the pairwise kinship coefficients among all individuals, σ_e^2 is the variance due to
824 individual-specific environmental effects and measurement error, and I is an identity matrix. Narrow
825 sense heritability is defined as the fraction of phenotypic variance σ_p^2 attributable to additive genetic
826 factors: $h^2 = \sigma_g^2 / \sigma_p^2$. The significance of the heritability estimate is tested by comparing the
827 likelihood of the model in which σ_g^2 is constrained to zero with that of a model in which σ_g^2 is
828 estimated. Prior to the heritability estimation, all phenotype (brain and behavioral phenotypes) were
829 adjusted for covariates including age, gender, handedness and head motion.

830 The heritability estimate was applied on 1070 subjects from HCP S1200 release with available
831 behavioral and imaging data for WM tasks, which consist of 447 unique families, including 143
832 monozygotic-twin pairs, 83 dizygotic-twin pairs and 290 non-twin siblings.

833 We further performed the bivariate genetic analyses to quantify the shared genetic variance and the
834 phenotypic correlation between brain responses and behavioral measures, relying on the following
835 model:

$$836 \quad \rho_p = \sqrt{h_a^2} \sqrt{h_b^2} \cdot \rho_g + \sqrt{1 - h_a^2} \sqrt{1 - h_b^2} \cdot \rho_e \quad (\text{Eq. 8})$$

837 where Pearson's phenotypic correlation ρ_p is decomposed into ρ_g and ρ_e , where ρ_g is the proportion
838 of variability due to shared genetic effects and ρ_e is that due to the environment, while h_a^2 and h_b^2
839 correspond to the narrow sense heritability for phenotypes a (representation of brain response) and b
840 (behavioral scores).

841

842 **Acknowledgment**

843 This work was partially supported by the Major Scientific Project of Zhejiang Lab (No.
844 2020ND8AD01), Courtois foundation through the Courtois NeuroMod Project and the IVADO
845 Postdoctoral Scholarships Program. PB is supported by a salary award of “Fonds de recherche du
846 Québec - Santé”, chercheur boursier junior 2.

847 **Author contributions**

848 Conceptualization: YZ, PB; Methodology: YZ, PB; Visualization: YZ, PB;
849 Investigation: YZ, NF, AD, PB;
850 Writing—original draft: YZ, NF, PB
851 Writing—review & editing: YZ, NF, AD, PB

852 **Competing interests**

853 The authors declare no competing financial interests.

854 **Data and materials availability**

855 We used the block-design task-fMRI dataset from the Human Connectome Project S1200 release,
856 downloaded from https://db.humanconnectome.org/data/projects/HCP_1200. In total, fMRI data from
857 1095 unique subjects under six different task domains and resting-state were used in this study. The
858 minimal preprocessed fMRI data of the CIFTI format were used, which maps individual fMRI time-
859 series onto the standard surface template with 32k vertices per hemisphere. Our decoding pipeline, as
860 well as the optimized decoding models and the construction of brain graphs, were made publicly
861 available in the following repository: https://github.com/zhangyu2ustc/gcn_tutorial_test.git

862 **References**

- 863 1. Haxby JV, Gobbini MI, Furey ML, Ishai A, Schouten JL, Pietrini P. Distributed and
864 overlapping representations of faces and objects in ventral temporal cortex. *Science*.
865 2001 Sep 28;293(5539):2425–30.
- 866 2. Wang X, Liang X, Jiang Z, Nguchu BA, Zhou Y, Wang Y, et al. Decoding and mapping
867 task states of the human brain via deep learning. *Hum Brain Mapp*. 2020 Apr
868 15;41(6):1505–19.
- 869 3. Li H, Fan Y. Interpretable, highly accurate brain decoding of subtly distinct brain states
870 from functional MRI using intrinsic functional networks and long short-term memory
871 recurrent neural networks. *NeuroImage*. 2019 Nov 15;202:116059.
- 872 4. Cohen JR, D'Esposito M. The Segregation and Integration of Distinct Brain Networks
873 and Their Relationship to Cognition. *J Neurosci*. 2016 Nov 30;36(48):12083–94.
- 874 5. Deco G, Tononi G, Boly M, Kringelbach ML. Rethinking segregation and integration:
875 contributions of whole-brain modelling. *Nat Rev Neurosci*. 2015 Jul;16(7):430–9.
- 876 6. Friston KJ. Functional and effective connectivity in neuroimaging: A synthesis. *Hum*
877 *Brain Mapp*. 1994;2(1–2):56–78.
- 878 7. Tononi G, Sporns O, Edelman GM. A measure for brain complexity: relating functional
879 segregation and integration in the nervous system. *Proc Natl Acad Sci*. 1994 May
880 24;91(11):5033–7.
- 881 8. Mottolese C, Richard N, Harquel S, Szathmari A, Sirigu A, Desmurget M. Mapping
882 motor representations in the human cerebellum. *Brain*. 2013 Jan 1;136(1):330–42.
- 883 9. Haxby JV, Connolly AC, Guntupalli JS. Decoding Neural Representational Spaces
884 Using Multivariate Pattern Analysis. *Annu Rev Neurosci*. 2014 Jul 8;37(1):435–56.
- 885 10. Eriksson J, Vogel EK, Lansner A, Bergström F, Nyberg L. Neurocognitive Architecture
886 of Working Memory. *Neuron*. 2015 Oct 7;88(1):33–46.
- 887 11. Harrison SA, Tong F. Decoding reveals the contents of visual working memory in early
888 visual areas. *Nature*. 2009 Apr;458(7238):632–5.
- 889 12. Riggall AC, Postle BR. The Relationship between Working Memory Storage and
890 Elevated Activity as Measured with Functional Magnetic Resonance Imaging. *J*
891 *Neurosci*. 2012 Sep 19;32(38):12990–8.
- 892 13. Christophel TB, Hebart MN, Haynes J-D. Decoding the Contents of Visual Short-Term
893 Memory from Human Visual and Parietal Cortex. *J Neurosci*. 2012 Sep
894 19;32(38):12983–9.
- 895 14. Sligte IG, van Moorselaar D, Vandenbroucke ARE. Decoding the Contents of Visual
896 Working Memory: Evidence for Process-Based and Content-Based Working Memory
897 Areas? *J Neurosci*. 2013 Jan 23;33(4):1293–4.
- 898 15. Poldrack RA, Mumford JA, Schonberg T, Kalar D, Barman B, Yarkoni T. Discovering
899 Relations Between Mind, Brain, and Mental Disorders Using Topic Mapping. Sporns O,
900 editor. *PLoS Comput Biol*. 2012 Oct 11;8(10):e1002707.

- 901 16. Poldrack RA, Halchenko Y, Hanson SJ. Decoding the Large-Scale Structure of Brain
902 Function by Classifying Mental States Across Individuals. *Psychol Sci.* 2009
903 Nov;20(11):1364–72.
- 904 17. Zhang Y, Tetrel L, Thirion B, Bellec P. Functional annotation of human cognitive states
905 using deep graph convolution. *NeuroImage.* 2021 May 1;231:117847.
- 906 18. Defferrard M, Bresson X, Vandergheynst P. Convolutional Neural Networks on Graphs
907 with Fast Localized Spectral Filtering. *Adv Neural Inf Process Syst* 29 [Internet]. 2016;
908 Available from: <http://arxiv.org/abs/1606.09375>
- 909 19. Van Essen DC, Smith SM, Barch DM, Behrens TEJ, Yacoub E, Ugurbil K. The WU-
910 Minn Human Connectome Project: An overview. *NeuroImage.* 2013 Oct 15;80:62–79.
- 911 20. Springenberg JT, Dosovitskiy A, Brox T, Riedmiller M. Striving for Simplicity: The All
912 Convolutional Net. 2014 Dec 21 [cited 2021 Jan 27]; Available from:
913 <https://arxiv.org/abs/1412.6806v3>
- 914 21. Kornblith S, Norouzi M, Lee H, Hinton G. Similarity of Neural Network Representations
915 Revisited. *ArXiv190500414 Cs Q-Bio Stat* [Internet]. 2019 Jul 19 [cited 2021 Jan 15];
916 Available from: <http://arxiv.org/abs/1905.00414>
- 917 22. Nie Q-Y, Müller HJ, Conci M. Hierarchical organization in visual working memory: From
918 global ensemble to individual object structure. *Cognition.* 2017 Feb 1;159:85–96.
- 919 23. Raut RV, Snyder AZ, Raichle ME. Hierarchical dynamics as a macroscopic organizing
920 principle of the human brain. *Proc Natl Acad Sci.* 2020 Aug 25;117(34):20890–7.
- 921 24. Bressler SL, Menon V. Large-scale brain networks in cognition: emerging methods and
922 principles. *Trends Cogn Sci.* 2010 Jun 1;14(6):277–90.
- 923 25. Maaten L van der, Hinton G. Visualizing Data using t-SNE. *J Mach Learn Res.*
924 2008;9(86):2579–605.
- 925 26. Saggat M, Sporns O, Gonzalez-Castillo J, Bandettini PA, Carlsson G, Glover G, et al.
926 Towards a new approach to reveal dynamical organization of the brain using
927 topological data analysis. *Nat Commun.* 2018 Dec;9(1):1399.
- 928 27. Glasser MF, Coalson TS, Robinson EC, Hacker CD, Harwell J, Yacoub E, et al. A
929 multi-modal parcellation of human cerebral cortex. *Nature.* 2016 Aug;536(7615):171–8.
- 930 28. Haxby JV, Guntupalli JS, Connolly AC, Halchenko YO, Conroy BR, Gobbini MI, et al. A
931 Common, High-Dimensional Model of the Representational Space in Human Ventral
932 Temporal Cortex. *Neuron.* 2011 Oct;72(2):404–16.
- 933 29. Nishimoto S, Vu AT, Naselaris T, Benjamini Y, Yu B, Gallant JL. Reconstructing visual
934 experiences from brain activity evoked by natural movies. *Curr Biol.* 2011 Oct
935 11;21(19):1641–6.
- 936 30. Mitchell TM, Shinkareva SV, Carlson A, Chang K-M, Malave VL, Mason RA, et al.
937 Predicting Human Brain Activity Associated with the Meanings of Nouns. *Science.* 2008
938 May 30;320(5880):1191–5.

- 939 31. Huth AG, Nishimoto S, Vu AT, Gallant JL. A continuous semantic space describes the
940 representation of thousands of object and action categories across the human brain.
941 *Neuron*. 2012 Dec 20;76(6):1210–24.
- 942 32. Varoquaux G, Schwartz Y, Poldrack RA, Gauthier B, Bzdok D, Poline J-B, et al. Atlases
943 of cognition with large-scale human brain mapping. Diedrichsen J, editor. *PLOS*
944 *Comput Biol*. 2018 Nov 29;14(11):e1006565.
- 945 33. Yoshimura N, Tsuda H, Kawase T, Kambara H, Koike Y. Decoding finger movement in
946 humans using synergy of EEG cortical current signals. *Sci Rep*. 2017 Sep
947 12;7(1):11382.
- 948 34. Poldrack RA. Inferring Mental States from Neuroimaging Data: From Reverse Inference
949 to Large-Scale Decoding. *Neuron*. 2011 Dec;72(5):692–7.
- 950 35. Cole MW, Ito T, Bassett DS, Schultz DH. Activity flow over resting-state networks
951 shapes cognitive task activations. *Nat Neurosci*. 2016;12.
- 952 36. Penfield W, Boldrey E. SOMATIC MOTOR AND SENSORY REPRESENTATION IN
953 THE CEREBRAL CORTEX OF MAN AS STUDIED BY ELECTRICAL STIMULATION1.
954 *Brain*. 1937 Dec 1;60(4):389–443.
- 955 37. Golarai G, Ghahremani DG, Whitfield-Gabrieli S, Reiss A, Eberhardt JL, Gabrieli JD, et
956 al. Differential development of high-level visual cortex correlates with category-specific
957 recognition memory. *Nat Neurosci*. 2007 Apr;10(4):512–22.
- 958 38. Barch DM, Burgess GC, Harms MP, Petersen SE, Schlaggar BL, Corbetta M, et al.
959 Function in the human connectome: Task-fMRI and individual differences in behavior.
960 *NeuroImage*. 2013 Oct 15;80:169–89.
- 961 39. Bijsterbosch J, Harrison SJ, Jbabdi S, Woolrich M, Beckmann C, Smith S, et al.
962 Challenges and future directions for representations of functional brain organization.
963 *Nat Neurosci*. 2020 Dec;23(12):1484–95.
- 964 40. Glasser MF, Sotiropoulos SN, Wilson JA, Coalson TS, Fischl B, Andersson JL, et al.
965 The minimal preprocessing pipelines for the Human Connectome Project. *NeuroImage*.
966 2013 Oct 15;80:105–24.
- 967 41. Bullmore E, Sporns O. Complex brain networks: graph theoretical analysis of structural
968 and functional systems. *Nat Rev Neurosci*. 2009 Mar;10(3):186–98.
- 969 42. Newman MEJ. Modularity and community structure in networks. *Proc Natl Acad Sci*.
970 2006 Jun 6;103(23):8577–82.
- 971 43. Fan L, Li H, Zhuo J, Zhang Y, Wang J, Chen L, et al. The Human Brainnetome Atlas: A
972 New Brain Atlas Based on Connectional Architecture. *Cereb Cortex*. 2016 Jan
973 8;26(8):3508–26.
- 974 44. McInnes L, Healy J, Saul N, Großberger L. UMAP: Uniform Manifold Approximation
975 and Projection. *J Open Source Softw*. 2018 Sep 2;3(29):861.
- 976 45. Moon KR, van Dijk D, Wang Z, Gigante S, Burkhardt DB, Chen WS, et al. Visualizing
977 structure and transitions in high-dimensional biological data. *Nat Biotechnol*. 2019
978 Dec;37(12):1482–92.

979

980

Supplementary Materials for

Representational learning on high-order human cognition using deep graph convolutions

Authors: Yu Zhang^{1,2,3,*}, Nicolas Farrugia⁴, Alain Dagher⁵ and Pierre Bellec^{2,3,*}

* Corresponding Author:

Yu Zhang, Research Center for Healthcare Data Science, Zhejiang Lab. yuzhang2bic@gmail.com

Pierre Bellec, Département de Psychologie, Université de Montréal. pierre.bellec@gmail.com

This PDF file includes:

Eight Supplementary Figures: **Figure 1-S1** to **Figure 7-S1**

Four Supplementary Tables: **Table S1** to **Table S4**

Two Supplementary Videos: **Figure 6-S3** and **Figure 6-S4**

Supplementary Figures

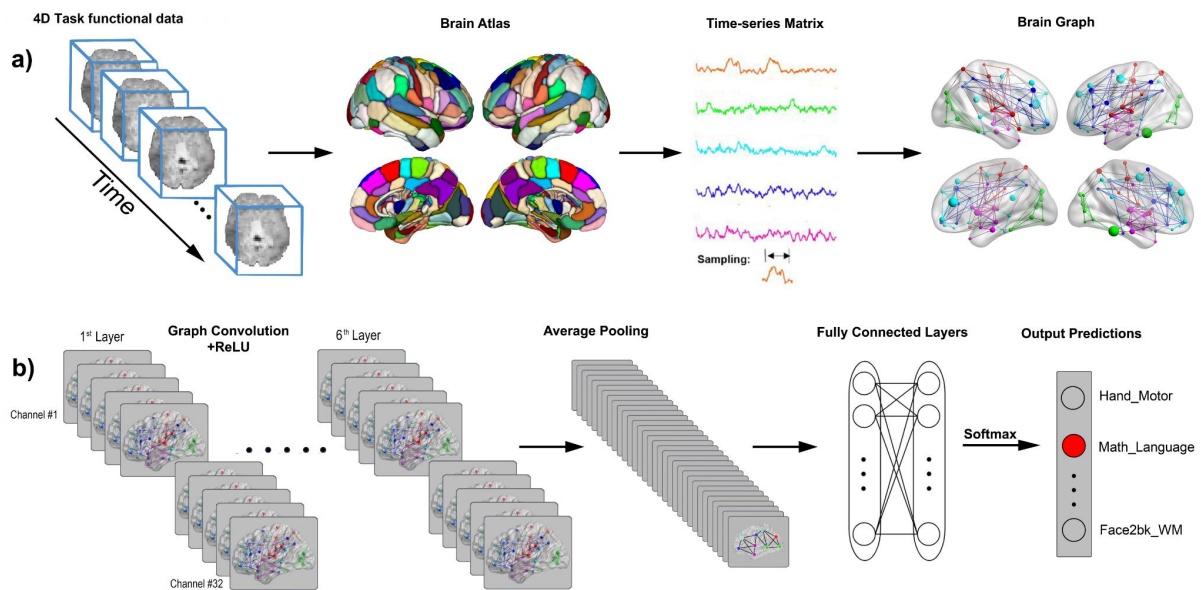


Figure 1-S1. Pipeline of brain decoding using graph convolution network.

The decoding model consists of six ChebNet graph convolutional layers with 32 graph filters at each layer, followed by a flatten layer and 2 fully connected layers. Specifically, for a short series of fMRI volumes, the measured brain activity was first mapped onto a predefined brain atlas consisting of hundreds of brain regions (e.g. 246 regions from Brainnetome atlas (43)). A functional graph was then constructed by calculating group-averaged resting-state functional connectivity for each pair of brain regions. Next, a new representation of task-evoked neural activity was generated through a multi-layer graph convolutional network, taking into account the segregation of localized brain activity and information integration among multiple brain networks. These representations were then used to predict the corresponding cognitive state associated with the short time window. The implementation of the ChebNet graph convolution was based on PyTorch 1.1.0, and was made publicly available in the following repository: https://github.com/zhangyu2ustc/gcn_tutorial_test.git.

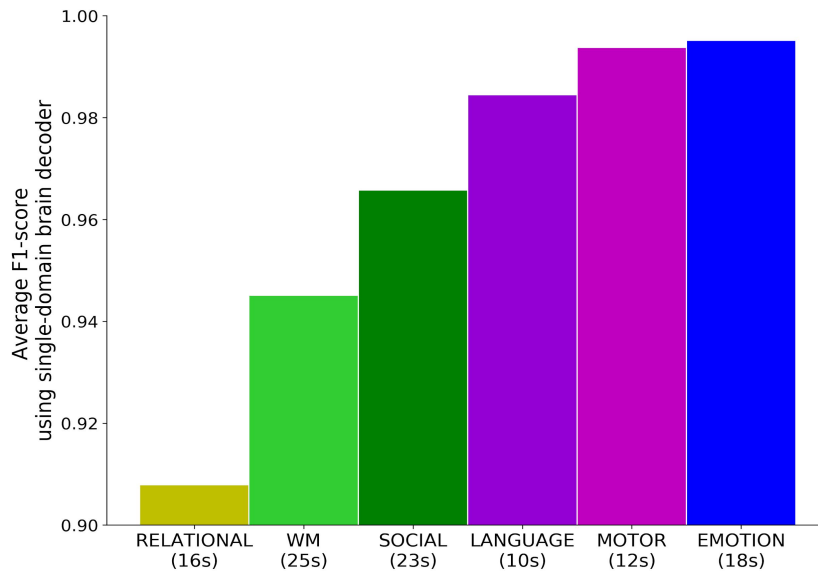


Figure 1-S2. Decoding accuracy of single-domain brain decoders for each of the six cognitive domains.

The same decoding pipeline was used as in Figure 1 except that the decoding model here was trained by using task-fMRI data exclusively from a single domain. Besides, variable temporal durations were used for each cognitive domain, according to the maximum length of event trials/blocks among all experimental tasks, for instance 12s for MOTOR tasks and 25s for WM tasks. Among the six cognitive domains, the emotion tasks (in blue, fearful face vs shape) and motor tasks (in magenta, distinguishing five types of body movements) were the most easily recognizable conditions, with F1-score reaching around 99%.

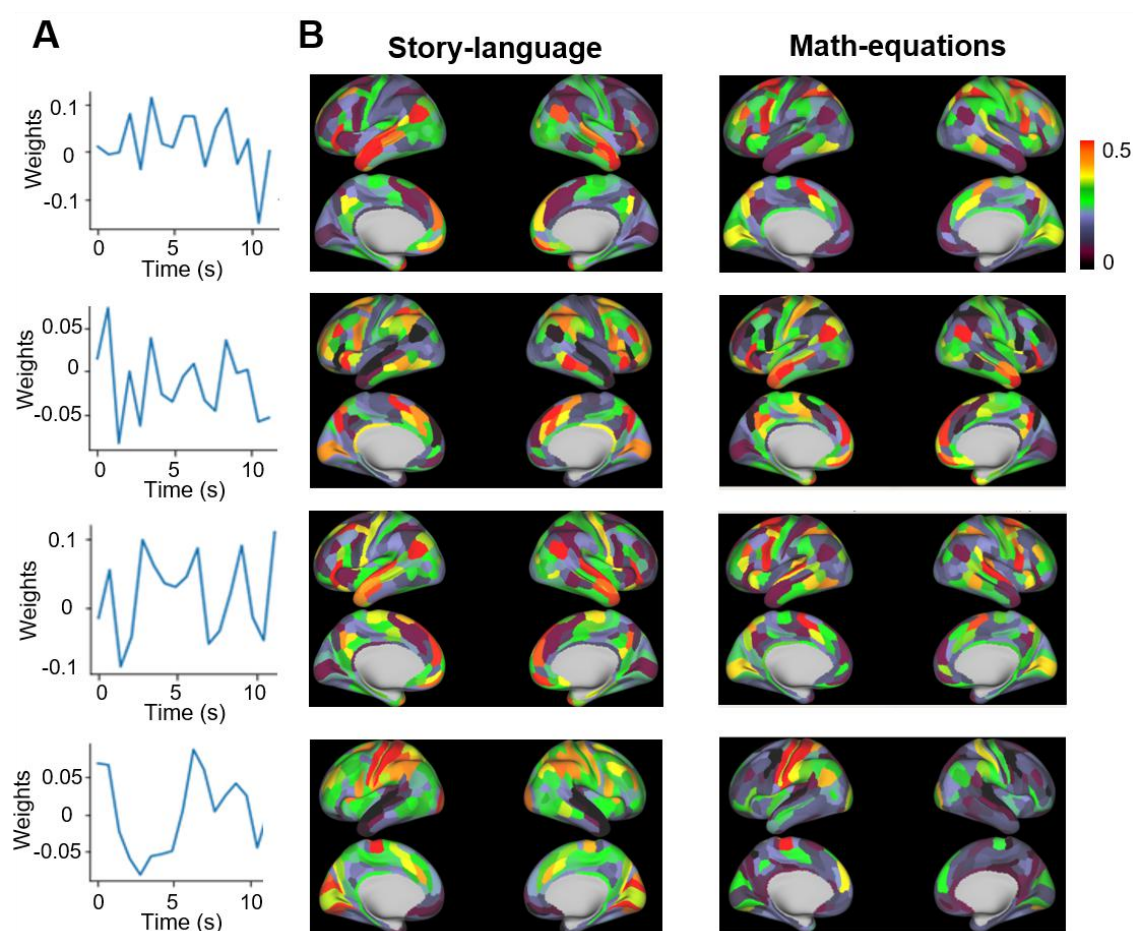


Figure 3-S1. Temporal filters and the corresponding activation maps extracted from the first ChebNet layer for the Language tasks.

We used the Language tasks as an example to illustrate that ChebNet graph convolutions captured hemodynamic response in the temporal domain and brain activations in the spatial domain. First of all, various temporal convolutional kernels were learned at the first ChebNet layer (1st column), which resembled the hemodynamic response function in BOLD signals. Second, using these temporal filters, the corresponding “activation maps” in the first ChebNet layer were extracted and grouped according to the task conditions, e.g. story (2nd column) and math (3rd column) for the Language tasks.

These activation maps demonstrated a possible explanation of the biological basis behind spatiotemporal graph convolutions. For instance, when the temporal filter only focused on brain activity within 0-6s after task onset (3rd row), the classical language network and frontoparietal network were detected for the story and math condition respectively, responsible for language comprehension and numerical processing during the cognitive task. By contrast, when the temporal filter focused on brain activity 5-10s after task onset (4th row), the motor and sensory cortex were

detected for both conditions, responsible for button press during the task. When both time windows were taken into account (1st row, 0-10s after task onset), the brain response in the language network and frontoparietal network were further enhanced while the activity in the auditory cortex was weakened. Similar analysis on the *Motor* tasks was shown in Figure 1-Supplement 2 in (17).

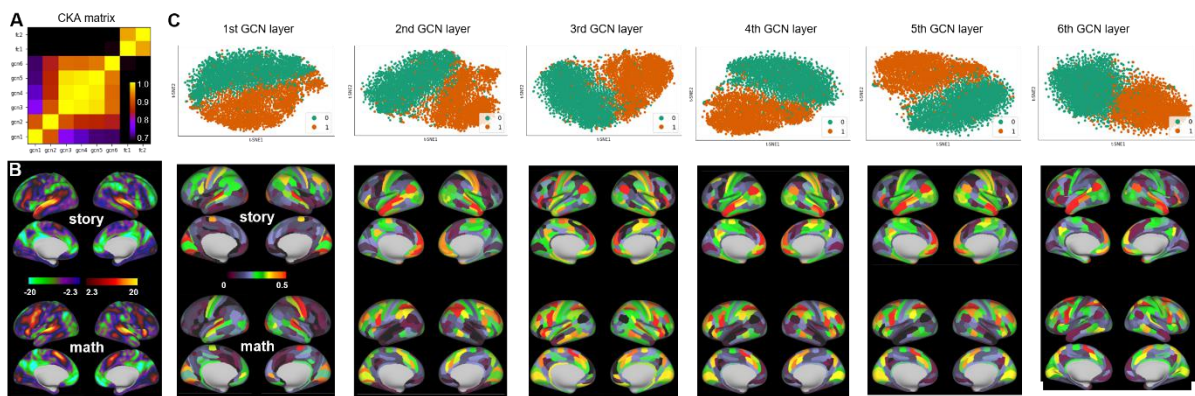


Figure 3-S2. Layer representations in ChebNet resembled brain activation maps.

We used the Language tasks as an example to illustrate the hierarchical organization of layer representations in ChebNet. At each ChebNet layer, the layer activations were extracted and saved as feature representations for the following analysis. First, the similarity of representations was calculated using CKA with a linear kernel (A), which illustrated a hierarchical organization of layer representations in ChebNet. These representations were then projected onto a 2-dimensional space using t-SNE (1st row in C) which indicated a nice disassociation between different task conditions (e.g. story vs math for the Language task). After that, the “activation map” associated with each task condition (2nd row in C) were calculated by averaging the layer representations across all data samples of the same category and mapped back onto the cortical surface. These representations resembled the actual brain activation maps detected by the canonical GLM approach (B), provided by (38), downloaded from neurovault (<https://neurovault.org/collections/457/>).

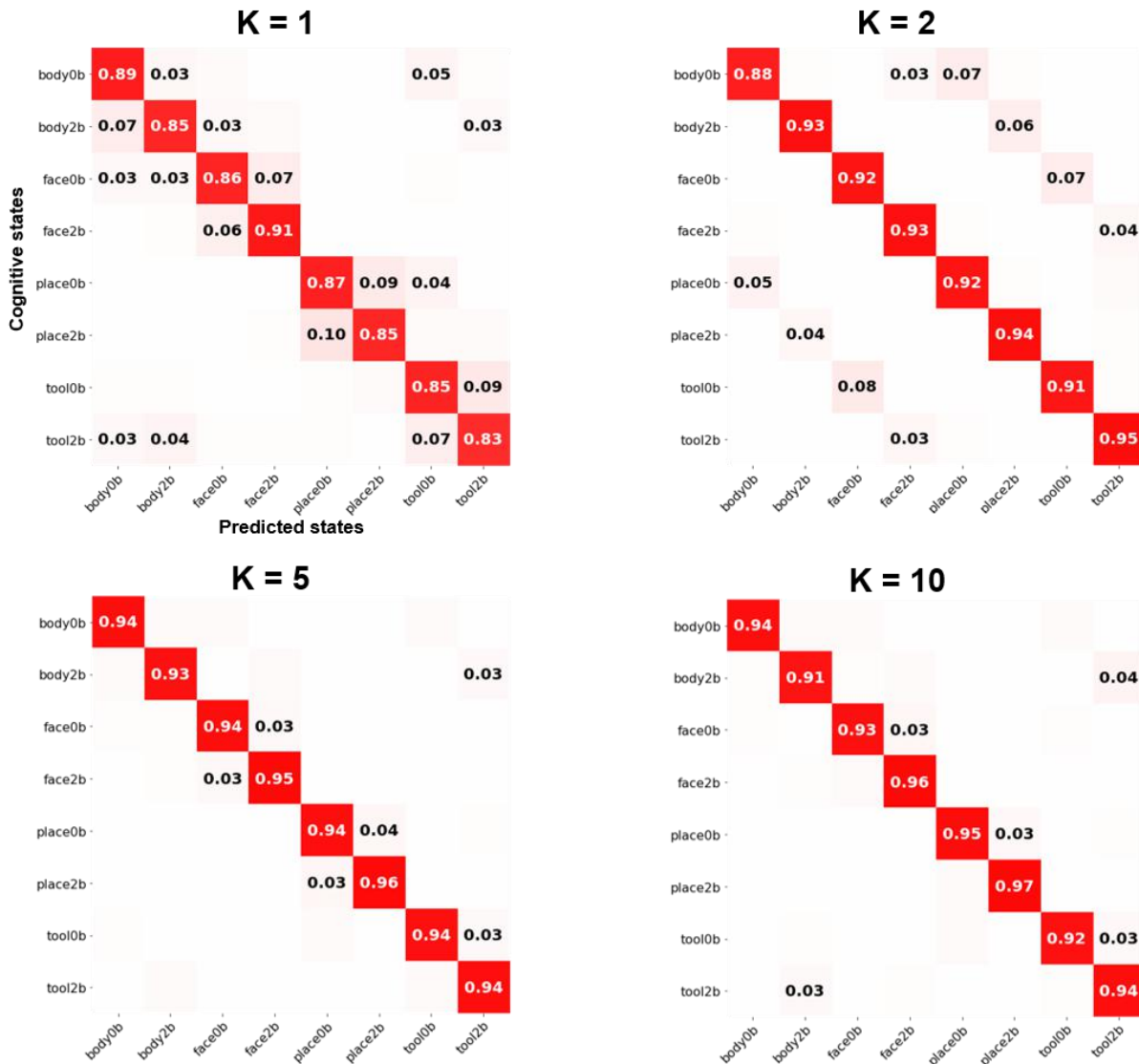


Figure 4-S1. Confusion matrix of Working-memory tasks using ChebNet graph convolution with different K-orders.

The confusion matrix was normalized by each task condition (row) such that each element in the matrix shows the recall score, i.e. among all predictions how many of them are positive predictions. The confusion matrix showed a nice block diagonal architecture, indicating that the majority of the cognitive tasks were accurately identified in all models. ALL decoding models were trained using 25s of fMRI time series. The classification errors were largely reduced when using high-order graph convolutions, e.g. K=1 vs K=5, especially for 0back vs 2back tasks. Our results indicated that large-scale functional integration of brain dynamics played an important role in the decoding of working memory tasks.

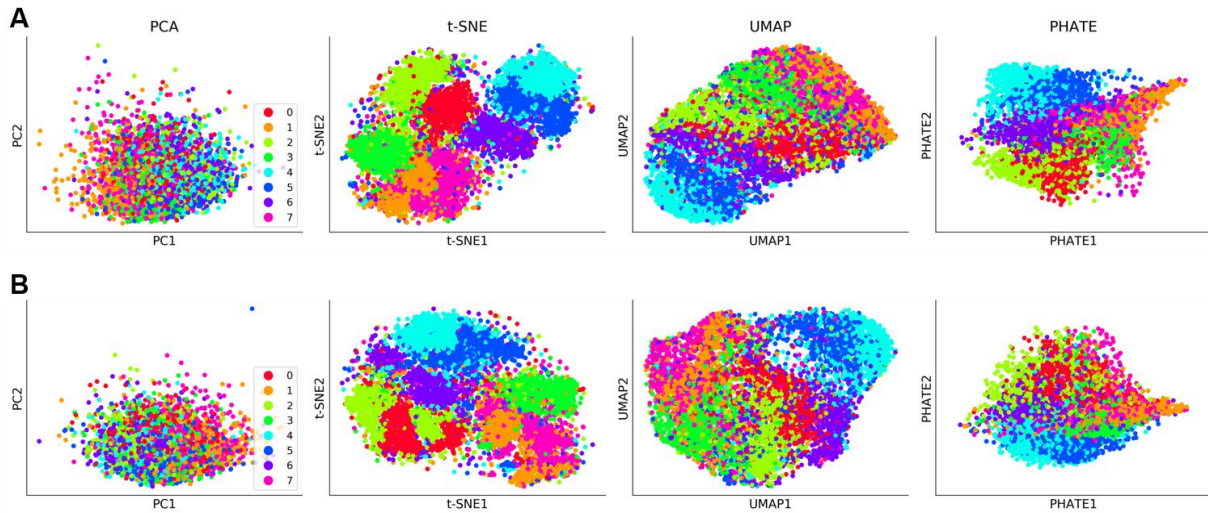


Figure 6-S1. Projection of features representations from the ChebNet decoding model for WM tasks.

Graph representations were mapped onto a 2-dimensional space by using different dimension reduction techniques, including PCA (first column), t-SNE (second column) (25), UMAP (third column) (44), and PHATE (last column)(45). The data samples included eight task conditions from Working-memory tasks, i.e. 0-back and 2-back on images of body parts (class 0 and 1), 0-back and 2-back on face images (class 2 and 3), 0-back and 2-back on place images (class 4 and 5), 0-back and 2-back on images of tools (class 6 and 7). Best visualization was provided by t-SNE. Two different decoding models were evaluated, including ChebNet-K5 (A) and ChebNet-K1 (B). The ChebNet-K5 model showed higher distinctions among eight WM task conditions.

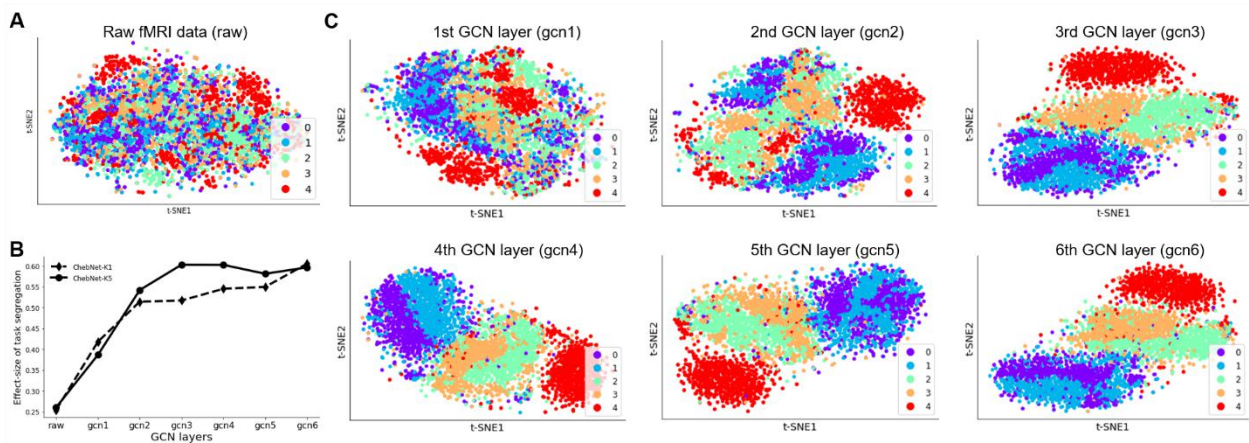


Figure 6-S2. Layer representations of the ChebNet-K5 on the Motor task.

Both raw fMRI data and layer representations of each ChebNet layer were projected onto 2-dimensional space by using t-SNE. No clear structure of task conditions was observed in the raw fMRI data (A), which slightly improved in the low-level representations, e.g. 1st and 2nd ChebNet layer (C). Since the 4th ChebNet layer, the tongue movement (in red) was easily distinguished from other motor tasks in the representations but still showed a mixture effect between left and right movements. In the last ChebNet layer, the representations for the five types of body movements were highly clustered and easily separated from each other (task segregation index $Q = 0.60$). Similar level of task segregation in the last ChebNet layer when using different K-orders in ChebNet, but a faster convergence speed was detected in the ChebNet-K5 model (B). The Motor task data includes five types of body movements, i.e. the movement of right foot (class 0, in purple), left foot (class 1, in green), right hand (class 2, in cyan), left hand (class 3, in orange), and tongue (class 4, in red).

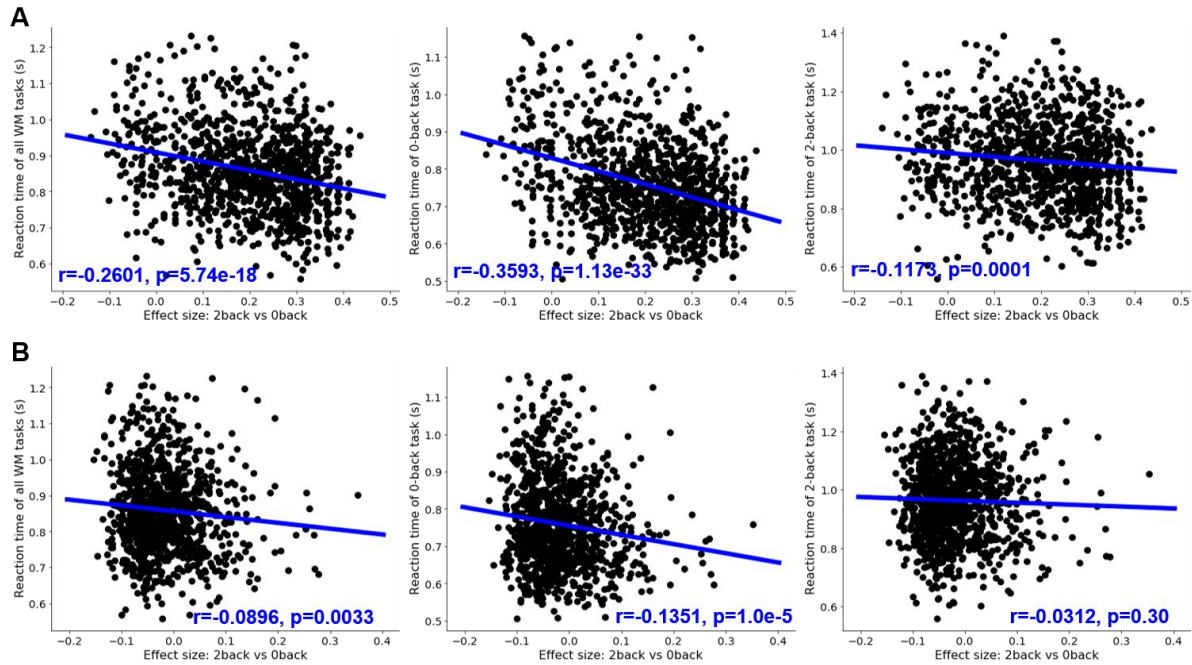


Figure 7-S1. Modularity scores in the state-transition graph significantly correlated with the reaction time of Working-memory tasks.

The modularity score was calculated based on the state-transition graph of each subject, as proposed by (26). We found significant correlations between the modularity scores of graph representations (A) and average reaction time during working-memory task (1st panel), 0back task conditions (2nd panel) and 2back task conditions (3rd panel). Much weaker associations were detected in the raw fMRI data (B). The blue lines indicated the linear regression models between the modularity score of the state-transition graph and the reaction time during task performance. The analysis was done among all subjects from HCP S1200 release, with complete records of behavioral and imaging data for working memory tasks (N=1074).

Supplementary Tables

Table S1. Scanning parameters and experimental designs of HCP task-fMRI dataset.

The entire dataset includes in total of 14,895 functional runs across the six cognitive domains, and resulted in 138,662 data samples of fMRI signals when using a 10s time window (i.e. 15 functional volumes at TR=0.72s)

| Task Domains | #Subjects | #Runs | #Volumes per run | #Trials per run | #Conditions | Minimal duration per block (sec) |
|-----------------------|------------------|--------------|-------------------------|------------------------|--------------------|---|
| Working memory | 1085 | 2 | 405 | 8 | 8 | 25 |
| Motor | 1083 | 2 | 284 | 10 | 5 | 12 |
| Language | 1051 | 2 | 316 | 8 | 2 | 10 |
| Social Cognition | 1051 | 2 | 274 | 5 | 2 | 23 |
| Relational processing | 1043 | 2 | 232 | 6 | 2 | 16 |
| Emotion | 1047 | 2 | 176 | 6 | 2 | 18 |

Table S2. Decoding accuracy from the single-domain decoders.

Six single-domain decoders were trained by using fMRI responses from each cognitive domain exclusively and to predict the cognitive states on a trial basis. Different temporal durations were used, according to the maximum length of event trials on the target cognitive domain, for instance 12s for MOTOR tasks and 25s for WM tasks.

| Task Domains | #Subjects | #Samples (single trials) | #Condition n | Time windows (sec) | Decoding accuracy (F1-score) |
|-----------------------|------------------|---------------------------------|---------------------|---------------------------|-------------------------------------|
| Working memory | 1085 | 17,360 | 8 | 25 | 0.9451 |
| Motor | 1083 | 21,660 | 5 | 12 | 0.9938 |
| Language | 1051 | 16,816 | 2 | 10 | 0.9845 |
| Social Cognition | 1051 | 10,510 | 2 | 23 | 0.9658 |
| Relational processing | 1043 | 12,516 | 2 | 16 | 0.9079 |
| Emotion | 1047 | 12,564 | 2 | 18 | 0.9952 |

Table S3: Comparison of decoding performance between different models.

We reported the best performance for the baseline models after a grid search of the hyperparameters. For SVC approaches, we used the one-vs-rest ('ovr') decision function to handle multi-classes and reported the highest accuracy after the grid search for the hyper-parameter ($C = [0.0001, 0.001, 0.1, 1, 10, 100]$). For Random Forest, we reported the highest accuracy after evaluating different settings of the classifier including depth of trees: $[4, 16, 64, 256, 1024]$ and number of trees: $[100, 2000]$. For MLP (multilayer perceptron), GCN (using first-order graph convolution, (17)) and ChebNet (using 5-order graph convolution), we reported the mean and standard deviation of the decoding accuracies among 10 fold cross-validation with shuffle splits. All models were evaluated on the task of decoding 21 task states by using 10s of fMRI signals (in total of 138,662 data samples).

| Models | Train Accuracy | Validation Accuracy | Test Accuracy |
|----------------|-----------------------|----------------------------|------------------------|
| SVC-linear | 67.2% | 63.3% | 64.1% |
| SVC-rbf | 99.7% | 73.5% | 73.8% |
| Random Forest | 100% | 48.0% | 47.5% |
| MLP(256-64) | 87.9%(+/-1.83%) | 83.2%(+/-3.28%) | 76.1%(+/-0.41%) |
| GCN | 96.3%(+/-0.42%) | 90.2%(+/-0.21%) | 90.7%(+/-0.20%) |
| ChebNet | 90.91%(+/-0.18%) | 92.73(+/-0.12%) | 93.43(+/-0.44%) |

Table S4: Heritability analysis of brain responses and behavioral scores.

Heritability estimates were conducted for both representations of brain responses and behavioral performance in-scanner, associated with WM tasks, after controlling for confounding effects of age, gender, handedness and head motion. The average accuracy (Acc) and reaction time (RT) showed high heritability estimates of additive genetic effects. For graph representations and raw fMRI signals, the high-dimensional data was first projected onto a 2-dimensional space using t-SNE and then the task segregation effect was estimated based on individual state-transition graph (see Method section). Significant heritability estimates were also detected in the ChebNet graph representations but not in raw fMRI signals.

| Traits | h^2 | SE | p-value | FDR corrected p-value | Covariance Explained |
|--------------------------------------|---------------|--------|-----------------|-----------------------|----------------------|
| ChebNet graph representations | 0.2882 | 0.0588 | 3.00E-07 | 3.43E-07 | 0.0017 |
| Raw fMRI signals | 0.0008 | 0.0545 | 0.4935 | 0.4935 | 0.0029 |
| WM_Task_Acc | 0.5624 | 0.0435 | 8.56E-27 | 3.42E-26 | 0.0434 |
| WM_Task_2bk_Acc | 0.5887 | 0.0425 | 6.97E-29 | 5.58E-28 | 0.0446 |
| WM_Task_0bk_Acc | 0.3215 | 0.0564 | 6.21E-09 | 8.29E-09 | 0.0222 |
| WM_Task_RT | 0.4118 | 0.0560 | 4.01E-13 | 8.01E-13 | 0.0105 |
| WM_Task_2bk_RT | 0.4534 | 0.0556 | 5.40E-15 | 1.44E-14 | 0.0094 |
| WM_Task_0bk_RT | 0.3294 | 0.0583 | 6.06E-09 | 8.29E-09 | 0.0085 |