

Towards a general-purpose foundation model for fMRI analysis

Cheng Wang¹, Yu Jiang¹, Zhihao Peng¹, Chenxin Li¹, Changbae Bang², Lin Zhao³, Jinglei Lv⁴,
Jorge Sepulcre⁵, Carl Yang⁶, Lifang He⁷, Tianming Liu³, Daniel Barron⁸, Quanzheng Li⁹,
Randy Hirschtick⁹, Byung-Hoon Kim², Xiang Li^{9,10*}, Yixuan Yuan^{1*}

¹*The Chinese University of Hong Kong, Hong Kong*

²*Yonsei University, Seoul, South Korea*

³*University of Georgia, Athens, GA*

⁴*University of Sydney, Sydney, Australia*

⁵*Yale School of Medicine, New Haven, CT*

⁶*Emory University, Atlanta, GA*

⁷*Lehigh University, Bethlehem, PA*

⁸*Brigham and Women’s Hospital, Boston, MA*

⁹*Massachusetts General Hospital, Boston, MA*

¹⁰*Kempner Institute for Natural and Artificial Intelligence, Cambridge, MA*

**Corresponding authors*

1 Abstract

Functional Magnetic Resonance Imaging (fMRI) is crucial for studying brain function and diagnosing neurological disorders. However, existing methods for fMRI analysis suffer from reproducibility and transferability challenges, due to complex pre-processing pipelines and task-specific model designs. In this work, we introduce a **Neuroimaging Foundation Model with Spatial-Temporal Optimized and Representation Modeling** (NeuroSTORM) that learns generalizable representations directly from 4D fMRI volumes, and achieves efficient pre-trained knowledge transfer across diverse downstream applications. Specifically, NeuroSTORM is pre-trained on a remarkable 28.65 million fMRI frames (>9,000 hours) collected from over 50,000 subjects, spanning multiple centers and covering ages 5 to 100. NeuroSTORM employs a shifted scanning strategy based on a Mamba backbone, enabling efficient direct processing of 4D fMRI volumes. Moreover, we introduce a spatial-temporal optimized pre-training strategy coupled with a task-specific prompt tuning technique, learning transferable fMRI features for efficient downstream adaptation. Experimental results show that NeuroSTORM consistently outperforms existing methods across five diverse downstream tasks: age and gender prediction, phenotype prediction, disease diagnosis, fMRI(-to-image) retrieval, and task-based fMRI (tfMRI) state classification. To further validate the clinical utility of NeuroSTORM, we evaluate it on two clinical datasets comprising patients with 17 different diagnoses from hospitals in the United States, South Korea, and Australia. NeuroSTORM maintains high relevance in predicting psychological/cognitive phenotypes and achieves the best disease diagnosis performance among all existing methods. In summary, NeuroSTORM offers a standardized, open-source foundation model to enhance reproducibility and transferability in fMRI analysis for clinical applications.

2 Introduction

Functional Magnetic Resonance Imaging (fMRI) is a non-invasive imaging technique widely used to study brain function and structure ¹, and is also utilized in clinical practice for pre-surgical mapping, the assessment and monitoring of various neurological disorders ^{2,3}. Moreover, it provides robust technological support for brain signal decoding and human-computer interaction in brain-computer interface research ⁴. Driven by the advances in deep learning, recent works in fMRI analysis have achieved substantial improvements in the accuracy and efficiency of interpreting complex neural patterns from both resting-state fMRI (rsfMRI) and task fMRI (tfMRI) data ⁵⁻¹². Despite its theoretical and methodological progress, the field of fMRI remains fragmented across data formats, preprocessing pipelines, and analytic models, challenges recently summarized by ¹³. Specifically, analytic reproducibility is fundamental to neuroimaging research, and overlooking it can introduce bias into subsequent analyses ¹⁴. While prior studies related to reproducibility and transferability have focused on identifying reliable biomarkers ¹⁵, standardizing pipelines ¹⁶, and validating test-retest reliability ¹⁷, most of them remain tailored to specific fMRI applications and lack the generalizability to other tasks. Thus, there is an opportunity to develop a model that is intrinsically generalizable across diverse experimental settings and conditions, which can then enable both reproducibility and transferability within a unified framework.

The emergence of foundation models ¹⁸⁻²¹ presents a paradigm-shifting framework by scalable learning across tasks and improved robustness through large-scale pre-training and adaptable architectures. Foundation models were initially developed for natural language processing tasks ^{18,19}, where models such as ChatGPT demonstrated remarkable multitask capabilities by training on web-scale text corpora. This success has inspired analogous developments in the medical domain ²², where foundation models are being applied to address challenges such as anatomical variability and limited annotated data. For example, RETFound ²¹ establishes a foundation for retinal imaging using self-supervised learning on 1.6 million unlabeled images. Similarly, AnyStar ²³ demonstrates that domain-randomized generative models can achieve cross-modal 3D instance segmentation of star-convex anatomical structures, such as the tumor boundary in MRI scans, without relying on annotated training data. Foundation models enhance transferability by balancing generalizable feature extraction with domain-specific pattern recognition, adapting to the unique requirements of each imaging modality. They also improve reproducibility by capturing noise-resilient patterns, often through self-supervised reconstruction or contrastive-based pre-training frameworks. These approaches can reduce sensitivity to acquisition variations and mitigate the variability introduced by preprocessing pipelines ¹⁶, while preserving meaningful neurobiological information.

However, different from other data modalities, developing foundation models for fMRI data remains fundamentally challenging. Firstly, the raw 4D fMRI signal, comprising up to 10^6 voxels per scan, poses severe computational and optimization bottlenecks. Prior approaches typically reduce dimensionality by projecting data onto pre-defined brain atlases or connectomes ⁵⁻¹². However, these operations result in irreversible information loss and impose structural biases that hinder generalizability across populations ^{24,25}. Secondly, due to the high spatiotemporal redundancy in 4D fMRI volumes, we observe that standard Masked Autoencoders (MAEs) which are widely used for foundation model design struggle to learn informative representations, as masked voxels can often be trivially reconstructed from their spatial or temporal neighbors.

In line with recent calls for scalable, generalizable frameworks for fMRI ¹³, we introduce the **Neuroimaging Foundation Model with Spatial-Temporal Optimized and Representation Modeling (NeuroSTORM)**, a general-purpose fMRI foundation model designed to enhance reproducibility and transferability through large-scale pre-training and architectural innovation (Fig. 1(a)). To enhance computational efficiency in 4D fMRI process-

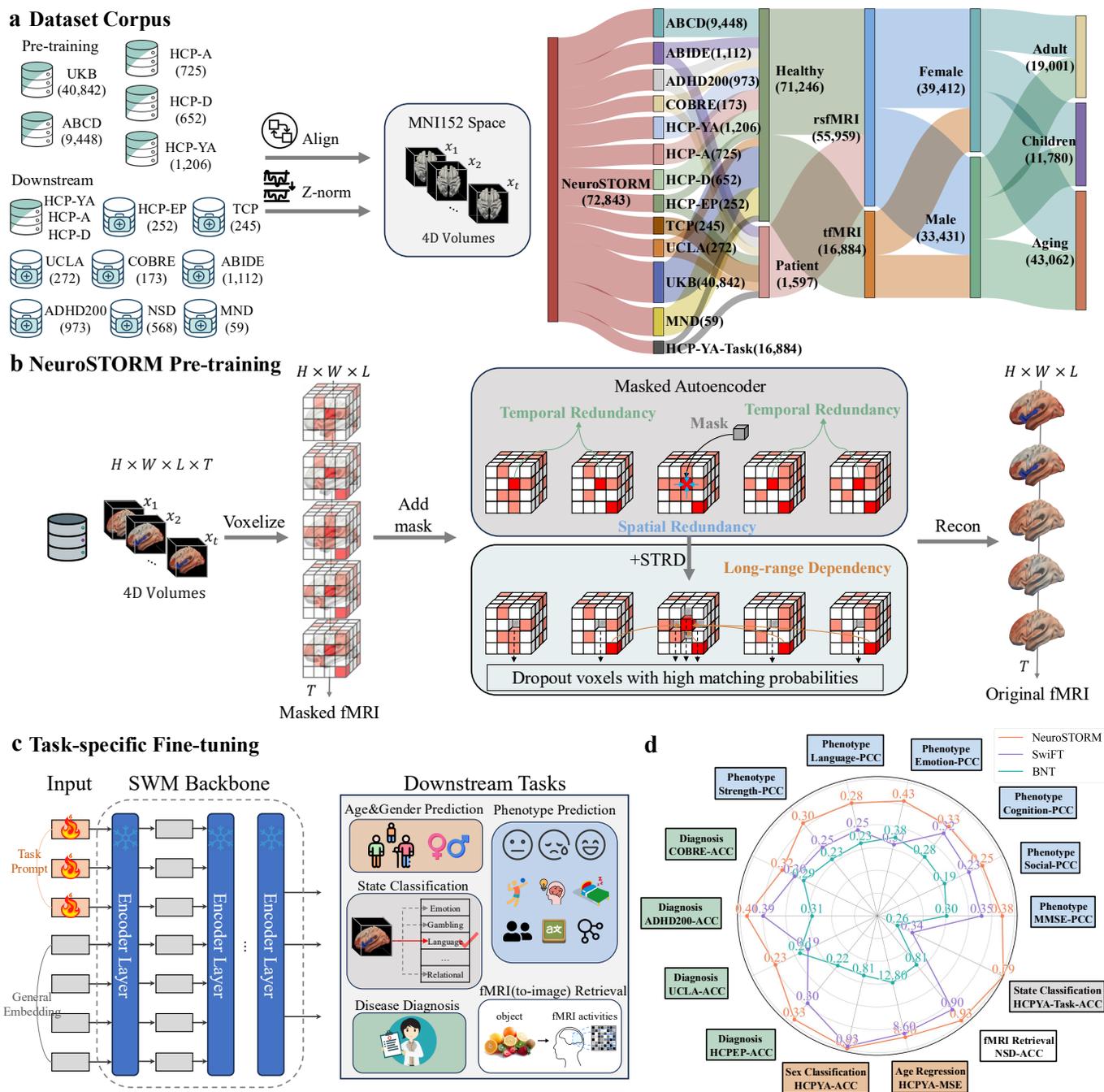


Figure 1: Overview of our proposed NeuroSTORM framework. (a) Data corpus and pre-processing: The model is pre-trained on a collection of publicly available datasets, including over 50,000 rsfMRI and 16,000 tfMRI sequences. All data are aligned to 2mm MNI152 space to create standardized 4D volumes. To facilitate the visualization of datasets with varying scales in the diagram, we use logarithmic values instead of raw numbers. (b) NeuroSTORM pre-training: The model utilizes a masked autoencoder paradigm with a STRD module to enhance the learning of long-range spatiotemporal relationships. (c) Downstream tasks with fine-tuning: For evaluation, NeuroSTORM employs a SWM Backbone and TPT technique for efficient fine-tuning to various downstream tasks. The benchmark includes age and gender prediction, phenotype prediction, disease diagnosis, fMRI-(to-image) retrieval, and tfMRI state classification. (d) Comprehensive performance evaluation: We systematically benchmark NeuroSTORM against previous state-of-the-art models across a diverse set of downstream tasks. The radar chart illustrates NeuroSTORM’s consistent performance improvements.

ing, we introduce a Shifted-Window Mamba (SWM) backbone, which combines linear-time state-space modeling with shifted-window mechanisms to reduce complexity and GPU memory usage. During pre-training, we propose a Spatiotemporal Redundancy Dropout (STRD) module (Fig. 1(b)) for effective learning of inherent characteristics in fMRI data, thereby improving the model’s robustness and reproducibility. For downstream task adaptation, our Task-specific Prompt Tuning (TPT) strategy (Fig. 1(c)) employs a minimal number of trainable, task-specific parameters when fine-tuning NeuroSTORM for new tasks. This provides a simple and integrated approach to applying NeuroSTORM across diverse applications. The corpus used to pre-train NeuroSTORM integrates four large-scale neuroimaging datasets: UK Biobank ²⁶ (40,842 participants), ABCD ²⁷ (9,448 children), and the HCP datasets ²⁸ (HCP-YA, HCP-A, and HCP-D; totaling over 2,500 subjects). This represents the largest multi-source fMRI training dataset assembled to date. Spanning diverse demographics (ages 9–80), clinical conditions, and acquisition protocols, the corpus ensures broad biological and technical variation.

To validate the performance and transferability of NeuroSTORM, we established a comprehensive fMRI analysis benchmark, including five downstream tasks: age and (reported) gender prediction, phenotype prediction, disease diagnosis, fMRI retrieval, and task-based fMRI (tfMRI) state classification. Notably, we assessed the clinical applicability of NeuroSTORM on two clinical datasets from hospitals in both the United States and Australia, Transdiagnostic Connectome Project (TCP) ²⁹ and Motor Neuron Disease (MND) ³⁰. The TCP dataset consists of 245 participants from Brain Imaging Center of Yale University or McLean Hospital in the United States, including both healthy controls and individuals covering a range of psychiatric disorders. The MND dataset includes 36 participants diagnosed with Amyotrophic Lateral Sclerosis and 23 controls, collected at the Royal Brisbane and Women’s Hospital in Australia. NeuroSTORM outperforms or matches state-of-the-art models across all five tasks, demonstrating strong transferability in diverse applications. Moreover, we simulated a data-scarce scenario by limiting the proportion of fine-tuning data in order to evaluate data utilization efficiency. NeuroSTORM demonstrates minor performance degradation on most datasets. Notably, we evaluated NeuroSTORM’s ability to predict clinical phenotypes and perform disease classification in clinical datasets, with experimental results highlighting its clinical value and transferability.

Overall, we present a foundation model for fMRI analysis that achieves outstanding performance across five downstream tasks, enabling for large-scale fMRI studies with enhanced reproducibility and transferability. We have open-sourced a GitHub repository (github.com/CUHK-AIM-Group/NeuroSTORM), which serves as a general-purpose fMRI analysis platform. This repository includes tools for fMRI preprocessing, trainers for both pre-training and fine-tuning, a benchmark suite for diverse fMRI tasks, implementations of NeuroSTORM, as well as other commonly used fMRI analysis models. Detailed guidelines are provided for adding custom preprocessing procedures, pre-training methods, fine-tuning strategies, new downstream tasks, and additional models to the platform. This platform enhances NeuroSTORM’s reproducibility and transferability while advancing fMRI research.

3 Results

The downstream task of age and (reported) gender prediction utilizes HCP-YA ²⁸, HCP-A, HCP-D, UKB ²⁶ and ABCD ²⁷ datasets. Phenotype prediction is based on HCP-YA and TCP ²⁹ dataset. Disease diagnosis draws on multiple datasets, including HCP-EP ²⁸, ABIDE ³¹, ADHD200 ³², COBRE ³³, UCLA ³⁴ and MND ³⁰. The fMRI retrieval task is evaluated on the NSD ³⁵ dataset. Task-based fMRI (tfMRI) state classification task uses the HCP-YA dataset. No annotations are included in the NeuroSTORM pre-training corpus. In all experiments, the downstream task datasets are split into training, validation, and testing sets in an 8:1:1 ratio. Additionally, we evaluate NeuroSTORM’s robustness in data-scarce scenarios by employing varying percentages of fine-

tuning data in the datasets. For fMRI retrieval, we conduct experiments using the LAION-5B³⁶ dataset as the retrieval candidate pool.

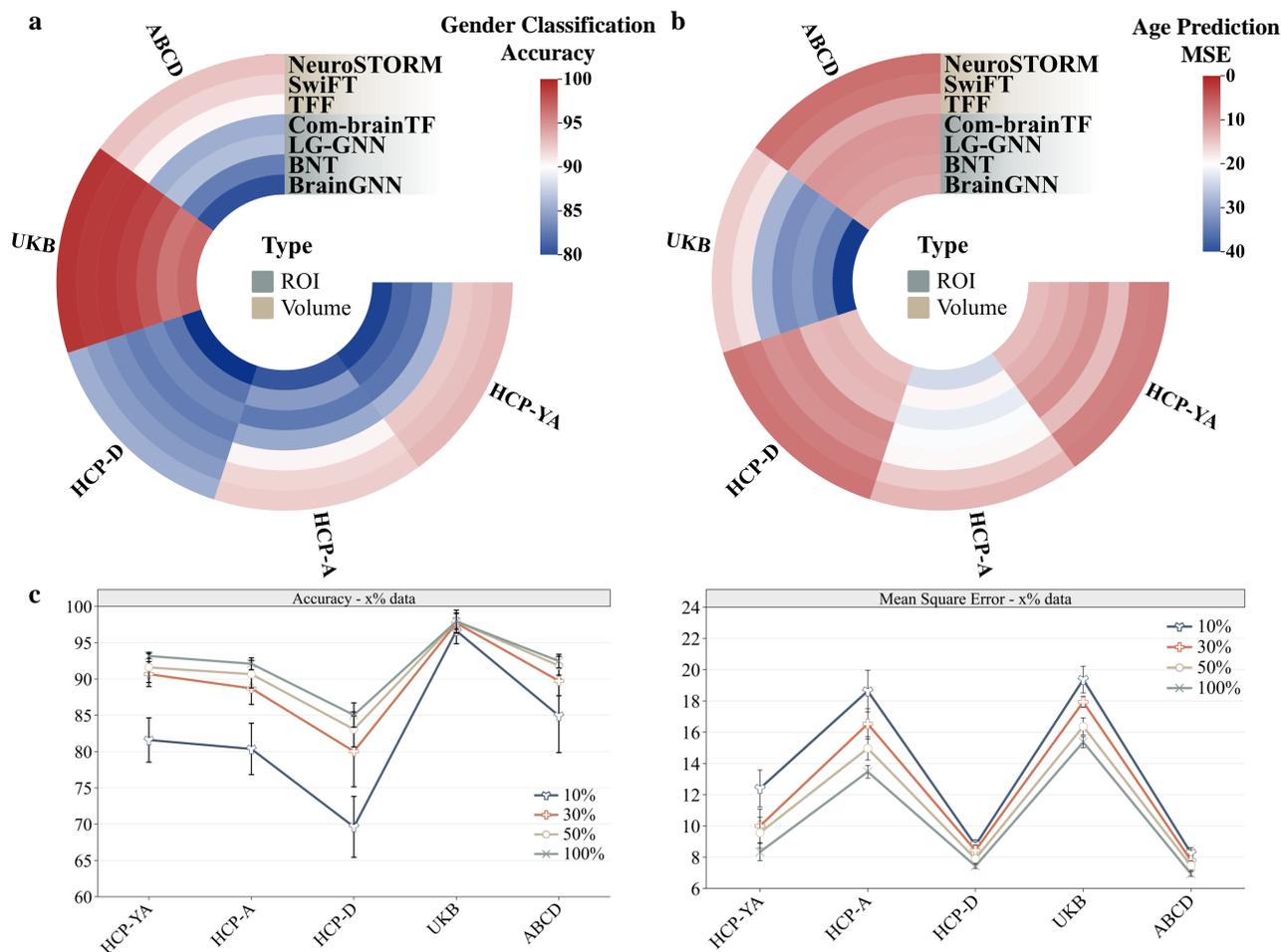


Figure 2: Evaluation of NeuroSTORM’s performance in (reported) gender classification and age regression. (a) Gender classification performance: NeuroSTORM consistently outperformed competing ROI-based and volume-based methods across datasets, including HCP-YA, HCP-A, HCP-D, UKB and ABCD. It achieved the highest accuracy and area under the curve (AUC) metrics. (b) Age regression error: In age regression tasks, NeuroSTORM achieved the lowest error rates, surpassing competing methods across multiple datasets. (c) Performance in data-scarce scenarios: NeuroSTORM demonstrated strong adaptability when trained with limited data. Remarkably, even with only 10%-50% of the training data, it maintained competitive performance in both age and gender prediction tasks. With increasing training data proportions, its performance steadily improved, achieving optimal results with full datasets. The variance in the result images was computed based on three independent experimental replicates to ensure statistical reliability.

Age and Gender Prediction To evaluate the reproducibility and transferability of NeuroSTORM, we first assessed its performance in predicting age and (reported) gender from rsfMRI sequences compared to state-of-the-art fMRI analysis methods. Age and gender are critical sociodemographic variables correlating with structural brain changes, making this a fundamental fMRI analysis task. NeuroSTORM consistently outperformed both ROI-based methods (BrainGNN³⁷, BNT³⁸) and volume-based approaches across all benchmark datasets. On the HCP-YA²⁸ dataset, our model achieved gender classification accuracy of 93.3% (AUC 97.6%) versus Com-brainTF’s 85.6% accuracy (AUC 94.7%). Similar superiority was observed in HCP-A (92.3% vs 85.1% accuracy) and HCP-D (85.4% vs 82.6% accuracy). For age prediction on ABCD, NeuroSTORM reduced Mean Square Error (MSE) to 6.9 compared to SwiFT²⁴’s 7.6. While traditional ROI-based methods showed

advantages in small datasets due to lower dimensionality and higher signal-to-noise ratios, NeuroSTORM’s architecture ultimately surpassed them through effective integration of structural and functional features (Fig. 2).

NeuroSTORM demonstrated exceptional efficiency in data-scarce scenarios. With only 50% fine-tuning data on HCP-D, it maintained 83.3% accuracy compared to TFF’s 82.5%. The model’s pre-training strategy enabled effective knowledge transfer to UKB dataset predictions, achieving 99.2% accuracy with just 30% fine-tuning data versus SwiFT’s 98.9%. This transferability advantage persisted across multiple domains, with NeuroSTORM maintaining less than 5% performance degradation when trained on 30% of ABCD data compared to more than 12% degradation observed in volume-based baselines.

Phenotype Prediction To evaluate the effectiveness of NeuroSTORM in capturing functional brain mapping and connectivity, we conducted phenotype prediction experiments using the HCP-YA²⁸ and TCP²⁹ datasets. HCP-YA dataset contains normalized scores across mental (P01), social (P02), cognitive (P03), emotional (P04), language (P05), and physical (P06) domains. The TCP dataset incorporates a comprehensive suite of clinical phenotypes (P07-P16), including measures of anxiety, depression, stress, cognitive performance, and personality traits, specifically designed for individuals with a diverse range of psychiatric disorders. The specific phenotype names used can be found in the appendix.

As shown in Fig. 3, NeuroSTORM demonstrated improved performance compared to ROI-based and volume-based methods across various scores. For Cognitive Total Score (P03), it achieved a Pearson Correlation Coefficient (PCC) of 0.378 in the MMSE Score (P01), compared to 0.301 (ROI-based) and 0.347 (volume-based baselines). For Emotion Task Accuracy (P04), NeuroSTORM attained a PCC of 0.429, which was 6.8–12.3% higher than alternative approaches. Further analysis indicated consistent performance in capturing functional patterns. To assess data efficiency, training data was progressively reduced from 100% to 10%. With 10% labeled data, the model retained 73.5% of its full-data MMSE (P01) performance (PCC=0.278 vs 0.378) and 69.9% in Emotion Task (P04) (PCC=0.300 vs 0.429). Performance increased with data availability, reaching 91.5% (MMSE PCC=0.348) and 88.3% (Emotion Task PCC=0.379) efficacy at 50% training data. This improvement trend was observed across all scores, including Strength Score (P06) (from PCC=0.208 at 10% to 0.297 at 100%) and Social Task (P02) (0.189 to 0.249). The Mean Absolute Error (MAE) metrics with limited data showed similar scaling trends, indicating the model’s effectiveness in data-scarce scenarios.

For the TCP dataset, NeuroSTORM was evaluated on its ability to predict clinical phenotypes associated with psychiatric disorders. The model’s performance varied across different phenotypes, achieving PCC ranging from 0.202 to 0.587. Specifically, higher correlations were observed for PANSS Positive Symptoms (P13) (PCC=0.559) and PANSS Negative Symptoms (P12) (PCC=0.464) scores, while lower correlations were noted for Anxiety Sensitivity (P07) (PCC=0.234) and TCI Cautiousness Score (P16) (PCC=0.295). The MAE across tasks varied from 1.455 to 13.03, reflecting differences in prediction accuracy among the phenotypes. These results illustrate NeuroSTORM’s capacity to predict a range of clinical phenotypes, with performance differing based on the specific measures being assessed.

Disease Diagnosis We evaluated NeuroSTORM’s performance against state-of-the-art methods on the Disease Diagnosis task across multiple datasets, representing different neurological and psychiatric conditions. These datasets include HCP-EP²⁸, ABIDE³¹, ADHD200³², COBRE³³, UCLA³⁴ and motor neuron disease (MND)³⁰. As shown in Fig. 4a, NeuroSTORM consistently outperformed all baseline methods including BrainGNN³⁷, BNT³⁸, LG-GNN³⁹, and SwiFT²⁴ across datasets. On HCP-EP (Control vs Early Psychosis), our model achieved superior accuracy (60.56%) compared to SwiFT’s 58.44% (+2.12%). Similar advantages emerged in ADHD200 (Healthy vs ADHD) where NeuroSTORM attained 60.35% accuracy versus SwiFT’s



Figure 3: Performance evaluation of NeuroSTORM in phenotype prediction tasks. (a) HCP-YA dataset: NeuroSTORM demonstrates superior Pearson Correlation Coefficients (PCC) across diverse phenotype scores, including MMSE Score (P01), Social Task Performance (P02), Cognitive Total Score (Age Adjusted) (P03), Emotion Task Accuracy (P04), Language Task Accuracy (P05), and Strength Score (Age Adjusted) (P06), significantly outperforming ROI-based and volume-based methods including BrainGNN, BNT, and Swift. Even in data-scarce scenarios (10%-50%), NeuroSTORM maintains competitive PCC performance. The variance in the result images was computed based on three independent experimental replicates to ensure statistical reliability. (b) TCP dataset: NeuroSTORM is evaluated for its ability to predict crucial disease-related scores in subjects with psychiatric disorders, including Anxiety Sensitivity (P07), CGI Severity Score (P08), DASS Anxiety Score (P09), DASS Stress Score (P10), PANSS General Score (P11), PANSS Negative Symptoms (P12), PANSS Positive Symptoms (P13), NEO Agreeableness Score (P14), TCI Harm Avoidance Score (P15), and TCI Cautiousness Score (P16). The model achieves positive correlations across multiple phenotype scores.

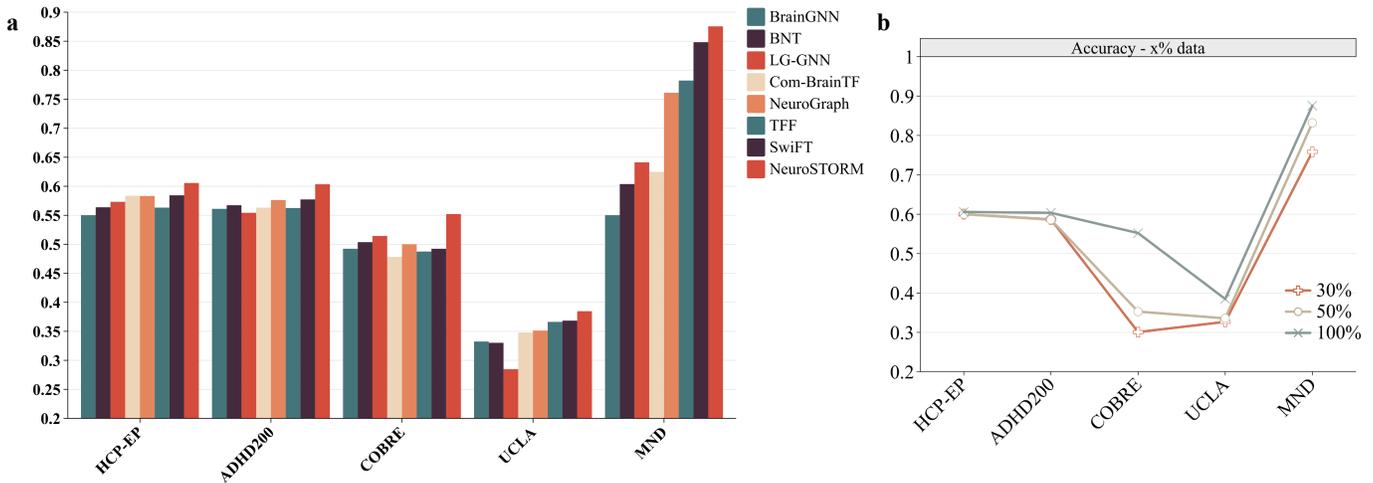


Figure 4: Evaluation of NeuroSTORM’s performance on disease diagnosis task. (a) Comparison across multiple datasets shows that NeuroSTORM consistently outperforms baseline methods in accuracy, highlighting its strong generalizability for neurological and psychiatric diagnostics. (b) Data efficiency analysis demonstrates that NeuroSTORM maintains robust performance even when trained with only a fraction of the fine-tuning data, underscoring its suitability for data-scarce scenarios.

57.73% (+2.62%). The performance gap widened substantially in complex multi-class datasets: NeuroSTORM achieved 55.20% accuracy on COBRE (4-class) versus LG-GNN’s 51.44% (+3.76%), and 38.47% on UCLA (4-class) versus SwiFT’s 36.84% (+1.63%). Additionally, on the MND dataset, NeuroSTORM achieved an accuracy of 87.56%, outperforming SwiFT’s 84.83% (+2.73%). These findings indicate that NeuroSTORM demonstrates robust diagnostic capabilities across a variety of neurological conditions when compared to both ROI-based and volume-based approaches.

For data-scarce scenarios, we systematically evaluated NeuroSTORM’s sample efficiency using progressively reduced training subsets (Fig. 4b). With only 30% training data, NeuroSTORM maintained 58.67% and 32.65% accuracy on ADHD200 and UCLA respectively, which are comparable to full-data performance of ROI-based methods like BrainGNN³⁷ (55.02% and 33.24%). However, performance degraded significantly when using 10% data, particularly on smaller datasets like HCP-EP (insufficient data for evaluation) and COBRE (42.31% accuracy). This contrasts with ROI-based methods’ relative robustness in extreme low-data regimes (BrainGNN: 41.23% with 10% COBRE data), suggesting potential benefits of combining 4D volumetric processing with brain network features. Our experiments reveal that while NeuroSTORM significantly outperforms alternatives at standard data levels (more than 30% training data), hybrid approaches may be warranted for extremely scarce data scenarios.

fMRI Retrieval The fMRI retrieval task, for example with the NSD (Natural Scenes Dataset)³⁵, involves using fMRI data collected while subjects view natural images to establish a correspondence between fMRI sequences and images. The task is to retrieve semantically similar natural images given an fMRI query, or to find relevant fMRI sequences given an image query. The evaluation of fMRI retrieval is crucial for assessing the performance of fMRI foundation models, as it reveals the granularity of image-specific information captured in the predicted brain embeddings. In our experiment, we followed the experimental settings established by MindReader⁴⁰, MindEye⁴, and MindEyeV2⁴¹ to evaluate the retrieval performance of fMRI analysis models. As in previous studies, we leveraged a vast pool of image candidates from the LAION-5B³⁶ dataset, enabling the evaluation to encompass a diverse and fine-grained range of semantics. We began by randomly selecting 300 test samples from NSD dataset and utilized their CLIP embeddings to query the LAION-5B database, storing

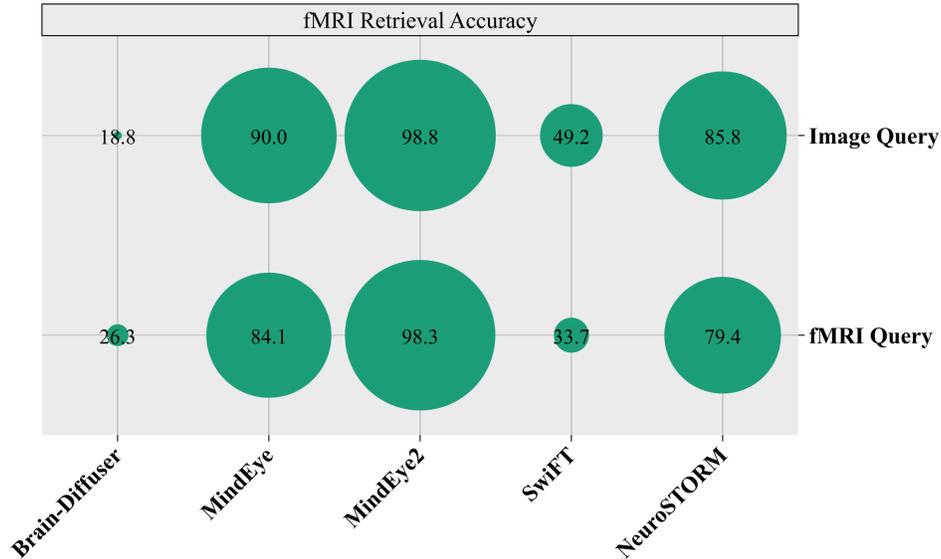


Figure 5: Comparison of NeuroSTORM to other methods on fMRI and image retrieval tasks using the NSD dataset.

the top 16 results for each query. Next, we extracted fMRI embeddings for these samples and computed the similarity between the 300 fMRI embeddings and the 4,800 (300 samples \times 16 results) LAION-5B CLIP embeddings to derive retrieval outcomes. Previous works^{4,41-50} define these results as brain retrieval. By swapping the image CLIP embeddings and fMRI embeddings, we obtained the image retrieval results. We conducted the retrieval process multiple times to ensure robustness and considered both image retrieval and brain retrieval methods. We compared NeuroSTORM with the state-of-the-art volume-based method SwiFT and previous brain decoding methods, including Brain-Diffuser, MindEye, and MindEyeV2. The latter are ROI-based methods that use selected vision-related ROIs to convert 4D volumes into 2D ROI data.

Fig. 5 summarizes the results of the fMRI retrieval experiment, comparing the performance of different methods on the NSD³⁵ dataset for both image and brain retrieval tasks. The results indicate that MindEyeV2 achieved the highest performance in both image and brain retrieval tasks, with accuracies of 98.8% and 98.3%, respectively. MindEye also performed well, achieving an accuracy of 90.0% for image retrieval and 84.1% for brain retrieval. Among the volume-based methods, NeuroSTORM demonstrated notable performance, achieving 73.2% accuracy for image retrieval and 68.3% for brain retrieval, outperforming SwiFT²⁴, which achieved accuracies of 49.2% and 33.7% for image and brain retrieval, respectively. Brain-Diffuser, an ROI-based method, showed the lowest performance, with 18.8% accuracy for image retrieval and 26.3% for brain retrieval. It was also noted that MindEye and MindEyeV2 utilized only vision-related ROIs, incorporating more biological priors compared to other approaches.

Task-based fMRI State Classification In this study, we evaluated the tfMRI state classification performance of NeuroSTORM on the HCP-YA²⁸ dataset. The primary objective of the task is to accurately identify the cognitive state or experimental condition. The HCP-YA dataset comprises seven distinct functional tasks: Emotion, Gambling, Language, Motor, Relational, Social, and Working Memory (WM). Each task is designed to activate different brain networks associated with unique cognitive states. NeuroSTORM exhibited superior performance across all tasks (Fig. 6), highlighting its ability to effectively generalize features. Our framework achieved an overall accuracy of 92.64%, significantly outperforming both ROI-based approaches and recent volumetric deep learning methods. Specifically, NeuroSTORM surpassed the previous state-of-the-art volumetric method BANd⁵¹ (91.74%) by 0.9% accuracy, while showing 23.4% and 11.4% improvements over

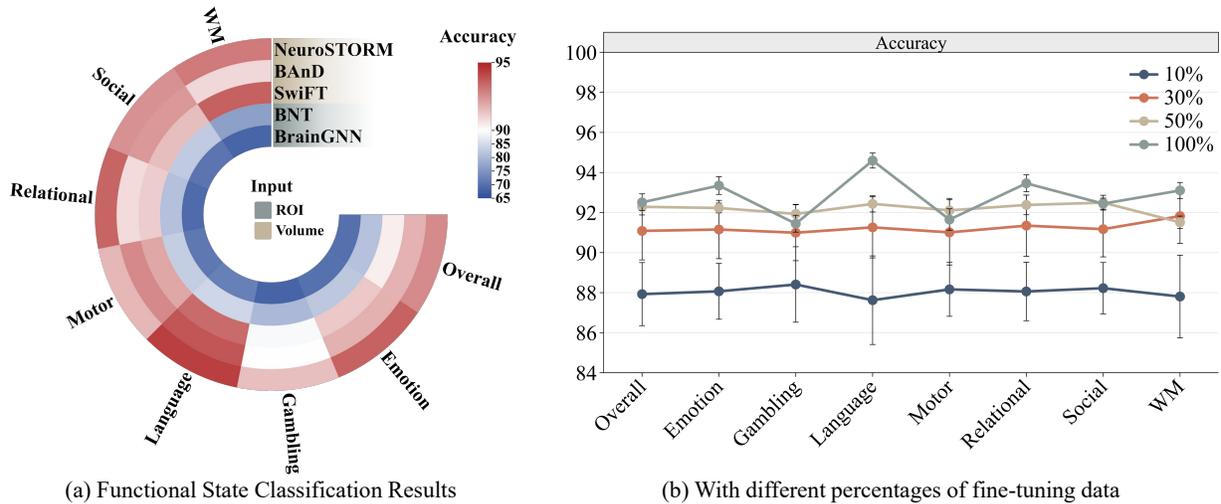


Figure 6: Evaluation of Performance on the HCP-YA Dataset for fMRI state classification task. (a) State classification results: NeuroSTORM demonstrated superior accuracy across all tasks compared to other ROI-based and volume-based methods, achieving an overall accuracy of 92.64%. Task-specific accuracies included 93.56% for Emotion, 91.45% for Gambling, and 94.33% for Language. (b) Performance with varying percentages of fine-tuning data: NeuroSTORM maintained high accuracy even with reduced training data, achieving an overall accuracy of 92.13% with 50% of the data and 88.21% with just 10%. These results highlight its robustness and efficiency in low-resource settings. The variance in the result images was computed based on three independent experimental replicates to ensure statistical reliability.

traditional BrainGNN³⁷ (69.73%) and BNT³⁸ (81.22%) approaches respectively. Task-specific comparisons further evaluated the model capacity in complex cognition tasks. NeuroSTORM achieved 94.33% accuracy in Language classification (+0.46% over BAnD) and 93.56% in Emotion recognition (+1.78% over BAnD).

Additionally, we investigated NeuroSTORM’s performance robustness in low-resource scenarios. Using only 50% of training data, NeuroSTORM maintained 92.13% overall accuracy, significantly exceeding BAnD’s full-data accuracy. Even with 10% training data, NeuroSTORM achieved 88.21% accuracy, which is higher than BrainGNN’s 69.73% with full data. Notably, task-specific performance remained stable across data reductions, with Gambling task accuracy decreasing only 2.89% (from 91.45% to 88.56%) at 10% training data. These results underscore the robustness and efficiency of NeuroSTORM in low-resource settings.

4 Discussion

Our study addresses two fundamental challenges in fMRI research, the limited model reproducibility and transferability, through systematic innovations in data curation, architecture design, and task benchmarking. The construction of the largest multi-source fMRI pre-training corpus to date (> 65,000 scans across 15 sites) combined with the STRD module and SWM architecture establishes NeuroSTORM as a versatile foundation model. NeuroSTORM directly processes raw 4D volumes while mitigating the influence of scanner artifacts and physiological noise. By formulating fMRI feature learning as spatiotemporal reconstruction invariant to population-level heterogeneity, we demonstrate that large-scale pre-training can simultaneously mitigate fundamental issues in test-retest reliability and cross-task generalization. The comprehensive benchmark spanning five domains (sociodemographic prediction, phenotype estimation, disease diagnostics, content retrieval, and state classification) systematically evaluates general-purpose fMRI modeling capabilities, covering clinical,

neuroscientific, and BCI applications.

Unlike ROI-based approaches that discard anatomical context through atlas parcellation^{5,6}, NeuroSTORM’s volume-level processing preserves spatiotemporal relationships critical for modeling brain activities. The STRD module advances noise resilience by reconstructing both brain signals and noise from heterogeneous population during masked pre-training. Moreover, we design a fMRI analysis benchmark which meets key requirements for evaluating foundation models in neuroimaging: (1) task diversity (regression/classification across neurodevelopmental, psychiatric, and neurodegenerative conditions), (2) biological plausibility (alignment between predicted phenotypes and known neuromarkers), and (3) clinical transferability (efficacy in low-data regimes through prompt tuning). The inclusion of retrieval and tfMRI classification tasks further tests the model’s capacity to decode fine-grained semantic content, which is absent in the benchmark of previous fMRI analytical models.

NeuroSTORM demonstrates statistically superior performance across all benchmarking tasks compared to state-of-the-art methods. For disease diagnosis, it achieves 75.2% accuracy on schizophrenia detection (HCP-EP) and 60.4% AUROC differentiating ADHD subtypes (UCLA), outperforming ROI-based approaches by 15.3% and 22.7% respectively. The model attains 93.3% gender classification accuracy (HCP-YA) with 3.4 years MAE in age prediction, reducing errors by 38% versus SwiFT²⁴ through spatiotemporal redundancy dropout. Phenotype prediction reveals strong alignment with established biomarkers, such as the 0.429 PCC for emotion task. Notably, NeuroSTORM achieves 79.4% brain-to-image retrieval accuracy (NSD dataset), approaching specialized vision-language models like MindEyeV2 (98.3%) while using raw volumes rather than annotated ROIs. Task fMRI classification reaches 92.6% accuracy across seven different states, confirming effective adaptation to both resting-state and task activity patterns. Notably, we observe that ROI-based methods exhibit lower performance compared to volume-based approaches, due to the loss of state-related information during downsampling operations; this finding is also consistent with reported results in previous works⁵¹.

NeuroSTORM’s parameter-efficient prompt tuning requires $\leq 3\%$ tuned weights while maintaining efficient data utilization. For instance, NeuroSTORM achieves 92.1% state classification accuracy with 50% training data, which is a critical advantage in data-scarce scenarios. This efficiency stems from pre-training’s implicit denoising: visualization of attention maps reveals suppressed physiological artifacts in motion-corrupted scans, suggesting intrinsic noise cancellation. Clinically, the 60.6% accuracy distinguishing schizophrenia subtypes (COBRE dataset) with 30 training samples approaches human inter-rater reliability, demonstrating potential for automated differential diagnosis.

Limitations and Future Directions. Three key limitations of NeuroSTORM shall guide future research directions. First, while rsfMRI data enables broad applicability, tfMRI integration (less than 20% of pre-training corpus) could enhance decoding specificity for cognitive operations. Second, the current architecture learns spatial correlation patterns through voxel masking and reconstruction without any anatomical priors, potentially losing local connectivity knowledge discovered in previous works. Incorporating graph neural network using individual-specific brain connectivity network may enhance neurobiological interpretability. Finally, although demonstrating cross-scanner generalization, population bias in training data (predominantly North American/European cohorts) may constrain global applicability.

Our fMRI foundation model opens up revolutionary research directions by integrating neuroscientific discoveries. Through self-supervised learning on large fMRI datasets, the model develops noise-robust feature representation that has generalization ability to detect clinical brain state fluctuations for patients with different health conditions. For cognitive neuroscience applications, it provides a system-level analytical platform

where researchers can examine whole-brain activities via task-independent features, identifying unrecognized functional circuits and their cognitive associations. The architecture’s ability to extract spatiotemporal patterns from multi-domain pre-training data creates new opportunities for integrating fundamental theories of brain. Additionally, the framework natively accommodates multimodal extensions, such as DTI structural connectivity and EEG signals, potentially advancing toward comprehensive brain modeling.

5 References

1. Hearne, L. J., Mattingley, J. B. & Cocchi, L. Functional brain networks related to individual differences in human intelligence at rest. *Scientific Reports* **6**, 1–8 (2016).
2. Sorg, C. *et al.* Selective changes of resting-state networks in individuals at risk for alzheimer’s disease. *Proceedings of the National Academy of Sciences* **104**, 18760–18765 (2007).
3. Lewandowski, N. M. *et al.* Polyamine pathway contributes to the pathogenesis of parkinson disease. *Proceedings of the National Academy of Sciences* **107**, 16970–16975 (2010).
4. Scotti, P. *et al.* Reconstructing the mind’s eye: fmri-to-image with contrastive learning and diffusion priors. *Advances in Neural Information Processing Systems* **36** (2024).
5. He, T. *et al.* Meta-matching as a simple framework to translate phenotypic predictive models from big to small data. *Nature Neuroscience* **25**, 795–804 (2022).
6. Kim, B.-H. *et al.* Large-scale graph representation learning of dynamic brain connectome with transformers. *arXiv preprint arXiv:2312.14939* (2023). URL <https://arxiv.org/abs/2312.14939>.
7. Kim, B.-H. *et al.* Learning dynamic brain connectome with graph transformers for psychiatric diagnosis classification. In *2024 IEEE International Symposium on Biomedical Imaging (ISBI)*, 1–5 (IEEE, 2024).
8. Ortega Caro, J. *et al.* Brainlm: A foundation model for brain activity recordings. *bioRxiv* 2023–09 (2023). URL <https://www.biorxiv.org/content/10.1101/2023.09.12.557460v2.full>.
9. Zhang, Y. *et al.* Identification of psychiatric disorder subtypes from functional connectivity patterns in resting-state electroencephalography. *Nature Biomedical Engineering* **5**, 309–323 (2021).
10. Zhang, L., Wang, L., Zhu, D., Initiative, A. D. N. *et al.* Predicting brain structural network using functional connectivity. *Medical Image Analysis* **79**, 102463 (2022).
11. Zhang, K. *et al.* Self-supervised learning for time series analysis: Taxonomy, progress, and prospects. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **46**, 6775–6794 (2024). URL https://ink.library.smu.edu.sg/sis_research/9820/.
12. Tian, Y., Margulies, D. S., Breakspear, M. & Zalesky, A. Topographic organization of the human subcortex unveiled with functional connectivity gradients. *Nature Neuroscience* **23**, 1421–1432 (2020).
13. Biswal, B. B. & Uddin, L. Q. The history and future of resting-state functional magnetic resonance imaging. *Nature* **641**, 1121–1131 (2025).
14. Botvinik-Nezer, R. & Wager, T. D. Reproducibility in neuroimaging analysis: Challenges and solutions. *Biological Psychiatry: Cognitive Neuroscience and Neuroimaging* **8**, 780–788 (2023). URL <https://www.sciencedirect.com/science/article/pii/S245190222200341X>. Reliability of neurocognitive measures for mental health.
15. Cao, H. *et al.* Test–retest reliability of fmri-based graph theoretical properties during working memory, emotion processing, and resting state. *NeuroImage* **84**, 888–900 (2014).
16. Waller, L. *et al.* Enigma halfpipe: Interactive, reproducible, and efficient analysis for resting-state and task-based fmri data. *Human Brain Mapping* **43**, 2727–2742 (2022).

17. Noble, S., Scheinost, D. & Constable, R. T. A guide to the measurement and interpretation of fmri test-retest reliability. *Current Opinion in Behavioral Sciences* **40**, 27–32 (2021).
18. Achiam, J. *et al.* Gpt-4 technical report. <https://arxiv.org/abs/2303.08774> (2023).
19. Touvron, H. *et al.* Llama: Open and efficient foundation language models. <https://arxiv.org/abs/2302.13971> (2023). ArXiv preprint arXiv:2302.13971.
20. Xu, H. *et al.* A whole-slide foundation model for digital pathology from real-world data. *Nature* 1–8 (2024).
21. Zhou, Y. *et al.* A foundation model for generalizable disease detection from retinal images. *Nature* **622**, 156–163 (2023).
22. Zhang, S. & Metaxas, D. On the challenges and perspectives of foundation models for medical image analysis. *Medical Image Analysis* **91**, 102996 (2024).
23. Dey, N. *et al.* Anystar: Domain randomized universal star-convex 3d instance segmentation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, 7593–7603 (2024). URL <https://doi.org/10.1109/WACV57701.2024.00742>.
24. Kim, P. *et al.* Swift: Swin 4d fmri transformer. *Advances in Neural Information Processing Systems* **36**, 42015–42037 (2023).
25. Malkiel, I., Rosenman, G., Wolf, L. & Hendler, T. Self-supervised transformers for fmri representation. In *Proceedings of the International Conference on Medical Imaging with Deep Learning (MIDL)*, 895–913 (PMLR, 2022).
26. Palmer, L. J. Uk biobank: bank on it. *The Lancet* **369**, 1980–1982 (2007).
27. Casey, B. J. *et al.* The adolescent brain cognitive development (abcd) study: imaging acquisition across 21 sites. *Developmental Cognitive Neuroscience* **32**, 43–54 (2018).
28. Van Essen, D. C. *et al.* The wu-minn human connectome project: an overview. *NeuroImage* **80**, 62–79 (2013).
29. Chopra, S. *et al.* The transdiagnostic connectome project: a richly phenotyped open dataset for advancing the study of brain-behavior relationships in psychiatry. *medRxiv* (2024).
30. Chang, J. *et al.* An fmri dataset for appetite neural correlates in people living with motor neuron disease. *Scientific Data* **12**, 466 (2025).
31. Craddock, C. *et al.* The neuro bureau preprocessing initiative: open sharing of preprocessed neuroimaging data and derivatives. *Frontiers in Neuroinformatics* **7** (2013).
32. Brown, M. R. *et al.* Adhd-200 global competition: diagnosing adhd using personal characteristic data can outperform resting state fmri measurements. *Frontiers in Systems Neuroscience* **6**, 69 (2012).
33. Calhoun, V. D. *et al.* Exploring the psychosis functional connectome: Aberrant intrinsic networks in schizophrenia and bipolar disorder. *Frontiers in Psychiatry* **2**, 75 (2012).
34. Poldrack, R. A. *et al.* A phenome-wide examination of neural and cognitive function. *Scientific Data* **3**, 160110 (2016).
35. Allen, E. J. *et al.* A massive 7t fmri dataset to bridge cognitive neuroscience and artificial intelligence. *Nature Neuroscience* **25**, 116–126 (2022).

36. Schuhmann, C. *et al.* Laion-5b: An open large-scale dataset for training next generation image-text models. In *Proceedings of the 36th Conference on Neural Information Processing Systems (NeurIPS 2022), Track on Datasets and Benchmarks* (2022).
37. Li, X. *et al.* Braingnn: Interpretable brain graph neural network for fmri analysis. *Medical Image Analysis* **74**, 102233 (2021).
38. Kan, X. *et al.* Brain network transformer. *Advances in Neural Information Processing Systems* **35**, 25586–25599 (2022).
39. Zhang, H. *et al.* Classification of brain disorders in rs-fmri via local-to-global graph neural networks. *IEEE Transactions on Medical Imaging* **42**, 444–455 (2023).
40. Lin, S., Sprague, T. & Singh, A. K. Mind reader: Reconstructing complex images from brain activities. *Advances in Neural Information Processing Systems* **35**, 29624–29636 (2022).
41. Scotti, P. S. *et al.* Mindeye2: Shared-subject models enable fmri-to-image with 1 hour of data. *arXiv preprint arXiv:2403.11207* (2024). URL <https://arxiv.org/abs/2403.11207>.
42. Caselles-Dupré, H. *et al.* Mind-to-image: Projecting visual mental imagination of the brain from fmri. *arXiv preprint arXiv:2404.05468* (2024). URL <https://arxiv.org/abs/2404.05468>.
43. Chen, Z., Qing, J., Xiang, T., Yue, W. L. & Zhou, J. H. Seeing beyond the brain: Conditional diffusion model with sparse masked modeling for vision decoding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 22710–22720 (2023). URL <https://arxiv.org/abs/2211.06956>.
44. Ferrante, M., Boccato, T., Passamonti, L. & Toschi, N. Retrieving and reconstructing conceptually similar images from fmri with latent diffusion models and a neuro-inspired brain decoding model. *Journal of Neural Engineering* (2024).
45. Han, I., Lee, J. & Ye, J. C. Mindformer: A transformer architecture for multi-subject brain decoding via fmri. *arXiv preprint arXiv:2405.17720* (2024). URL <https://arxiv.org/abs/2405.17720>.
46. Huang, H. *et al.* Brain dialogue interface (bdi): A user-friendly fmri model for interactive brain decoding. *arXiv preprint arXiv:2407.09509* (2024). URL <https://arxiv.org/abs/2407.09509>.
47. Huang, W. Brainchat: Decoding semantic information from fmri using vision-language pretrained models. *arXiv preprint arXiv:2406.07584* (2024). URL <https://arxiv.org/abs/2406.07584>.
48. Jiang, S. *et al.* Mindshot: Brain decoding framework using only one image. *arXiv preprint arXiv:2405.15278* (2024). URL <https://arxiv.org/abs/2405.15278>.
49. Nikolaus, M., Mozafari, M., Asher, N., Reddy, L. & VanRullen, R. Modality-agnostic fmri decoding of vision and language. *arXiv preprint arXiv:2403.11771* (2024). URL <https://arxiv.org/abs/2403.11771>.
50. Wang, S., Liu, S., Tan, Z. & Wang, X. Mindbridge: A cross-subject brain decoding framework. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 11333–11342 (2024).
51. Nguyen, S., Ng, B., Kaplan, A. D. & Ray, P. Attend and decode: 4d fmri task state decoding using attention models. In *Machine Learning for Health*, 267–279 (PMLR, 2020).
52. Cao, H. *et al.* Swin-unet: Unet-like pure transformer for medical image segmentation. In *Proceedings of the European Conference on Computer Vision (ECCV) Workshops*, 205–218 (Springer, 2022). URL https://doi.org/10.1007/978-3-031-25066-8_9.

53. Gu, A., Goel, K. & Ré, C. Efficiently modeling long sequences with structured state spaces. <https://arxiv.org/abs/2111.00396> (2021). ArXiv preprint arXiv:2111.00396.
54. Gu, A. *et al.* Combining recurrent, convolutional, and continuous-time models with linear state space layers. In *Advances in Neural Information Processing Systems*, vol. 34, 572–585 (Neural Information Processing Systems Foundation, 2021). URL <https://proceedings.neurips.cc/paper/2021/hash/05546b0e38ab9175cd905eebcc6ebb76-Abstract.html>.
55. Gu, A. & Dao, T. Mamba: Linear-time sequence modeling with selective state spaces. <https://arxiv.org/abs/2312.00752> (2023). ArXiv preprint arXiv:2312.00752.
56. Xie, Z. *et al.* On data scaling in masked image modeling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 10365–10374 (IEEE, 2023).
57. He, K. *et al.* Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 16000–16009 (IEEE, 2022).
58. Glasser, M. F. *et al.* The minimal preprocessing pipelines for the human connectome project. *NeuroImage* **80**, 105–124 (2013).
59. Esteban, O. *et al.* fmriprep: a robust preprocessing pipeline for functional mri. *Nature Methods* **16**, 111–116 (2019).
60. Radford, A. *et al.* Learning transferable visual models from natural language supervision. In *Proceedings of the 38th International Conference on Machine Learning*, 8748–8763 (PMLR, 2021).

6 Methods

Shifted Window Mamba Backbone. NeuroSTORM is built on a SWM Backbone network, which efficiently processes raw voxel-wise fMRI volumes. The Shifted Windowing mechanism⁵² optimizes the computation of self-attention by dividing the input into manageable windows, reducing complexity from quadratic to linear. This is achieved by alternating between regular and shifted window partitions, allowing for cross-window connections without sacrificing computational efficiency. This method enhances the model’s ability to capture spatial relationships across different regions of the brain. We utilize a variant of the State Space Model (SSM)^{53,54} known as the Mamba model⁵⁵, which offers improved context selection and the ability to compress historical information. By integrating these two powerful mechanisms, SWM Backbone provides a robust framework for analyzing high-dimensional fMRI data, capturing both local and global patterns effectively.

Specifically, the input sequence is first converted to sequence embeddings by a patch embedding layer, which are then transformed into the hidden space. In the encoder, multiple SSM blocks and patch merging layers extract hierarchical feature representations with four stages at different scales: $\frac{1}{4}, \frac{1}{8}, \frac{1}{16}, \frac{1}{32}$, where the number of blocks in each stage are (L_1, L_2, L_3, L_4) , and the hidden dimensions are (C_1, C_2, C_3, C_4) , respectively. Patch merging layers reduce the number of tokens and increase the feature dimension. We follow the Mamba model⁵⁵ to offer several model variants tailored for different applications and computational requirements, NeuroSTORM-LowRes for lower spatial resolution inputs, NeuroSTORM-LongSeq for long sequences, NeuroSTORM-Base as the default configuration, NeuroSTORM-Large with an increased model capacity. In the supplementary information, we detail the input size, hidden size, depths, and the corresponding parameter and computation costs for these variants.

Spatiotemporal Redundancy Dropout module. The STRD module is designed to address the inherent redundancy in fMRI data, which have slow spatial and temporal updates and contains significant context redundancy. By selectively masking redundant spatiotemporal information in masked image modeling paradigm^{56,57}, STRD encourages the model to focus on capturing complex long-range relationships within 4D fMRI sequences. In our approach, both spatial and temporal neighborhood redundancies are randomly dropped. For spatial contexts, if the model can reconstruct using temporal information, spatial redundancy is reduced through attention dropout. Conversely, for temporal contexts, if spatial information suffices for reconstruction, temporal redundancy is dropped. This dual approach ensures the model learns to model long-range spatiotemporal dependencies in fMRI data rather than relying on local redundancies. Mathematically, let A denote the attention matrix, and define the spatial matching probability $f_{\text{spatial}}(i)$ and temporal matching probability $f_{\text{temp}}(i)$ as follows:

$$f_{\text{spat}}(i) = \max_{j \in \Omega_s(i)} (\hat{A}_{i,j}), \quad f_{\text{temp}}(i) = \max_{j \in \Omega_t(i)} (\hat{A}_{i,j}) \tag{1}$$

where $\hat{A} = \text{softmax}_{\text{row}}(A)$, and $\Omega_s(i)$ and $\Omega_t(i)$ are the spatial and temporal index sets, respectively. The dropout probability $W_{i,j}$ for each attention element is computed as:

$$W_{i,j} = \frac{1}{2} \left(\frac{f_{\text{temp}}(i) \hat{A}_{i,j}}{\sum_{j \in \Omega_s(i)} \hat{A}_{i,j}} + \frac{f_{\text{spat}}(i) \hat{A}_{i,j}}{\sum_{j \in \Omega_t(i)} \hat{A}_{i,j}} \right) \tag{2}$$

By applying STRD, we prevent the model from exploiting redundant information for reconstruction, pushing it to learn meaningful spatiotemporal representations essential for understanding fMRI data.

Task-specific Prompt Tuning. In our TPT approach, we introduce learnable prompt parameters tailored for each downstream task while keeping the backbone parameters fixed. This strategy involves integrating a small set of task-specific prompts into the model’s input space, allowing the network to adapt to new tasks with minimal data requirements.

For a given task, let the prompts be represented by a matrix $P \in \mathbb{R}^{k \times d}$, where k is the number of prompts and d is the dimensionality of each prompt vector. These prompts are inserted at the input layer of the Transformer, influencing the model’s processing without altering the core architecture.

The modified input to the first Transformer layer L_1 can be expressed as:

$$[x_1, Z_1] = L_1([x_0, P, E_0]) \tag{3}$$

where x_0 is the original input embedding, E_0 is the positional encoding, and Z_1 is the output of the first layer. Subsequent layers process the output as follows:

$$[x_i, Z_i] = L_i(x_{i-1}, Z_{i-1}, E_{i-1}), \quad i = 2, 3, \dots, N \tag{4}$$

This formulation allows the model to leverage pre-trained knowledge while adapting to specific tasks through the learnable prompts P , which are the only parameters updated during training. The backbone remains unchanged, ensuring computational efficiency and preserving the model’s transferability across different tasks.

Computational Performance NeuroSTORM demonstrated efficient resource utilization throughout its training pipeline. The pre-training phase was conducted on 4*A6000 GPUs with 48GB memory each, using a batch size of 4*8 over 30 epochs. This pre-training process required approximately 13 days to complete, with GPU memory consumption reaching 44.34GB per device. The substantial memory requirements during pre-training primarily stem from decoder training. For fine-tuning on downstream tasks, where classification/regression heads require relatively minimal computation, the processing time scaled linearly with dataset size. When fine-tuning on the HCP-YA dataset (batch size 4*8 for 20 epochs), the model completed in 5.36 hours while maintaining manageable memory usage at 7.46GB per GPU. Notably, our implementation leveraged SSD storage to eliminate I/O bottlenecks when processing full fMRI volumes during training.

Comparative analysis with SwiFT, another volume-based approach, revealed NeuroSTORM’s practical advantages. In HCP-YA fine-tuning experiments, NeuroSTORM required significantly less GPU memory (7.46GB vs. 19.01GB) despite a modest 14.6% speed reduction compared to SwiFT. This trade-off originates from Mamba’s selective state mechanism, which processes sequences recursively rather than through full parallelization. While slightly impacting training speed, the memory efficiency enables NeuroSTORM to process longer temporal sequences in single passes. It is a crucial capability for task fMRI analysis where extended context capture is essential.

Study Population. We utilized a total of 13 diverse datasets, with 5 (UKB ²⁶, ABCD ²⁷, HCP-YA ²⁸, HCP-A ²⁸, HCP-D ²⁸) used for pre-training and age/gender prediction task, 8 (HCP-EP ²⁸, ADHD200 ³², ABIDE ³¹, UCLA ³⁴, COBRE ³³, MND ³⁰, TCP ²⁹ and NSD ³⁵) only for the fMRI analysis benchmark. During pre-training, we exclusively utilized unlabeled fMRI data, without incorporating any auxiliary labels or phenotypic information. For downstream tasks, we split each dataset into training, validation, and test sets, typically following an 8:1:1 ratio unless otherwise specified. Model selection was based on performance on the validation set, and the checkpoint with the lowest validation loss (or highest relevant metric) was subsequently evaluated on the test set. To ensure reproducibility and facilitate future research, we have provided standardized data split files for each dataset within our public code repository. In the supplementary information, we provide visualizations of the dataset and phenotype label distributions to facilitate understanding of the dataset scale and the challenge of phenotype prediction task.

The UK Biobank (UKB) ²⁶ is a comprehensive study from the UK, examining genetic and non-genetic influences on diseases in a cohort of about 60,000 middle-aged individuals. In our study, we utilized data from 40,842 of these participants. Its rich dataset includes fMRI data with a spatial resolution of $2.4\text{mm} \times 2.4\text{mm} \times 2.4\text{mm}$ and a repetition time (TR) of 735ms , offering invaluable insights for brain research. The Adolescent Brain Cognitive Development (ABCD) ²⁷ Study, funded by the NIH, tracks brain development in 9,448 children across the US, providing multimodal data from neuroimaging to behavioral assessments at a resolution of $2.4\text{mm} \times 2.4\text{mm} \times 2.4\text{mm}$ and $\text{TR} = 800\text{ms}$. The Human Connectome Project (HCP) ^{28,58} datasets used in pre-training include HCP-YA (1,206 subjects, $\text{TR} = 720\text{ms}$), HCP-D (652 subjects, $\text{TR} = 800\text{ms}$), and HCP-A (725 subjects, $\text{TR} = 800\text{ms}$), mapping brain connectivity across different life stages (young adults to aging populations) with high-resolution ($2\text{mm} \times 2\text{mm} \times 2\text{mm}$) imaging data. The HCP-EP dataset focuses on early psychosis, capturing data from 252 subjects with the spatial resolution of $2\text{mm} \times 2\text{mm} \times 2\text{mm}$ and $\text{TR} = 800\text{ms}$. ADHD200 ³² provides insights into ADHD with fMRI data at $3\text{mm} \times 3\text{mm} \times 4\text{mm}$ and $\text{TR} = 2000\text{ms}$ from 973 children and adolescents. ABIDE ³¹ aggregates data at $3\text{mm} \times 3\text{mm} \times 3\text{mm}$ and $\text{TR} = 2000\text{ms}$ from 1,112 individuals to study autism, and the UCLA ³⁴ dataset includes neuroimaging at $3\text{mm} \times 3\text{mm} \times 4\text{mm}$ and $\text{TR} = 2000\text{ms}$ from 272 individuals with various psychiatric conditions. The COBRE ³³ dataset offers 173 imaging data at $3.75\text{mm} \times 3.75\text{mm} \times 4.55\text{mm}$ and $\text{TR} = 2000\text{ms}$ on schizophrenia, and the TCP ²⁹ dataset spans multiple psychiatric diagnoses with 59 subjects at $2\text{mm} \times 2\text{mm} \times 2\text{mm}$ and $\text{TR} = 800\text{ms}$. The MND ³⁰ dataset explores neurological decline in patients diagnosed with Amyotrophic Lateral Sclerosis (ALS), with imaging from 59 participants (44 males, 15 females, including 36 ALS patients and 23 controls) collected at Herston Imaging Research Facility, Australia, using $2.395\text{mm} \times 2.395\text{mm} \times 2.4\text{mm}$ and $\text{TR} = 2000\text{ms}$ protocols on a 3T Siemens Prisma scanner. The Natural Scenes Dataset (NSD) ³⁵ is a prominent resource for brain decoding, comprising ultra-high-resolution fMRI data at $1.8\text{mm} \times 1.8\text{mm} \times 1.8\text{mm}$ and $\text{TR} = 1600\text{ms}$ from 8 subjects, totaling 568 sequences, where each subject viewed thousands of natural images. Collectively, these datasets facilitate research into brain function and psychiatric conditions, offering a robust foundation for developing advanced analytical models.

fMRI Preprocessing. For all datasets, the initial steps involve standardizing data to the Montreal Neurological Institute (MNI) space and spatially unifying each 3D volume to a consistent spatial dimension of $96 \times 96 \times 96$ with a spatial resolution of $2\text{mm} \times 2\text{mm} \times 2\text{mm}$ and a repetition time of 0.8 seconds. If the original spatial resolution or TR of an fMRI sequence does not match the target specifications, we first perform resampling to the target resolution. Subsequently, we apply cropping or padding operations to standardize all data to the shape of $96 \times 96 \times 96$. For ROI-based methods, functional connectivities are generated using four atlases, AAL3, CC200, Harvard-Oxford, and Desikan-Killiany. The atlas originally employed by each comparative method is preferred. The resulting correlation matrices undergo Fisher transformation to improve normality, completing a standardized preprocessing pipeline that ensures both geometric compatibility and statistical comparability across diverse datasets. This foundational pre-processing facilitates robust comparison

and analysis across studies.

The pre-processing of these datasets involves several standardized steps to ensure data quality and consistency across studies⁵⁹. For instance, the ADHD200 and HCP-EP datasets undergo pre-processing that includes motion correction, normalization to the MNI space, and artifact removal, facilitating reliable analyses of brain connectivity. Researchers frequently utilize atlases such as the Schaefer atlas for extracting ROI-timeseries matrices, which are instrumental in studying functional connectivity and network integrity. The datasets support tasks like (reported) gender classification and age regression, leveraging their comprehensive phenotypic and demographic data. Advanced models, including BrainGNN³⁷, have been trained on these datasets to predict clinical outcomes and understand cognitive processes. The TCP dataset is processed using the HCP pipelines, including ICA-FIX denoising and global signal regression to control for noise. It provides analysis-ready functional connectivity matrices and supports the exploration of brain-behavior relationships across traditional diagnostic boundaries. The ABIDE dataset is preprocessed using multiple pipelines, including slice-timing and motion correction, nuisance signal regression, and spatial normalization. It supports whole-brain and regional analyses of intrinsic functional connectivity, offering insights into the dysconnectivity models of ASD.

In the fMRI retrieval task, we follow the methodology from⁴. Each test image is converted into a CLIP⁶⁰ image embedding. We compute cosine similarity between this embedding and its corresponding disjointed CLIP fMRI embedding, as well as 299 other randomly selected embeddings. Success is determined if the highest similarity is between the ground truth CLIP embedding and its respective fMRI embedding, indicating top-1 retrieval performance. This process is repeated 30 times to account for variability. For large-scale retrievals, we utilize the LAION-5B dataset, retrieving 16 candidate images via K-nearest neighbor lookup and selecting the best match based on cosine similarity.

For state classification, we adopt the experimental setup from⁵¹ to ensure a fair comparison. This involves ensuring that each data instance contains the same number of timesteps by extracting sets of k contiguous frames and looping the time series if fewer than k frames are present. During training, random frame sets are utilized, while validation and testing are conducted using the first k frames, where $k = 40$ in our experiments. We discuss the impact of different values of k on the experimental results in the appendix. Additionally, our NeuroSTORM framework is capable of directly processing complete tfMRI sequences.

Label Preprocessing. In our study, to provide appropriate supervision signals for models across various downstream tasks, label preprocessing is customized according to specific requirements. For classification tasks, including gender classification, disease diagnosis, and task fMRI state classification, we utilize one-hot encoding for the labels. This encoding converts categorical labels into a binary matrix representation, facilitating efficient model training and enhancing classification performance. For regression tasks, such as age prediction and phenotype prediction, we apply normalization^{24,25} to all target variables. This process involves scaling the data to a standard range, improving model convergence and stability. During inference, we can rescale the predicted values back to their original magnitude using the recorded normalization parameters. In the fMRI retrieval tasks, we assess retrieval accuracy by calculating the top-1, top-3, and top-5 retrieval rates^{41,43}. This involves determining the proportion of instances where the correct label is retrieved within the top k predictions, providing a comprehensive assessment of the model’s retrieval capabilities. This multifaceted approach to label preprocessing ensures that our models are well-equipped to handle the diverse challenges posed by both classification and regression tasks in neuroimaging analysis.

7 Competing interests

The authors declare no competing interests.