

BrainMoE: Cognition Joint Embedding via Mixture-of-Expert Towards Robust Brain Foundation Model

Ziquan Wei, Tingting Dan, Tianlong Chen, Guorong Wu*

UNC Chapel Hill, USA

NeurIPS 2025



Motivation: Why we need a BrainMoE

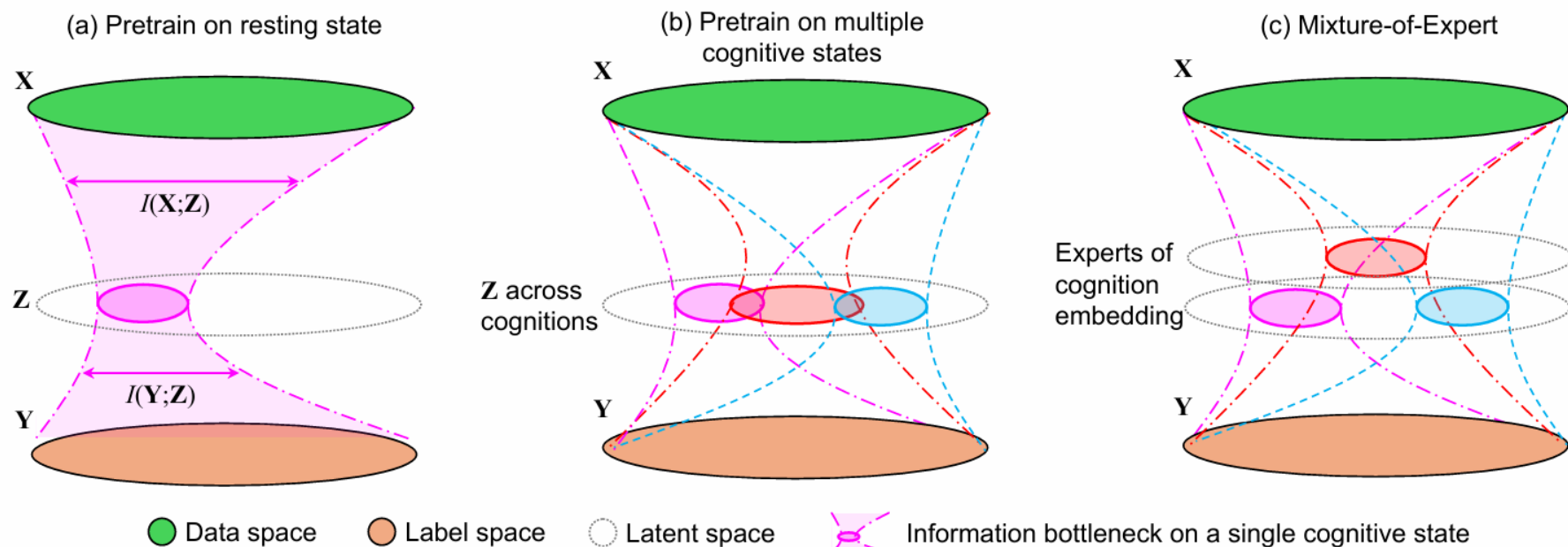


Figure 1: Motivation of BrainMoE through the lens of the information bottleneck theory. (a) Feature representation learning makes an information bottleneck between data and label space, where \mathbf{X} denotes data, \mathbf{Z} denotes latent feature representation, \mathbf{Y} is the label of data, and $I(\cdot; \cdot)$ is the mutual information. (b) A model pre-trained on multiple cognitive states may compromise the underlying heterogeneity between different states, where \mathbf{Z} cannot be optimal for all states. (c) The mixture of brain experts dedicated to diverse cognitive states leads to a joint cognition embedding so that the downstream applications can be advanced by stratified pre-training on rich behavioral tasks.

More data != more performance

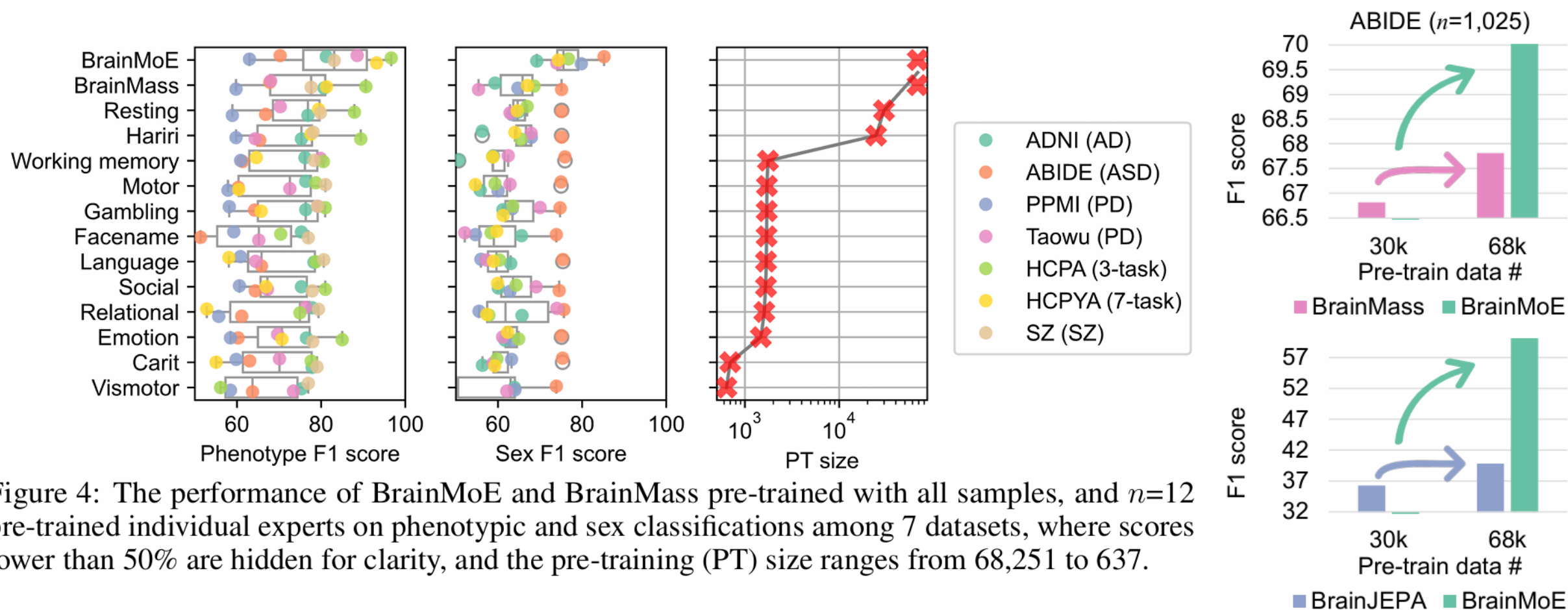
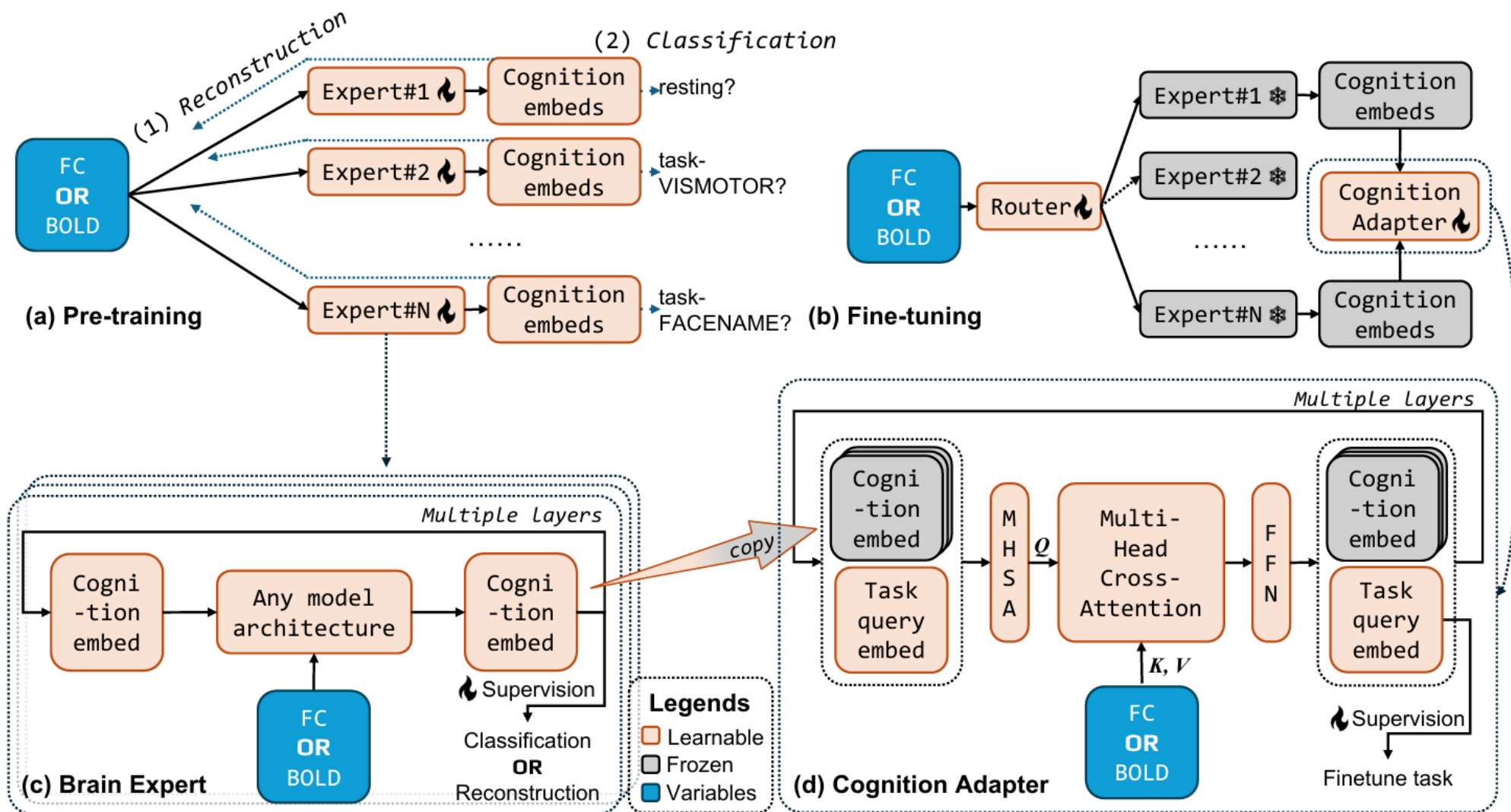


Figure 4: The performance of BrainMoE and BrainMass pre-trained with all samples, and $n=12$ pre-trained individual experts on phenotypic and sex classifications among 7 datasets, where scores lower than 50% are hidden for clarity, and the pre-training (PT) size ranges from 68,251 to 637.

Framework of BrainMoE



Accuracy scores: Phenotypic prediction

Table 1: MoE improvement on phenotypic classification F1 score compared to the baseline, where 30k is pre-trained on resting-state data ($n=29,951$), and 68k is pre-trained on all data ($n=68,251$). PT stands for pre-training. Colored text indicates the performance increase/decrease from using 68k.

Predictor	BrainMass					BrainJEPA		
	SVM	SVM	MLP	MLP	BrainMoE	ViT	ViT	BrainMoE
PT #	30k	68k	30k	68k	68k	30k	68k	68k
ADNI	75.32 \pm 7.06	75.32 \pm 7.06	76.86 \pm 7.26	80.70 \pm 7.85	81.23 \pm 11.00	74.16 \pm 8.55	74.16 \pm 8.55	77.11 \pm 6.64
↳AD		0.00		3.84 ↑	4.37 ↑		0.00	2.95 ↑
ABIDE	62.31 \pm 1.95	64.12 \pm 2.31	66.81 \pm 4.18	67.81 \pm 3.91	70.26 \pm 3.40	36.25 \pm 6.93	39.82 \pm 3.91	54.55 \pm 9.89
↳ASD		1.81 ↑		1.00 ↑	3.45 ↑		3.77 ↑	18.30 ↑
PPMI	54.87 \pm 15.76	56.52 \pm 14.86	58.90 \pm 14.29	59.77 \pm 14.22	62.97 \pm 13.94	38.69 \pm 13.91	38.69 \pm 13.91	60.49 \pm 11.59
↳PD (staged)		1.65 ↑		0.87 ↑	4.07 ↑		0.00	21.80 ↑
Taowu	58.33 \pm 34.78	65.67 \pm 20.55	70.29 \pm 17.97	68.00 \pm 21.46	88.57 \pm 12.51	36.08 \pm 26.38	36.94 \pm 20.98	79.86 \pm 14.46
↳PD (binary)		7.34 ↑		2.29 ↓	18.28 ↑		0.86 ↑	43.78 ↑
SZ	76.95 \pm 9.01	76.95 \pm 9.01	79.85 \pm 8.69	77.63 \pm 8.56	83.10 \pm 11.33	76.98 \pm 9.00	78.97 \pm 9.79	82.86 \pm 9.19
↳Schizophrenia		0.00		2.22 ↓	3.25 ↑		1.99 ↑	5.88 ↑
HCPA	85.16 \pm 0.41	89.73 \pm 0.58	87.91 \pm 0.48	90.63 \pm 0.74	96.67 \pm 0.77	59.54 \pm 15.47	53.12 \pm 14.19	81.74 \pm 0.51
↳3-task,rest		4.57 ↑		2.72 ↑	8.76 ↑		6.42 ↓	22.20 ↑
HCPYA	77.51 \pm 2.42	80.87 \pm 1.77	79.40 \pm 1.78	81.27 \pm 1.27	93.19 \pm 0.72	50.68 \pm 25.20	56.10 \pm 29.16	74.59 \pm 3.79
↳7-task		3.36 ↑		1.87 ↑	13.79 ↑		5.42 ↑	23.91 ↑

Accuracy scores: Sex classification

Table 2: MoE improvement on sex classification F1 score compared to the baseline, where the sample size of the downstream dataset is indicated. PT stands for pre-training.

Predictor PT #	BrainMass					BrainJEPA		
	SVM 30k	SVM 68k	MLP 30k	MLP 68k	BrainMoE 68k	ViT 30k	ViT 68k	BrainMoE 68k
ADNI ↳n=138	48.60 \pm 6.55	54.30 \pm 12.48 5.70 ↑	64.82 \pm 4.30	59.30 \pm 13.05 5.52 ↓	69.22 \pm 5.26 4.40 ↑	37.70 \pm 8.73	36.42 \pm 8.22 1.28 ↓	62.98 \pm 7.76 25.28 ↑
ABIDE ↳n=1025	73.84 \pm 3.49	73.84 \pm 3.49 0.00	75.12 \pm 5.27	75.12 \pm 6.32 0.00	85.21 \pm 3.77 10.09 ↑	78.08 \pm 5.84	78.08 \pm 5.84 0.00	78.08 \pm 6.15 0.00
PPMI ↳n=209	52.03 \pm 14.16	56.58 \pm 12.28 4.55 ↑	63.32 \pm 14.96	64.73 \pm 14.12 1.41 ↑	79.81 \pm 8.46 16.49 ↑	46.23 \pm 9.00	46.23 \pm 9.00 0.00	67.57 \pm 7.24 21.34 ↑
Taowu ↳n=40	46.24 \pm 23.96	46.24 \pm 23.96 0.00	62.86 \pm 28.20	55.38 \pm 20.93 7.48 ↓	74.00 \pm 27.79 11.14 ↑	46.24 \pm 23.97	51.67 \pm 21.12 5.43 ↑	90.00 \pm 12.96 43.76 ↑
HCPA ↳n=4863	66.20 \pm 1.58	68.25 \pm 1.70 2.05 ↑	66.93 \pm 0.63	68.58 \pm 0.99 1.65 ↑	76.76 \pm 0.93 9.83 ↑	40.32 \pm 4.07	40.32 \pm 4.07 0.00	44.24 \pm 7.01 3.92 ↑
HCPYA ↳n=3293	63.33 \pm 3.01	65.47 \pm 3.51 2.14 ↑	64.57 \pm 2.53	66.98 \pm 3.30 2.41 ↑	74.36 \pm 4.43 9.79 ↑	40.20 \pm 4.22	40.20 \pm 4.22 0.00	43.24 \pm 8.19 3.04 ↑

Performance interpretation

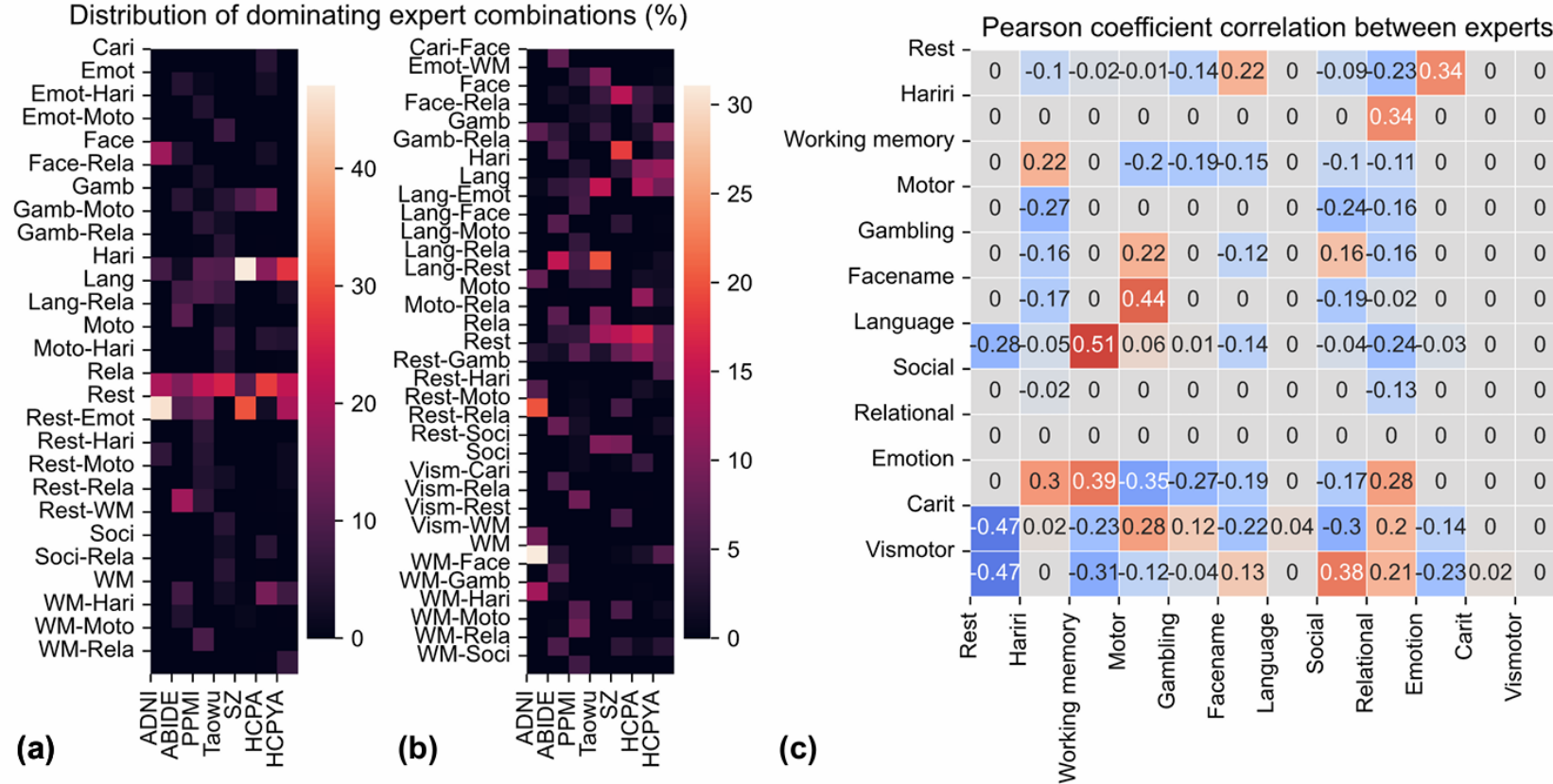


Figure 5: The distribution of dominating expert combinations of **(a)** late fusion MoE and **(b)** BrainMoE shows the router preference, where the first 4 letters of cognitive state are used as abbreviation. **(c)** The correlation between cognition embeddings of experts shows expert diversity.

More applications

Table 5: Age regression performance compared to the baseline, where the sample size of the downstream dataset is indicated, and unit is year.

MSE	ADNI	ABIDE	PPMI	Taowu	HCPA	HCPYA
BrainMass	$36.28_{\pm 19.83}$	$36.77_{\pm 17.2}$	$33.09_{\pm 21.87}$	$38.10_{\pm 28.33}$	$22.66_{\pm 7.51}$	$5.46_{\pm 2.69}$
BrainMoE	$36.27_{\pm 10.23}$	$4.86_{\pm 2.67}$	$29.89_{\pm 8.39}$	$30.72_{\pm 12.25}$	$10.56_{\pm 2.86}$	$3.45_{\pm 1.08}$

Table 6: BrainMoE applies on a multimodal dataset, NATVIEW [27].

NATVIEW	8-task (F1)	Sex (F1)	Age (MSE)
BrainMass (fMRI)	$67.66_{\pm 5.74}$	$63.67_{\pm 5.16}$	$8.05_{\pm 5.58}$
CBraMod (EEG)	$68.71_{\pm 1.46}$	$65.39_{\pm 2.33}$	$8.26_{\pm 5.97}$
BrainMoE (13 Ex.)	$68.73_{\pm 3.72}$	$65.47_{\pm 5.38}$	$7.99_{\pm 5.53}$

Ablation studies

