

FLORA: A Unified Generalist Model for Visual Brain Decoding via Multimodal

*Figure 1.* A comparative overview of visual decoding paradigms: (a) Single-modality decoding involves a neural encoder and a neural decoder, designed to process only one type of neural data at a time. (b) Multi-modality decoding employs distinct encoders tailored for various neural data types, with a unified decoder for subsequent tasks. (c) **Unified Multi-modality Model (Ours)**, featuring a shared, integrated encoding process that captures complementary representations across diverse neural data modalities. This architecture facilitates seamless interaction with multiple downstream tasks, embodying a **one-to-any** paradigm.

## Abstract

Decoding visual information from neural data using artificial intelligence enhances our understanding of the human visual system. However, simultaneously acquiring paired neural data across modalities is challenging, leading most existing approaches to process these signals independently. This neglects their complementary characteristics and hinders decoding performance. In this study, we introduce FLORA, an end-toend generalist model designed to integrate crossmodal neural data-including EEG, MEG, and fMRI-to construct a unified neural representation. FLORA employs multimodal large language models (MLLMs) alongside multimodal adapters and specialized diffusion model decoders, achieving superior performance on downstream tasks (e.g., neural signal retrieval and visual stimulus reconstruction) compared to single-modal approaches. By leveraging high-performance models, FLORA minimizes the number of parameters in the alignment and fusion layers, ensuring cost-effective fine-tuning. This design facilitates efficient training and seamless integration

of extra modalities and datasets. Our approach holds promise for advancing our understanding of the brain's visual mechanisms and fostering new insights within the cognitive science and brain-computer interface communities. Our code is available at https://anonymous.4open. science/r/FLORA-2C4A.

# 1 Introduction

The brain processes external stimuli by receiving and encoding information through the coordinated activity of largescale neural networks and intricate mechanisms. Gaining insight into these processes is essential, as it can reveal fundamental neural mechanisms and broaden potential applications in fields like brain-computer interfaces (BCIs). Recent advancements in computer science, particularly with generative models and pre-trained techniques, have driven substantial progress in decoding neural data from various neural modalities, including EEG (Li et al., 2024b; Kim et al., 2024), fMRI (Chen et al., 2023b; Caro et al., 2023), and MEG (Benchetrit et al., 2024; Cui et al., 2024). By leveraging the power of deep learning and unsupervised

algorithms, these methods are reshaping our capacity to interpret complex patterns of brain activity and enhancing the real-time decoding of neural signals.

One core challenge in achieving better decoding performance lies in establishing cross-modality representations for diverse neural data. Such data exhibit significant complexity and variability, influenced by factors like individual subjects, acquisition devices, visual content, and spatial-temporal dy-063 namics. Existing approaches often rely solely on a single-064 modality (Pan et al., 2024; Xia et al., 2024c) or alignment 065 among neural data, images, and text with contrastive learn-066 ing (Fu et al., 2024; Jiang et al., 2024; Lin et al., 2022). 067 These strategies either limit the decoding performance or 068 introduce unconfirmed prior knowledge. Recent works de-069 velop multi-encoder structures such as NeuroBind (Yang 070 et al., 2024a), which deploys distinct encoders tailored for specific types of neural data. Additionally, unified-encoder 072 architectures like UMBRAE (Xia et al., 2024b) have aimed to establish a uniform representation across different neural 074 data modalities. While these efforts demonstrate progress 075 in enhancing performance for various downstream tasks, they remain constrained by the inherent limitations of single 077 modalities and the lack of effective modal complementarity. 078 This challenge highlights the critical need for more cohesive 079 and innovative frameworks that not only integrate multiple modalities but also facilitate their unique characteristics to 081 improve overall decoding properties. 082

083 Inspired by the power of multimodal large language models (MLLMs) in processing the input information in a wider 085 variety of forms (Liu et al., 2024), we utilize it to integrate 086 multi-modal neural data, for enhancing the encoding and decoding performance of complex brain activity in a previ-087 088 ously unattainable way with single-modal data. In this work, 089 we develop FLORA, a multimodal large language model 090 that processes four modalities (i.e., EEG, MEG, fMRI, and 091 behavioral data) within a unified framework. As illustrated 092 in Fig. 1, FLORA comprises multimodal encoders, and a 093 trainable mixture of experts (MoE) to perfrom universal 094 projection across modalities. Unlike previous approaches, 095 the encoder in FLORA is trained on a diverse range of 096 multimodal data, including 16,540 paired EEG-image sam-097 ples (Grootswagers et al., 2022), 19848 paired MEG-image 098 samples, and 8540 paired fMRI-image samples (Hebart 099 et al., 2023), ensuring shared representations across all 100 modalities. To evaluate the FLORA model, we apply it for downstream tasks including retrieval, caption, and reconstruction compared to single-modality models. Through solid experiments, we find that FLORA outperforms main-104 stream single-modality-based decoder frameworks in all 105 downstream tasks, benefiting from its unified training strat-106 egy. The unification in FLORA advances neuroscientific insights into visual representation in the brain and broadens practical Brain-Computer Interface (BCI) applications. 109

Our work has three main contributions.

- We propose a novel unified encoder to capture diverse neural data embeddings, realizing feature complementarity across heterogeneous neural modalities.
- Our framework, FLORA, enables seamless alignment and integration with various downstream tasks, significantly enhancing its adaptability and achieving SOTA performance across joint-subjects visual decoding benchmarks.
- To the best of our knowledge, **FLORA** is the first work to concurrently unify multiple modalities (EEG, MEG and fMRI) and task paradigms, providing a cohesive computational foundation for non-invasive neural decoding and BCI applications.

#### **Related work** 2

#### 2.1 Visual decoding via uni-modal neural data

Significant progress has been made in visual decoding through single modalities such as EEG, fMRI, and MEG. Early research predominantly focused on leveraging deep learning techniques to classify (Zheng & Chen, 2021; Spampinato et al., 2019) and retrieve (Ye et al., 2022) visual stimuli based on neural recordings. In recent years, visual decoding has achieved unprecedented quality (Sun et al., 2023; Zhou et al., 2024; Benchetrit et al., 2024; Ma et al., 2024), driven by advancements in generative models (Ho et al., 2020; Rombach et al., 2022; Liu et al., 2022) and LLMs (Touvron et al., 2023; Brown et al., 2020). For instance, RealMind (Li et al., 2024a) utilizes multimodal models that combine the ability of semantic and geometric information, leading to improved decoding performance. Additionally, (Scotti et al., 2024b) demonstrates high-quality cross-subject reconstruction and retrieval capabilities, requiring minimal data, particularly valuable in scenarios with limited sample sizes. Moreover, (Benchetrit et al., 2024) successfully achieves real-time level retrieval and generation by capitalizing on the high temporal resolution of MEG, thereby enhancing the decoding of dynamic visual perception. These advancements underscore the practicality and effectiveness of utilizing artificial intelligence to gain insights into visual information processing and the functioning of the visual system.

#### 2.2 Multimodal Large Models

Multimodal Large Language Models (MLLMs) have transformed AI by enabling seamless integration and generation of information across diverse modalities, including text, images, audio, video, and 3D objects. Foundational architectures like single-tower models (e.g., ViT (Dosovitskiy



128 Figure 2. Overall Architecture of FLORA. FLORA comprises neural modality encoders, a universal projection, and a Mixture of 129 Experts (MoE) projection module, with separate outputs for each modality. Left: The modality encoders, implemented as a time series 130 backbone network with spatial-temporal convolutional layer, transforms the input signal into a sequence of tokens. The universal encoder is designed to align high-dimensional feature representations from heterogeneous data. Middle: The MoE projection module projects and 131 aligns diverse neural modalities into a unified semantic representation space. Right: Various task heads facilitate different downstream 132 tasks such as retrieval, captioning, and reconstruction. All modules are jointly trained during the same stage, optimizing computational 133 efficiency. 134

136 et al., 2021)) and multi-tower models (e.g., CLIP (Radford 137 et al., 2021)) pioneered image-text alignment, establishing 138 the basis for multimodal learning. Building on these innova-139 tions, models like BLIP (Li et al., 2022), BLIP-2 (Li et al., 140 2023), Qwen-VL (Bai et al., 2023), and DALL·E (Ramesh 141 et al., 2022) have significantly improved cross-modal under-142 standing and generative capabilities. Inspired by the success 143 of vision-language models, recent research has expanded 144 to include audio (e.g., Tacotron (Shen et al., 2018), Nat-145 uralSpeech (Tan et al., 2022)), video (e.g., Create-to-Tell 146 (Pan et al., 2018), DreamVideo (Wang et al., 2024a)), and 147 3D objects (e.g., DreamFusion (Poole et al., 2022)). Mod-148 els like NExT-GPT (Wu et al., 2024) exemplify versatile 149 multimodal capabilities, seamlessly interacting across text, 150 images, audio, and video using distinct encoders tailored 151 to various downstream tasks. Similarly, Meta-Transformer 152 (Zhang et al., 2023) utilizes a unified parameter set to en-153 code data from multiple modalities, effectively supporting a 154 wide range of task types. 155

#### 3 Method

156

157

161

163

164

135

158 **Formulation.** Let T represent the length of neural data, 159 C the number of channels, and N the total number of 160 data samples. Our objective is to derive neural embeddings  $Z = \mathcal{E}(X) \in \mathbb{R}^{N \times F}$  from the brain activity data 162  $X \in \mathbb{R}^{N \times C \times T}$ , where  $\mathcal{E}$  is the unified encoder and Fis the projection dimension of the embeddings. Concurrently, we use the CLIP model to extract image embeddings  $\hat{Z} \in \mathbb{R}^{N \times F}$  from images I. Our goal is to effectively align the multimodal representation with the image representation, as illustrated in Fig. 2.

#### 3.1 Model Architecture

Our method aims to process the combination of multiple modalities of neural data via one unified model. All crossmodal neural features are efficiently aligned through the MoE module, maximizing complementary strengths across modalities. Multi-stage training strategies are applied to the encoding and alignment procedure. In the training phase, encoders are trained with neural data and image pairs using a contrastive learning framework. Then each modality is aligned by entering the Unified Projector based on each pretrained encoder. In the inference phase, neural embeddings from the Unified Projector can be used for a variety of zeroshot tasks, including EEG/MEG/fMRI image classification, retrieval, and image reconstruction.

#### 3.1.1 NEURAL FEATURE EXTRACTION MODULE

This module is mainly employed to address the cross-subject challenges in multimodal neural data decoding as shown in Fig. 3. The core difficulty arises from the heterogeneous neural responses of different subjects to identical visual stimuli, leading to substantial inter-subject variability in



214

215

165

167

168 169



Figure 3. Architecture of the Neural Encoder. The original neural sequences from multiple variates are simultaneously embedded into tokens. Multi-granularity attention is applied to the embedded tokens of correlated variables, enhancing interpretability and revealing electrode-level correlations. The representations for each token are then assigned through the router layers. Subsequently, Temporal-Spatial convolution is employed to mitigate overfitting while improving the model's capacity for temporal-spatial representation learning.

the data. By incorporating multi-granular approaches, as 188 proposed in (Wang et al., 2024c), into the encoder architec-189 ture, inter-subject heterogeneity can be mitigated, thereby 190 facilitating the learning of shared joint embeddings across 191 subjects.

For input neural signal  $X \in \mathbb{R}^{T \times C}$ , the architecture first 193 performs Multi-Granularity Time Patching with exponentially increasing patch lengths  $\{2^1, 2^2, ..., 2^n\}$ :  $x_p^{(i)} \in$ 195 196  $\mathbb{R}^{N_i \times (2^i \cdot C)}, N_i = \lceil T/2^i \rceil$ . Then we get a group tokens  $x^{(i)} \in \mathbb{R}^{N_i \times D}$  by: 197 198

$$x^{(i)} = x_p^{(i)} W^{(i)} + W_{pos}[1:N_i] + W_{gr}^{(i)}$$
(1)

where  $W_{\text{pos}}[1:N_i] \in \mathbb{R}^{N_i \times D}$  is the first  $N_i$  rows of the 201 202 positional embedding  $W_{\rm pos}$ , and a learnable granularity em-203 bedding  $W_{\rm gr}^{(i)} \in \mathbb{R}^{1 \times D}$ . And the router is used in the multi-204 granularity self-attention (as described later) for each granularity: 206

$$u^{(i)} = W_{pos}[N_i + 1] + W_{gr}^{(i)}$$
(2)

The embeddings undergo a two-stage *Multi-Granularity* 208 Attention mechanism. In the first stage, intra-granularity 209 attention captures temporal dependencies within each scale 210 by concatenating patch embeddings with their correspond-211 ing router. The intermediate sequence of embeddings 212  $z^{(i)} \in \mathbb{R}^{(N_i+1) \times D}$  is formed as: 213

$$z^{(i)} = \operatorname{Concat}(x^{(i)}, u^{(i)}) \tag{3}$$

216 where  $x^{(i)} \in \mathbb{R}^{N_i \times D}$  represents patch embeddings at scale 217 i, and  $u^{(i)} \in \mathbb{R}^{1 \times D}$  is the router embedding. Self-attention 218 is applied on each  $z^{(i)}$  to update both  $x^{(i)}$  and  $u^{(i)}$ . The 219

collection of all scale-specific representations is:

2

$$Z = \text{Concat}(z^{(1)}, z^{(2)}, ..., z^{(n)})$$
(4)

Then, inter-granularity attention facilitates information exchange across scales through router interactions. All router embeddings are aggregated:

$$U = \operatorname{Concat}(u^{(1)}, u^{(2)}, ..., u^{(n)}) \in \mathbb{R}^{n \times D}$$
(5)

These router embeddings serve as query vectors to attend to features across scales. The attended features are concatenated as:

$$H_{\text{attn}} = \text{Concat}(x^{(1)}, x^{(2)}, ..., x^{(n)}) \in \mathbb{R}^{(\sum_{i=1}^{n} N_i) \times D}$$
(6)

Finally,  $H_{\text{attn}}$  is processed through a temporal-spatial convolution module to obtain  $H_{\text{conv}}$ , followed by an MLP to get the final encoding  $Y \in \mathbb{R}^D$ , which aligns with CLIP embeddings of the corresponding visual stimuli image.

### 3.1.2 UNIFIED PROJECTION MODULE

Aligning the multimodal neural representation space with the visual embedding space of the CLIP model is a crucial step, particularly given the complexity and heterogeneity of the multimodal neural imaging-image pairs. We employ a MoE module (Zhu et al., 2024), which has K projection experts  $E_1, E_2, \ldots, E_K$ , where each expert is a two-layer MLP, and a soft router  $R_{\text{soft}}$  to control the contribution of each expert.

Using a small network q, we use the soft router computes a score vector for each token, with a length equal to the number of experts K. Finally, the Softmax function is applied to obtain the final routing vector, yielding the final routing vector. It can be formulated as:

$$R_{\text{soft}}(x_i) = \frac{Sigmoid(g(x_i))}{\sum Sigmoid(g(x_i))}$$
(7)

To alleviate the heterogeneity of different neural modalities.  $R_{\text{soft}}$  is a lightweight MLP designed to receive the inputs of different encoders, and calculate the routing weights  $W_{\text{soft}} \in \mathbb{R}^{N \times F \times K}$  of each expert for each neural token, which can be formulated as:

$$W_{\text{soft}}(Z) = \sigma \cdot R_{\text{soft}}(Z) \tag{8}$$

where  $\sigma$  is the SoftMax function. Then we can obtain aligned neural tokens  $Z \in \mathbb{R}^{N \times F}$  through a weighted sum of all experts' output as follows:

$$Z = \sum_{k=1}^{K} W_{\text{soft},k} \cdot E_k(Z) \tag{9}$$

where  $w_{\text{soft},k}$  denotes the routing weight of the k-th projection expert.

### 3.2 Training Objective

We applied multiple approaches to loss functions, serving various downstream tasks. For retrieval tasks, we utilize CLIP Loss for contrastive learning(Radford et al., 2021). By aligning fused data F with corresponding image data I, FLORA could utilize potential of each modalities and achieve the identification of images seen by subjects with given neural data. Training consists of three progressive stages: visual-neural Data contrast training, visual-neural generation training, and supervised fine-tuning through Universal Projection.

#### 3.2.1 HIGH-LEVEL VISUAL MODELING

This pipline is primarily used to learn high level unified representation from the unified encoder. We use the contrastive learning functions inspired by (Scotti et al., 2024a) to train both the multi-tower encoder and *Unified Projection*:

$$\mathcal{L}_{\text{SoftCLIP}} = -\sum_{i,j=1}^{N} \frac{\exp\left(\frac{t_i \cdot t_j}{\tau}\right)}{Z_i} \cdot \log\left(\frac{\exp\left(\frac{p_i \cdot t_j}{\tau}\right)}{Z_i^{(p)}}\right)$$
(10)
$$Z_i = \sum_{m=1}^{N} \exp\left(\frac{t_i \cdot t_m}{\tau}\right), \quad Z_i^{(p)} = \sum_{m=1}^{N} \exp\left(\frac{p_i \cdot t_m}{\tau}\right)$$
(11)

Where  $p_i = \mathcal{E}(x_i), t_i = \text{CLIP}_{\text{Image}}(y_i)$ . For generative tasks, including captioning and reconstruction, we incorporate Mean Squared Error (MSE) loss to ensure consistent

model training for regression-based tasks. As a result, the overall loss function of FLORA is formulated as a composite of these individual loss components. The optimization of model parameters  $\theta$  is then governed by the following objective function:

$$\min_{\boldsymbol{\alpha}} \mathcal{L} = \mathcal{L}_{\text{BiMixCo|SoftCLIP}} + \boldsymbol{\alpha} \cdot \mathcal{L}_{\text{MSE}} + \boldsymbol{\beta} \cdot \mathcal{L}_{\text{prior}} \quad (12)$$

where  $\alpha$  and  $\beta$  are hyperparameters to control the balance of each loss function.

#### 3.2.2 LOW-LEVEL VISUAL MODELING

Previous works (Scotti et al., 2024b; Li et al., 2024b; Zhang et al., 2024) have extensively used prior diffusion to enhance the transformation of neural data embeddings to image embedding distributions. We continue this approach with prior diffusion, but use the neural embeddings directly output by the *Unified Projection* module for retrieval evaluation—rather than after image prior enhancement, because this can more objectively measure the degree of alignment of neural data with the neural manifold of the image without over-reliance on image priors.

To train the low-level pipline, the original unified encoder goes through additional CNN upsampler to get  $\hat{f}$ , the approximation of original f. Then a new image  $\hat{I}$  is reconstructed using the Stable Diffusion (Rombach et al., 2021) VAE decoder  $D(\cdot)$  given  $\hat{f}$ , and a compound loss  $\mathcal{L}$  is minimized:

$$\mathcal{L}_{\text{lowlevel}} = \|D(f) - D(\hat{f})\| + \|f - \hat{f}\| + \lambda_P \mathcal{L}_P(D(\hat{f}))$$
(13)

where  $\mathcal{L}_P(\cdot)$  is a perceptual loss such as LPIPS (Zhang et al., 2018), and  $\lambda_P$  is loss weights.

## 4 **Experiment**

### 4.1 Dataset and Implementation

**Dataset.** We jointly trained FLORA with the THINGS-EEG2 (Gifford et al., 2022), THINGS-MEG, and THINGSfMRI (Hebart et al., 2023) datasets to achieve a unified semantic representation. Notably, we noticed that relatively few decoding studies (Xue et al., 2024) have utilized the THINGS-EEG1 dataset (Grootswagers et al., 2022), which might be due to its low signal-to-noise ratio. This limitation is primarily attributed to two methodological factors: (i) the short presentation time of visual stimulus (50ms), and (ii) the lack of trial repetition (each stimulus viewed only once by each participant). Finally, we explored fine-tuning experiments using THINGS EEG1 in the Appendix B. The THINGS-EEG2 training set consists of 1654 concepts, each associated with 10 images, with 4 repetitions per image, while the testing set includes 200 concepts, each represented

275 Table 1. Multimodal Retrieval Performance. We present the retrieval performance of various methods on the THINGS-EEG2, THINGS-276 MEG and THINGS-fMRI datasets. Retrieval refers to the task of correctly identifying paired neural latent representations from the 200 or 277 100 class test set clip embeddings. Performance is evaluated across joint subject retrieval tasks, and we report results for several retrieval 278 configurations: 2-way, 4-way, 10-way, as well as Top-1 and Top-5 accuracy for the 200-way retrieval task (100 for fMRI). FLORA is 279 compared with previous SOTA methods, including NICE (Song et al., 2023), EEGNetV4 (Lawhern et al., 2018), B.D. (Benchetrit et al., 280 2024), MindEyeV2 (Scotti et al., 2024b), UMBRAE (Xia et al., 2024a), ATM-S (Li et al., 2024b), CongnitionCapturer (Zhang et al.,

281	2024), MindBridge	(Wang	g et al.,	2024b),	WAVE (	Wang et a	al., 202	24d) and	d MB2C	C (Wei et	al., 2024	).In sul	oject re	sults ref	ers to Ap	pendix F
282	-		Т	HINGS	EEG2			THINGS-MEG					THINGS-fMRI			
283	Model	2-way	4-way	10-way	200-way	200-way	2-way	4-way	10-way	200-way	200-way	2-way	4-way	10-way	100-way	100-way
205		Top-1	Top-1	Top-1	Top-1	Top-5	Top-1	Top-1	Top-1	Top-1	Top-5	Top-1	Top-1	Top-1	Top-1	Top-5
284	NICE	90.10	77.15	57.15	13.80	35.80	60.62	39.62	20.87	2.50	7.88	87.00	74.67	55.00	22.00	51.00
285	EEGNetV4	88.15	77.05	59.10	12.15	35.85	71.13	48.63	29.38	3.88	11.88	86.00	74.67	55.33	18.00	44.67
286	B.D.	90.05	76.85	55.85	12.45	34.15	62.12	35.50	17.63	1.88	6.50	72.67	58.33	38.00	6.67	26.00
287	MindEyeV2	90.25	79.85	64.30	17.35	43.65	64.25	40.25	20.37	2.13	8.38	92.00	80.33	64.33	21.33	<u>60.33</u>
207	UMBRAE	66.50	41.95	21.85	2.90	10.00	61.75	35.50	16.63	2.00	6.50	85.67	70.67	50.67	16.67	42.67
288	ATM-S	93.35	83.70	68.60	19.70	48.60	77.00	59.25	38.88	6.25	20.12	92.67	81.67	65.33	24.00	55.33
289	CongnitionCapturer	92.90	84.00	67.85	20.50	50.70	67.88	47.63	23.57	2.38	10.37	86.67	73.00	55.33	20.67	48.00
290	MindBridge	89.10	75.35	57.45	15.00	39.15	64.88	40.12	20.13	1.87	7.37	88.33	73.33	55.67	16.67	54.33
201	WAVE	89.40	75.85	58.80	15.00	39.05	67.50	44.75	24.00	3.00	11.38	88.67	71.67	52.33	20.33	45.33
291	MB2C	85.20	70.25	49.65	10.20	29.85	61.75	33.25	17.00	1.25	5.50	84.67	69.00	49.67	18.67	42.33
292	FLORA-unimodal	95.55	86.90	73.45	25.35	57.30	81.75	64.50	46.62	8.00	24.38	91.33	79.00	63.00	26.33	57.33
293	FLORA-multimodal	<u>94.05</u>	<u>87.30</u>	<u>73.15</u>	<u>25.05</u>	<u>56.35</u>	80.50	<u>61.88</u>	<u>39.75</u>	<u>6.88</u>	<u>23.38</u>	<u>92.33</u>	84.67	70.67	28.33	63.33

295 by a single image and 80 repeated trials. The original train-296 ing set for THINGS-MEG contains 1854 concepts, each 297 with 12 images (viewed once), and the testing set selects one unseen image from 200 of these concepts, repeated 12 299 times. The THINGS-fMRI training set comprises 720 con-300 cepts, each associated with 12 images (viewed once), and 301 the testing set includes 100 concepts, with a single image 302 selected from each and repeated 12 times. To facilitate zero-303 shot evaluation, we reclassified the data during experiments 304 (Appendix D for more details). 305

306 Implementation and training setup. All experiments, 307 including both training and inference procedures across 308 modalities, were conducted using two NVIDIA A100-80GB 309 GPUs. For visual encoding, we employed CLIP (ViT-L-14) 310 to generate image embeddings. Each modality encoder 311 was trained on the original training set of the THINGS 312 dataset for 50 epochs, with a learning rate of 3e-4 and a 313 batch size of 360, utilizing the AdamW optimizer. During 314 the forward pass, batches from different modalities were 315 processed sequentially, with each batch maintaining a size 316 equal to the specified batch size. Hyperparameters were 317 shared across all modalities. 318

#### 4.2 Results on Retrieval

294

319

320

321

322

323

324

325

327

328

329

Multimodal retrieval Image retrieval metrics are used to quantify the extent of fine-grained image information captured in neural embeddings. The effectiveness of the image retrieval task is evaluated by computing the cosine similarity between embeddings derived from EEG, MEG, and fMRI signals, and the CLIP embeddings for a test set consisting of images. Tab. 1 presents the average retrieval performance across all subjects for EEG, MEG and fMRI.

#### 4.3 **Results on Reconstruction**

Multimodal reconstruction Previous studies have made significant advancements in the task of reconstructing images from fMRI data (Takagi & Nishimoto, 2023; Chen et al., 2023a; Scotti et al., 2024a;b). The SDXL (Podell et al., 2023) + ip-adapter framework is utilized in the reconstruction pipline. Our approach builds upon the framework of MindEye2 (Scotti et al., 2024b), integrating state-of-theart models such as SDXL (Podell et al., 2023) and FLUX for image reconstruction. What distinguishes our work is the ability to leverage the complementary strengths of the unified encoder, enabling the reconstruction of images from latent representations of multimodal neural data across different modalities using a single model. We employ a lightweight UNet as a prior diffusion denoising network to map the embeddings into the image space across different neural modalities. The cross-modal reconstruction performance results are presented in Tab. 2.

### 4.4 Results on Captioning

Multimodal caption Prior work (Li et al., 2024a) has demonstrated that directly reconstructing images from latent representations yields superior performance compared to first reconstructing the images and then generating captions. Building on this insight, we further employ Prior Diffusion to modulate the distribution of model features, thereby enhancing caption generation performance. Subsequently, we evaluated the captioning performance of our model by inputting the neural latents from the unified encoder to OpenFlamingo (Awadalla et al., 2023). The results are presented in Tab. 3.

*Table 2.* **Multimodal Visual Reconstruction.** Quantitative evaluation of the reconstruction quality for EEG, MEG and fMRI from all subjects. Comparative methods include MindEye (Scotti et al., 2024a), BrainDiffuser (Ozcelik & VanRullen, 2023), B.D. (Benchetrit Y, 2023), MindEye2(Scotti et al., 2024b), UMBRAE (Xia et al., 2024a) and ATM-S (Li et al., 2024b). It is noteworthy that, in contrast to other THINGS dataset-based works that conduct training and inference within single subjects, FLORA implements training and reconstruction in a cross-subject and cross-modality setting. For detailed descriptions of the evaluation details, please refer to the Appendix D.

Dataset	Method		Low-lev	el	Н			
Dutuset		PixCorr ↑	SSIM↑	AlexNet(2)↑	$\overline{\text{AlexNet}(5)}\uparrow$	Inception <sup>↑</sup>	CLIP↑	SwAV↓
	MindEye	0.130	0.308	0.917	0.974	0.936	0.942	0.369
	BrainDiffuser	0.254	0.356	0.942	0.962	0.872	0.915	0.423
NSD-fMRI	B.D.	0.305	0.366	0.962	0.977	0.910	0.917	0.410
	MindEye2	0.322	0.431	0.961	0.986	0.954	0.930	0.344
	UMBRAE	0.283	0.341	0.955	0.970	0.917	0.935	0.393
THINGS-fMRI	FLORA-multimodal	0.090	0.362	0.535	0.548	0.551	0.624	0.676
	B.D.	0.058	0.327	0.695	0.753	0.593	0.700	0.630
THINCS MEC	B.D. (averaged)	0.076	0.336	0.736	0.826	0.671	0.767	0.584
THINGS-MEG	ATM-S	0.104	0.340	0.613	0.672	0.619	0.603	0.651
	FLORA-multimodal	0.070	0.291	0.553	0.613	0.567	0.619	0.663
TUINCS EEC2	ATM-S	0.160	0.345	0.776	0.866	0.734	0.786	0.582
THINGS-EEG2	FLORA-multimodal	0.117	0.362	0.702	0.772	0.666	0.725	0.604

*Table 3.* **Multimodal Brain Captioning.** We compare the performance of various methods for brain captioning. To more clearly illustrate the superiority of our approach, the baseline method was trained and evaluated on the NSD dataset, while our method was trained and evaluated on the THINGS-fMRI dataset, in cross-subject settings.

Method	BLEU1 $\uparrow$	$\textbf{BLEU2} \uparrow$	BLEU3 $\uparrow$	BLEU4 $\uparrow$	METEOR $\uparrow$	$\textbf{ROUGE} \uparrow$	$\textbf{CIDEr} \uparrow$	SPICE $\uparrow$	CLIP-S↑	<b>RefCLIP-S</b> ↑
SDRecon (Takagi & Nishimoto, 2023)	36.21	17.11	7.72	3.43	10.03	25.13	13.83	5.02	61.07	66.36
OneLLM (Han et al., 2024)	47.04	26.97	15.49	9.51	13.55	35.05	22.99	6.26	54.80	61.28
UniBrain (Mai & Zhang, 2023)	-	-	-	-	16.90	22.20	-	-	-	-
BrainCap (Ferrante et al., 2023)	55.96	36.21	22.70	14.51	16.68	40.69	41.30	9.06	64.31	69.90
UMBRAE (Xia et al., 2024a)	57.84	38.43	25.41	17.17	18.70	42.14	53.87	12.27	66.10	72.33
FLORA-unimodal w/o prior fMRI*	36.71	30.58	25.36	18.80	41.14	40.52	4.46	16.19	17.55	73.37
FLORA-unimodal w/ prior fMRI*	33.73	27.77	22.41	16.05	41.73	39.21	4.74	15.05	17.71	59.81

<sup>\*</sup> The evaluation metrics for FLORA were computed using the **nltk** evaluation framework.

#### 4.5 In-depth Analysis

Visual Concept Localization We adopted a methodology similar to that presented in (Shen et al., 2024) for localizing semantic concepts within neural signals. Specifically, we employed Honeybee (Cha et al., 2024) to extract the target concepts from natural language with the prompt "Describe the main concept in the picture in three words". These concepts were encoded by the CLIP ViT-L-14 text encoder and used as target representations for integrated Grad-CAM, facilitating the spatial localization of these concepts within brain signals. We trained a unified encoder incorporating all three input modalities. The final layers of this model were leveraged to extract semantic features. Fig. 4 visualizes 376 the localization outcomes on the THINGS-fMRI dataset. showing the discrimination of various semantic information 378 379 within brain signals in response to novel visual stimuli.

We conducted ablation studies on the localized semantic concepts. After identifying the relevant concepts within the original fMRI voxels from visual cortex, we selectively masked the corresponding voxels. Feature extraction and

visual reconstruction were subsequently performed using these modified brain signals. In Fig. 5, the targeted removal of fMRI voxels from specific brain regions associated with particular semantic concepts led to the exclusion of these semantics in the reconstructed visual representation.

## 5 Discussion and Conclusion

Here, we present FLORA, a unified generalist model that employs a shared encoder to extract cohesive neural representations across three distinct modalities: EEG, MEG, and fMRI. To the best of our knowledge, this represents the first work to align visual representations across multiple neural modalities using a unified encoder. FLORA achieves substantial performance gains across a range of downstream neural decoding tasks, representing a significant step forward in the development of a unified framework for cross-modal visual neural decoding.

Our framework has several limitations that warrant further investigation. First, the samples from different neural data modalities are imbalanced. Such imbalanced distri-



*Figure 4.* Gradient heatmaps of brain activity generated by FLORA. It shows the neural representation of different semantic information in the brain in response to the same visual stimulus, such as the semantics 'green' and 'tall' from the bamboo image.



*Figure 5.* Concept analysis on THINGS-fMRI using FLORA through semantic signal nullification. We show the impacts on visual reconstruction by ablating certain concepts from visual stimuli.

420 butions of samples across modalities and classes pose a 421 challenge for the unified encoder in effectively capturing 422 modality-specific representations, thereby constraining the 423 overall performance of the framework. To mitigate this is-424 sue, one potential solution could involve pre-training the 425 modality-specific encoders on larger and more balanced neu-426 ral datasets, which may enhance their ability to generalize 427 across modalities. Second, due to the shared parameters 428 and representation spaces, our framework underperforms 429 compared to some SOTA methods in certain single-modality 430 downstream tasks, such as reconstruction and caption gener-431 ation. This performance gap could potentially be addressed 432 by leveraging larger datasets and more sophisticated model 433 architectures. Finally, inherent limitations of THINGS se-434 ries datasets, including their limited size and repetitive na-435 ture, restrict FLORA's ability to enforce consistent category 436 distributions between training and test sets across modalities. 437 For example, the fMRI training and test sets exhibit an inclu-438 sion relationship in their category sets, which may introduce

biases. To overcome these challenges, larger-scale, crossmodal datasets with more consistent category distributions are of high need.

Despite these limitations, FLORA demonstrates the potential of multimodal grand unified neural decoding by effectively leveraging the complementary information from different modalities to enhance performance across various tasks. Future work may prioritize two key directions: (i) the utilization of larger neural datasets, both paired and unpaired, to train robust multimodal foundation models, and (ii) the exploration of cutting-edge unified model architectures that incorporate advanced components for representation, alignment, and fusion stages, tailored to diverse neural decoding tasks. With such a unified framework, beyond the visual concept analysis (Fig. 4), more scientific questions can be studied, such as how these semantic representations across modalities are unified and adapted to personal experiences and knowledge.

404

405

406

412 413 414

415 416

417

418 419

### 440 **References**

441

484

- Alayrac, J.-B., Donahue, J., Luc, P., Miech, A., Barr, I., Hasson, Y., Lenc, K., Mensch, A., Millican, K., Reynolds, M., et al. Flamingo: a visual language model for fewshot learning. *Advances in neural information processing systems*, 35:23716–23736, 2022.
- Awadalla, A., Gao, I., Gardner, J., Hessel, J., Hanafy, Y.,
  Zhu, W., Marathe, K., Bitton, Y., Gadre, S., Sagawa,
  S., et al. Openflamingo: An open-source framework
  for training large autoregressive vision-language models. *arXiv preprint arXiv:2308.01390*, 2023.
- Bai, J., Bai, S., Yang, S., Wang, S., Tan, S., Wang, P.,
  Lin, J., Zhou, C., and Zhou, J. Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond, 2023. URL https: //arxiv.org/abs/2308.12966.
- 459 Benchetrit, Y., Banville, H., and King, J.-R. Brain decoding:
  460 toward real-time reconstruction of visual perception. In
  461 *The Twelfth International Conference on Learning Repre-*462 *sentations*, 2024.
- Benchetrit Y, Banville H, K. J. R. Brain decoding: toward
  real-time reconstruction of visual perception. *arXiv*, 2310:
  19812, 2023.
- 467 Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, 468 J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., 469 Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., 470 Henighan, T., Child, R., Ramesh, A., Ziegler, D. M., Wu, 471 J., Winter, C., Hesse, C., Chen, M., Sigler, E., Litwin, M., 472 Gray, S., Chess, B., Clark, J., Berner, C., McCandlish, 473 S., Radford, A., Sutskever, I., and Amodei, D. Language 474 models are few-shot learners, 2020. URL https:// 475 arxiv.org/abs/2005.14165. 476
- 477 Caro, J. O., de O. Fonseca, A. H., Averill, C., Rizvi, S. A., 478 Rosati, M., Cross, J. L., Mittal, P., Zappala, E., Levine, D., Dhodapkar, R. M., Abdallah, C. G., and van Dijk, D. 480 Brainlm: A foundation model for brain activity record-481 ings. *bioRxiv*, 2023. doi: 10.1101/2023.09.12.557460.
  482 URL https://www.biorxiv.org/content/ 483 early/2023/09/13/2023.09.12.557460.
- Cha, J., Kang, W., Mun, J., and Roh, B. Honeybee: Localityenhanced projector for multimodal llm. In *Proceedings*of the IEEE/CVF Conference on Computer Vision and
  Pattern Recognition, pp. 13817–13827, 2024.
- Chen, Z., Qing, J., Xiang, T., Yue, W. L., and 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*, pp. 22710–22720, 2023a.

- Chen, Z., Qing, J., Xiang, T., Yue, W. L., and Zhou, J. H. Seeing beyond the brain: Conditional diffusion model with sparse masked modeling for vision decoding, 2023b. URL https://arxiv.org/abs/2211.06956.
- Cui, W., Jeong, W., Thölke, P., Medani, T., Jerbi, K., Joshi, A. A., and Leahy, R. M. Neuro-gpt: Towards a foundation model for eeg, 2024. URL https://arxiv.org/ abs/2311.03764.
- Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., and Houlsby, N. An image is worth 16x16 words: Transformers for image recognition at scale, 2021. URL https://arxiv.org/abs/2010.11929.
- Ferrante, M., Ozcelik, F., Boccato, T., VanRullen, R., and Toschi, N. Brain captioning: Decoding human brain activity into images and text. arXiv preprint arXiv:2305.11560, 2023.
- Ferrante, M., Boccato, T., and Toschi, N. Towards neural foundation models for vision: Aligning eeg, meg and fmri representations to perform decoding, encoding and modality conversion. In *ICLR 2024 Workshop on Repre*sentational Alignment, 2024.
- Fu, H., Shen, Z., Chin, J. J., and Wang, H. Brainvis: Exploring the bridge between brain and visual signals via image reconstruction, 2024. URL https://arxiv.org/abs/2312.14871.
- Gifford, A. T., Dwivedi, K., Roig, G., and Cichy, R. M. A large and rich eeg dataset for modeling human visual object recognition. *NeuroImage*, 264:119754, 2022.
- Grootswagers, T., Zhou, I., Robinson, A. K., Hebart, M. N., and Carlson, T. A. Human EEG recordings for 1,854 concepts presented in rapid serial visual presentation streams. *Sci Data*, 9(1):3, January 2022. ISSN 2052-4463. doi: 10.1038/ s41597-021-01102-7. URL https://www.nature. com/articles/s41597-021-01102-7.
- Han, J., Gong, K., Zhang, Y., Wang, J., Zhang, K., Lin, D., Qiao, Y., Gao, P., and Yue, X. Onellm: One framework to align all modalities with language. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 26584–26595, 2024.
- Hebart, M. N., Contier, O., Teichmann, L., Rockter, A. H., Zheng, C. Y., Kidder, A., Corriveau, A., Vaziri-Pashkam, M., and Baker, C. I. Things-data, a multimodal collection of large-scale datasets for investigating object representations in human brain and behavior. *Elife*, 12:e82580, 2023.

- Ho, J., Jain, A., and Abbeel, P. Denoising diffusion probabilistic models, 2020. URL https://arxiv.org/
  abs/2006.11239.
- Jiang, S., Meng, Z., Liu, D., Li, H., Su, F., and Zhao, Z.
  Mindshot: Brain decoding framework using only one image, 2024. URL https://arxiv.org/abs/2405.
  15278.
- Kim, J. W., Alaa, A., and Bernardo, D. Eeg-gpt: Exploring capabilities of large language models for eeg classification and interpretation, 2024. URL https://arxiv.org/abs/2401.18006.
- Lawhern, V. J., Solon, A. J., Waytowich, N. R., et al. Eegnet:
  a compact convolutional neural network for eeg-based
  brain–computer interfaces. *Journal of Neural Engineering*, 15(5):056013, 2018.
- 513 Li, D., Qin, H., Wu, M., Cao, Y., Wei, C., and Liu, Q.
  514 Realmind: Zero-shot eeg-based visual decoding and
  515 captioning using multi-modal models, 2024a. URL
  516 https://arxiv.org/abs/2410.23754.

533

534

535

536

Li, D., Wei, C., Li, S., Zou, J., and Liu, Q. Visual decoding and reconstruction via eeg embeddings with guided
diffusion. *arXiv preprint arXiv:2403.07721*, 2024b.

Li, J., Li, D., Xiong, C., and Hoi, S. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International conference on machine learning*, pp. 12888–12900. PMLR, 2022.

- Li, J., Li, D., Savarese, S., and Hoi, S. Blip-2: bootstrapping language-image pre-training with frozen image encoders and large language models. In *Proceedings of the* 40th International Conference on Machine Learning, pp. 19730–19742, 2023.
  - Lin, S., Sprague, T., and Singh, A. K. Mind reader: Reconstructing complex images from brain activities, 2022. URL https://arxiv.org/abs/2210.01769.
- Liu, X., Gong, C., and Liu, Q. Flow straight and fast: Learning to generate and transfer data with rectified flow, 2022.
  URL https://arxiv.org/abs/2209.03003.
- Liu, Z., Dong, Y., Liu, Z., Hu, W., Lu, J., and Rao, Y.
  Oryx mllm: On-demand spatial-temporal understanding at arbitrary resolution, 2024. URL https://arxiv. org/abs/2409.12961.
- Ma, S., Wang, L., Hou, S., and Yan, B. Aligned with llm: a new multi-modal training paradigm for encoding fmri activity in visual cortex, 2024. URL https://arxiv. org/abs/2401.03851.

- Mai, W. and Zhang, Z. Unibrain: Unify image reconstruction and captioning all in one diffusion model from human brain activity. *arXiv preprint arXiv:2308.07428*, 2023.
- Ozcelik, F. and VanRullen, R. Natural scene reconstruction from fmri signals using generative latent diffusion. *Scientific Reports*, 13(1):15666, 2023.
- Pan, H., Li, Z., Fu, Y., Qin, X., and Hu, J. Reconstructing visual stimulus images from eeg signals based on deep visual representation model, 2024. URL https:// arxiv.org/abs/2403.06532.
- Pan, Y., Qiu, Z., Yao, T., Li, H., and Mei, T. To create what you tell: Generating videos from captions, 2018. URL https://arxiv.org/abs/1804.08264.
- Podell, D., English, Z., Lacey, K., Blattmann, A., Dockhorn, T., Müller, J., Penna, J., and Rombach, R. Sdxl: Improving latent diffusion models for high-resolution image synthesis. arXiv preprint arXiv:2307.01952, 2023.
- Poole, B., Jain, A., Barron, J. T., and Mildenhall, B. Dreamfusion: Text-to-3d using 2d diffusion, 2022. URL https://arxiv.org/abs/2209.14988.
- Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al. Learning transferable visual models from natural language supervision. *arXiv preprint arXiv:2103.00020*, 2021.
- Ramesh, A., Dhariwal, P., Nichol, A., Chu, C., and Chen, M. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 1(2):3, 2022.
- Rombach, R., Blattmann, A., Lorenz, D., Esser, P., and Ommer, B. High-resolution image synthesis with latent diffusion models, 2021.
- Rombach, R., Blattmann, A., Lorenz, D., Esser, P., and Ommer, B. High-resolution image synthesis with latent diffusion models, 2022. URL https://arxiv.org/ abs/2112.10752.
- Scotti, P., Banerjee, A., Goode, J., et al. Reconstructing the mind's eye: fmri-to-image with contrastive learning and diffusion priors. *Advances in Neural Information Processing Systems*, 36, 2024a.
- Scotti, P. S., Tripathy, M., Torrico, C., Kneeland, R., Chen, T., Narang, A., Santhirasegaran, C., Xu, J., Naselaris, T., Norman, K. A., et al. Mindeye2: Shared-subject models enable fmri-to-image with 1 hour of data. In *Forty-first International Conference on Machine Learning*, 2024b.
- Shen, G., Zhao, D., He, X., Feng, L., Dong, Y., Wang, J., Zhang, Q., and Zeng, Y. Neuro-vision to language:

550

577

578 579

580

581

582 583

> 584 585

587

590

591

592 593

595 596 597

594

Decoding natural images from eeg for object recognition. arXiv preprint arXiv:2308.13234, 2023.

org/abs/1712.05884.

Spampinato, C., Palazzo, S., Kavasidis, I., Giordano, D., Shah, M., and Souly, N. Deep learning human mind

Image reconstruction and interaction via non-invasive

brain recordings. arXiv preprint arXiv:2404.19438, 2024.

Yang, Z., Chen, Z., Zhang, Y., Wang, Y., Skerry-Ryan,

R., Saurous, R. A., Agiomyrgiannakis, Y., and Wu, Y.

Natural tts synthesis by conditioning wavenet on mel

spectrogram predictions, 2018. URL https://arxiv.

Song, Y., Liu, B., Li, X., Shi, N., Wang, Y., and Gao, X.

Shen, J., Pang, R., Weiss, R. J., Schuster, M., Jaitly, N.,

for automated visual classification, 2019. URL https: //arxiv.org/abs/1609.00344.

Sun, J., Li, M., Chen, Z., Zhang, Y., Wang, S., and Moens, M.-F. Contrast, attend and diffuse to decode high-resolution images from brain activities, 2023. URL https://arxiv.org/abs/2305.17214.

Takagi, Y. and Nishimoto, S. High-resolution image reconstruction with latent diffusion models from human brain activity. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 576 14453-14463, 2023.

Tan, X., Chen, J., Liu, H., Cong, J., Zhang, C., Liu, Y., Wang, X., Leng, Y., Yi, Y., He, L., Soong, F., Qin, T., Zhao, S., and Liu, T.-Y. Naturalspeech: End-to-end text to speech synthesis with human-level quality, 2022. URL https://arxiv.org/abs/2205.04421.

Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.-A., Lacroix, T., Rozière, B., Goyal, N., Hambro, 586 E., Azhar, F., Rodriguez, A., Joulin, A., Grave, E., and Lample, G. Llama: Open and efficient foundation lan-588 guage models, 2023. URL https://arxiv.org/ 589 abs/2302.13971.

Wang, C., Gu, J., Hu, P., Xu, S., Xu, H., and Liang, X. Dreamvideo: High-fidelity image-to-video generation with image retention and text guidance, 2024a. URL https://arxiv.org/abs/2312.03018.

Wang, S., Liu, S., Tan, Z., and Wang, X. Mindbridge: A cross-subject brain decoding framework. In Proceedings of the IEEE/CVF Conference on Computer Vision and 598 Pattern Recognition, pp. 11333–11342, 2024b. 599

600 Wang, Y., Huang, N., Li, T., Yan, Y., and Zhang, X. 601 Medformer: A multi-granularity patching transformer 602 for medical time-series classification. arXiv preprint 603 arXiv:2405.19363, 2024c. 604

Wang, Y., Turnbull, A., Xiang, T., Xu, Y., Zhou, S., Masoud, A., Azizi, S., Lin, F. V., and Adeli, E. Decoding visual experience and mapping semantics through whole-brain analysis using fmri foundation models. arXiv preprint arXiv:2411.07121, 2024d.

Wei, Y., Cao, L., Li, H., and Dong, Y. Mb2c: Multimodal bidirectional cycle consistency for learning robust visual neural representations. In Proceedings of the 32nd ACM International Conference on Multimedia, pp. 8992–9000, 2024.

Wu, S., Fei, H., Qu, L., Ji, W., and Chua, T.-S. Next-gpt: Any-to-any multimodal llm. In Forty-first International Conference on Machine Learning, 2024.

Xia, W., de Charette, R., Öztireli, C., and Xue, J.-H. Umbrae: Unified multimodal decoding of brain signals. arXiv preprint arXiv:2404.07202, 2024a.

Xia, W., de Charette, R., Oztireli, C., and Xue, J.-H. Umbrae: Unified multimodal brain decoding. In European Conference on Computer Vision, pp. 242-259. Springer, 2024b.

Xia, W., de Charette, R., Öztireli, C., and Xue, J.-H. Dream: Visual decoding from reversing human visual system, 2024c. URL https://arxiv.org/abs/ 2310.02265.

Xue, S., Jin, B., Jiang, J., Guo, L., and Liu, J. A hybrid local-global neural network for visual classification using raw eeg signals. Scientific Reports, 14(1):27170, 2024.

Yang, F., Feng, C., Wang, D., Wang, T., Zeng, Z., Xu, Z., Park, H., Ji, P., Zhao, H., Li, Y., and Wong, A. Neurobind: Towards unified multimodal representations for neural signals, 2024a. URL https://arxiv.org/abs/ 2407.14020.

Yang, F., Feng, C., Wang, D., Wang, T., Zeng, Z., Xu, Z., Park, H., Ji, P., Zhao, H., Li, Y., et al. Neurobind: Towards unified multimodal representations for neural signals. arXiv preprint arXiv:2407.14020, 2024b.

Ye, Z., Yao, L., Zhang, Y., and Gustin, S. See what you see: Self-supervised cross-modal retrieval of visual stimuli from brain activity, 2022. URL https: //arxiv.org/abs/2208.03666.

Zhang, K., He, L., Jiang, X., Lu, W., Wang, D., and Gao, X. Cognitioncapturer: Decoding visual stimuli from human eeg signal with multimodal information. arXiv preprint arXiv:2412.10489, 2024.

Zhang, R., Isola, P., Efros, A. A., Shechtman, E., and Wang, O. The unreasonable effectiveness of deep features as a perceptual metric. In Proceedings of the IEEE conference

- 605 on computer vision and pattern recognition, pp. 586–595,606 2018.
- <sup>607</sup>
  <sup>608</sup> Zhang, Y., Gong, K., Zhang, K., Li, H., Qiao, Y.,
  <sup>609</sup> Ouyang, W., and Yue, X. Meta-transformer: A uni<sup>610</sup> fied framework for multimodal learning. *arXiv preprint*<sup>611</sup> *arXiv:2307.10802*, 2023.
- <sup>612</sup>
  <sup>613</sup>
  <sup>614</sup>
  <sup>614</sup>
  <sup>615</sup>
  <sup>616</sup>
  <sup>616</sup>
  <sup>617</sup>
  <sup>617</sup>
  <sup>618</sup>
  <sup>618</sup>
  <sup>619</sup>
  <sup>619</sup>
  <sup>619</sup>
  <sup>611</sup>
  <sup>611</sup>
  <sup>612</sup>
  <sup>612</sup>
  <sup>612</sup>
  <sup>612</sup>
  <sup>612</sup>
  <sup>613</sup>
  <sup>614</sup>
  <sup>615</sup>
  <sup>616</sup>
  <sup>617</sup>
  <sup>617</sup>
  <sup>618</sup>
  <sup>618</sup>
  <sup>619</sup>
  <sup>619</sup>
  <sup>619</sup>
  <sup>619</sup>
  <sup>619</sup>
  <sup>619</sup>
  <sup>611</sup>
  <sup>612</sup>
  <sup>612</sup>
  <sup>612</sup>
  <sup>612</sup>
  <sup>612</sup>
  <sup>613</sup>
  <sup>614</sup>
  <sup>615</sup>
  <sup>615</sup>
  <sup>616</sup>
  <sup>617</sup>
  <sup>617</sup>
  <sup>618</sup>
  <sup>618</sup>
  <sup>617</sup>
  <sup>618</sup>
  <sup>618</sup>
  <sup>619</sup>
  <
- <sup>619</sup> Zhou, Q., Du, C., Wang, S., and He, H. Clip-mused: Clipguided multi-subject visual neural information semantic decoding, 2024. URL https://arxiv.org/abs/ 2402.08994.
  - Zhu, X., Hu, Y., Mo, F., Li, M., and Wu, J. Unimed: A unified medical generalist foundation model for multi-task learning via connector-moe. *arXiv preprint arXiv:2409.17508*, 2024.

## A Appendix Overview

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

- Sec. B: Additional Ablation Experiments.
- Sec. C: Additional Implementation Details.
- Sec. D: Evaluation Details.
- Sec. E: Comparison with Prior Works.
- Sec. F: Additional Qualitative Results.

## **B** Additional Ablation Experiments

## B.1 Ablation Study of FLORA

To assess the contribution of each component within our 645 framework, we conducted ablation studies on various model 646 modules in Tab. 4. Using Medformer (Wang et al., 2024c) as 647 the baseline, we evaluated the impact of the temporal-spatial 648 convolution module and the Mixture of Experts (MoE) mod-649 ule. Additionally, we explored replacing the MoE module 650 with the Perceiver as the universal projection module; how-651 ever, this approach resulted in a prohibitively large number 652 of parameters, slower computation times, and suboptimal 653 654 performance.

The experimental results substantiate that our proposed architecture effectively optimizes the trade-off between model
capacity and computational efficiency. The temporal-spatial
convolutional module preserves predictive performance

while substantially reducing the parameter count, thereby facilitating the scalable processing of high-dimensional neural datasets. More critically, our soft-routing universal projection mechanism, when trained with heterogeneous multimodal neural data, learns invariant and transferable feature representations. Traditional MoE module selects a single expert based on the maximum score, while softrouter implements a weighted combination of experts. This architectural paradigm is particularly advantageous for joint-subject modeling, as the exposure to diverse neurophysiological patterns across modalities during training enhances FLORA's robustness in addressing inter-subject variability while maintaining a significantly lower computational burden compared to alternatives such as Perceiver (Alayrac et al., 2022). The empirical results demonstrate that this routing strategy in universal projection facilitates knowledge transfer across modalities, contributing to strong generalization capabilities in zero-shot cross-modal inference.

The ablation study corroborates our design rationale, highlighting the efficacy of multimodal training in bolstering model generalization—an attribute that becomes increasingly pivotal when scaling to expansive neural datasets encompassing extensive subject pools and multiple trial repetitions.

### B.2 Supervised Scaling Up on Data Size

We analyzed the effect of training set size on joint subject training performance using THINGS-EEG2, as illustrated in 6A. This factor is critical for model efficacy. The training set comprised 1654 distinct concepts  $\times$  10 image conditions, each repeated four times. We preserved all four repetitions for training. Expanding the dataset size by increasing the number of conditions yielded substantial improvements in decoding accuracy, particularly in transitions from 20% to 40% and from 40% to 60%. Additionally, augmenting repetitions contributed significantly to performance gains, notably from 80% to 100%. These findings indicate that, under joint-subject training, larger datasets facilitate further performance enhancements. Our approach demonstrates strong scalability, underscoring its potential for improving neural decoding with increased data availability.

A key advantage of our cross-modality joint-subject training framework is its capacity to enable subject adaptation with minimal training data. To rigorously assess the generalization capability of our approach, we leverage the pretrained unified encoder on THINGS-EEG2, THINGS-MEG, and THINGS-fMRI to adapt the joint-subject model to THINGS-EEG1, utilizing varying amounts of training data, as illustrated in Fig. 6B. The dataset partitioning strategy for training and evaluation in THINGS-EEG1 remains consistent with its original split. We perform adaptation evaluation on a 200-way image retrieval task (chance level: 0.5%). Given

Table 4. Retrieval performance of FLORA under varying input projection mechanisms. The first two rows compare time-series Transformers trained using an attention-based algorithm without any projection or feature fusion. The following rows illustrate the impact of FLORA enhancements. "Medformer": A multi-granularity time-series model designed for medical data. "Medformer+TSConv": Medformer with the addition of temporal-spatial convolution. "Re. EEG": Retrieval performance on the THINGS-EEG2 dataset. "Re. MEG": Retrieval performance on the THINGS-MEG dataset. "Re. fMRI": Retrieval performance on the THINGS-fMRI dataset.

#	Description	Para.	Model	Uni. Proj.	Softrouter	Re. EEG ↑	Re. MEG ↑	Re. fMRI †
1	Med. (Linear)	97.9M	Time series	X	×	18.15	3.50	25.00
2	Med. (TSConv)	18.67M	FLORA-uni.	×	×	26.50	3.00	24.00
3	+Uni. Proj.	162.0M	FLORA-multi.	1	X	26.45	8.38	24.33
4	+Softrouter	162.0M	FLORA-uni.	1	1	25.35	8.00	26.33
	+Softrouter	162.0M	FLORA-multi.	1	1	25.05	6.88	28.33

the inherently low signal-to-noise ratio in EEG data, the zero-shot performance of our model approximates chance level. However, after 50 epochs of fine-tuning, the model exhibits notable scalability, underscoring its potential for efficient subject adaptation and robust cross-modality transfer.

660

661

662

663

664

675

676

677

678

679

698

699

700

701

704 705

706



*Figure 6.* Effect of data size in training and test set (THINGS-EEG2), and weakly-supervised subject adaptation (THINGS-EEG1).

## C Additional Implementation Details

#### C.1 Neural Modality Encoders

707**EEG Encoder.** The shape of an EEG signal is  $R^{63\times250}$ .708With Medformer, we get multiple sets of tokens output709 $x \in \mathbb{R}^{221\times250}$ . By temporal-spatial convolution layers,710we then resize the tensor  $x \in \mathbb{R}^{221\times250}$  into a 1D tensor711 $x \in \mathbb{R}^{1\times1024}$ .

712 713 **MEG Encoder.** The shape of an MEG signal is  $R^{271 \times 201}$ . 714 With Medformer, we get multiple sets of tokens output  $x \in \mathbb{R}^{178 \times 250}$ . By temporal-spatial convolution layers, we then resize the tensor  $x \in \mathbb{R}^{178 \times 250}$  into a 1D tensor  $x \in \mathbb{R}^{1 \times 1024}$ .

**fMRI Encoder.** In the THINGS-fMRI data set, the original voxel number of visual ROI (Region of Interest) of each of the three subjects was 6036, 5944 and 5238. In order to facilitate model processing, pad was unified to 7000 when data was loaded. So the shape of an fMRI signal is  $R^{7000}$ . We tokenize it with a linear layer: Linear( $C_{\rm in} = 7000, C_{\rm out} = 8192$ ). We then resize the output tensor  $x \in \mathbb{R}^{8192}$  into a 2D tensor  $x \in \mathbb{R}^{8 \times 1024}$  to align with the input of the transformer encoder. With Medformer, we get multiple sets of tokens output  $x \in \mathbb{R}^{899 \times 250}$ . By temporal-spatial convolution layers, we then resize the tensor  $x \in \mathbb{R}^{899 \times 250}$  into a 1D tensor  $x \in \mathbb{R}^{1 \times 1024}$ .

#### C.2 Low-Level Pipeline

Compared to vision-centric pretraining paradigms such as ViT, ResNet, and DINO, the CLIP vision model exhibits a deficiency in capturing fine-grained low-level visual features. To mitigate this limitation, our framework integrates a dedicated low-level visual reconstruction pipeline. Specifically, we aim to reconstruct fundamental perceptual attributes—such as contour, posture, and orientation—by leveraging EEG-derived representations. This is achieved through an alignment mechanism with the latent space of a variational autoencoder (VAE), facilitating the recovery of pixel-level structural information. By enforcing this latent alignment, our approach enhances the model's capacity to infer and preserve crucial low-level visual semantics from neural signals, thereby improving the fidelity of EEG-tovision mappings.

We trained the low-level reconstruction pipeline for 200 epochs, employing a latent mean squared error (MSE) loss in conjunction with a contrastive learning loss and a variational autoencoder (VAE) image reconstruction loss. The objective was to align the  $4 \times 64 \times 64$  EEG-derived latent representations—obtained via a projection layer and an up-

sampling CNN—with the VAE latent space. However, we
observed that the reconstruction loss and contrastive learning loss underperformed compared to solely optimizing the
latent space loss, while also imposing significantly higher
GPU memory demands.

Furthermore, our investigation revealed that incorporating a low-level visual model for knowledge distillation in the lowlevel reconstruction pipeline not only failed to enhance VAE latent alignment but also exacerbated overfitting. These findings suggest that zero-shot low-level visual reconstruction from EEG signals lacks stability and may introduce misleading artifacts in the generated outputs.

In our framework, when utilizing the low-level reconstruction pipeline, we typically set the inference steps of SDXL to 10 (or SDXL-Turbo to 4) and configure the image-to-image denoising strength to 0.5.

# **D** Evaluation Details

728

729

730

731

732 733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

756

757

758

759

760

761

762

763

764

765

766

767

769

In this section, we will give more evaluation details.

Visual Retrieval tasks. We use forward retrieval (all image embeddings in test set is retrieved using one neural data embedding) for evaluation. For the neural data of each modality, we output it as  $x \in \mathbb{R}^{1 \times 1024}$  and calculate the dot product similarity with the image embedding, taking the image with the greatest similarity as the image output by zero-shot. For THINGS-EEG1, the test set has 12 conditions per concept, so we only take the first image condition among the 200 concepts in the test set, that is, we use 200 images from 200 categories to test. For THINGS-EEG2, the test set has only 1 condition per concept, so we test with 200 images from 200 categories. For THINGS-MEG, there are 200 concepts in the test set, each of which has 12 conditions, so we only use the first condition of each concept as the test, with a total of 200 images from 200 categories to test. For THINGS-fMRI, the test set is 100 concepts, where each concept has 1 condition, so a total of 200 images from 200 categories are tested.

Visual Reconstruction tasks. The test set configuration of the visual reconstruction task is exactly the same as that of the retrieval task, except that the image reconstruction has more indicators to measure the performance of the image reconstruction. It should be emphasized that, in contrast to existing approaches on the THINGS dataset which train models independently for each subject's neural data and conduct visual reconstruction, FLORA-multimodal necessitates addressing highly heterogeneous multimodal and cross-subject neural data in both the training phase and visual reconstruction inference evaluation. Moreover, within the fMRI modality, FLORA employs the THINGS-fMRI dataset, which is smaller in scale than the NSD dataset used in existing major works, thus requiring the model to more efficiently utilize limited data resources.

Visual Captioning tasks. The idea of employing Prior diffusion to alter the distribution of model features is inspired by the two-stage image reconstruction approach described in ATM-S (Li et al., 2024b), with the expectation of achieving similar effects in the domain of caption generation. Although the final quantitative results show only marginal differences, the incorporation of Prior diffusion leads to more structurally coherent and semantically accurate captions that better reflect the content of the original images. The test set configuration for the visual captioning task is exactly the same as for the retrieval task. The THINGS image data set lacks the original text captions annotation. To generate captions for the test set images as ground truth, we used the prompt "Please provide a concise and formal sentence to describe the image, beginning with 'An image of,' and keeping it to approximately 10 words." from the large model Kimi to obtain descriptions for each image in the test set. By using different MMLMs as interfaces, we can generate captions of different qualities from the latent, and we give an assessment of captioning quality for different MMLMs adapted to FLORA.

# E Comparison with Prior Works

The key distinction between FLORA and previous neural decoding approaches lies in its demonstration that specialized neural encoders, when coupled with a universal projection layer, can effectively align diverse neural modalities with visual representations in a unified framework. As evidenced in Tab. 5, while prior works typically develop independent architectures for each modality, FLORA achieves multimodal integration across EEG, MEG, and fMRI signals while maintaining a compact parameter footprint (7.36M). The empirical results (Tab. 1) demonstrate competitive performance across modalities compared to modality-specific models, while ablation studies (Tab. 4) reveal that joint training within this unified framework provides particular benefits for modalities with limited data availability.

The effectiveness of FLORA's architecture in handling diverse neural signals suggests broader implications for scalable neural decoding. By successfully integrating multiple modalities through a shared projection space while maintaining computational efficiency, FLORA provides insights into architectural design principles for larger-scale neural models. The framework's demonstrated capabilities in addressing fundamental challenges - from cross-subject variability to zero-shot generalization - while preserving modality-specific characteristics, establish a foundation for developing more comprehensive neural decoding systems. These results not only validate the feasibility of unified multimodal frameworks but also suggest promising directions

for scaling neural decoding models toward more ambitious 771 multimodal applications.

773

774

803

804

805

806

807

808 809

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

775 Table 5. Comparison between FLORA and related works on visual retrieval, reconstruction, and caption tasks. Methods compared 777 include MindEye2 (Scotti et al., 2024b), NICE (Song et al., 2023), 778 B.D. (Benchetrit et al., 2024), N.F. (Ferrante et al., 2024), NeuroVL (Shen et al., 2024), UMBRAE (Xia et al., 2024b), V.D. (Li 779 et al., 2024b), and NeuroBind (Yang et al., 2024b). 780

Method	Modalities	Tasks
MindEye2	<b>8</b> 2	رج 📭 🧟
NICE	a Mara	
B.D.	~* <b>\</b> [~~ &\$	<b>F</b>
N.F.	11 11	
NeuroVL	820 B	
UMBRAE		
V.D.	~ Mar ~ Mar	چ 📭 م
NeuroBind	1 Fara 1 Fara	ج 🗖 کې
FLORA	+ 1 For + 1 For &	
• · · · · EEG	Meg 💥 i fMRI	: CI 💛 :
👝 : Retrieval	Fee : Caption C : Grour	iding 🎺 : Reconst

#### F **Additional Qualitative Results**

To achieve cross-modal concept-level zero-shot retrieval, we realigned the THINGS-MEG dataset split to match THINGS-EEG2. Tab.6 demonstrates the impact of this realignment and trial averaging on retrieval performance. Notably, due to the realignment, some concepts in both training and test sets inherit different numbers of trials from the original MEG dataset structure. To maintain consistency, we standardized to one trial per concept in the "unave." setting. However, for concepts that were originally in the MEG test set (which had 12 trials per concept), we also evaluated a "averaged" setting where these 12 trials were averaged. The performance improvement from FLORA-unimodal (unave.) to FLORA-unimodal demonstrates the benefit of this trial averaging, which effectively improves the signal-to-noise ratio for these concepts.

Table 6. Impact of average retrieval performance across multiple conditions on the THINGS-MEG test set. We report the retrieval performance of FLORA on the THINGS-MEG dataset under various evaluation settings.

Madal	2-way	4-way	10-way	200-way	200-way
Model	Top-1	Top-1	Top-1	Top-1	Top-5
FLORA-unimodal (unave.)	72.38	54.62	36.88	5.75	17.87
FLORA-unimodal	81.75	64.50	46.62	8.00	24.38
FLORA-multimodal (unave.)	74.75	52.00	36.38	6.12	18.25
FLORA-multimodal	<u>80.50</u>	<u>61.88</u>	<u>39.75</u>	<u>6.88</u>	<u>23.38</u>

The original THINGS-MEG dataset presents this evaluation scenario in its data partitioning: the training set contains 1,654 concepts (12 images per concept, single repetition), while the test set comprises 200 concepts (1 image per concept, 12 repetitions per image) in a concept-level zero-shot setting. In this experiment, FLORA-unimodal was trained without the MoE-based universal projection and multimodal training, using an independent model for each subject's MEG recordings in an in-subject evaluation protocol. While our model has demonstrated superior performance in crosssubject zero-shot scenarios with richer data availability (as shown in Tab.1), it maintains competitive performance in this more constrained setting, as evidenced in Tab.7. The results validate that FLORA's core architecture provides robust performance even when trained on single-subject data with limited per-concept samples, without leveraging the benefits of multimodal training data and cross-subject information.

Table 7. Retrieval performance on the original THINGS-MEG dataset. We trained and evaluated various methods on the original THINGS-MEG training and test sets. Each method was independently trained and assessed for each subject, with results reported for different retrieval configurations: 2-way, 4-way, 10way, as well as Top-1 and Top-5 accuracy for the 200-way retrieval task.The methods compared include NICE (Song et al., 2023), EEGNetV4 (Lawhern et al., 2018), B.D., MindEyeV2 (Scotti et al., 2024b) and ATM-S (Li et al., 2024b).

Model	2-way	4-way	10-way	200-way	200-way
Widdel	Top-1	Top-1	Top-1	Top-1	Top-5
NICE	89.29	76.94	60.19	17.12	40.15
EEGNetV4	85.34	71.78	52.85	12.76	32.38
B.D.	77.16	57.76	37.82	6.06	19.38
ATM-S	90.25	79.50	63.66	17.84	44.73
FLORA-unimodal	87.75	75.25	57.50	15.62	39.37

Our evaluation demonstrates that FLORA achieves strong performance in both in-subject and cross-subject settings on the THINGS-EEG2 dataset. When tested in the in-subject paradigm, where models are independently trained and tested on individual subjects' data as is shown in Tab.8, FLORA maintains competitive performance.

Notably, our cross-subject evaluation results shown in Tab.1 show better retrieval performance compared to the in-subject setting, validating our model's effectiveness in addressing 825 Table 8. In-subject Retrieval on THINGS-EEG2. We report the 826 retrieval performance of different methods on the THINGS-EEG2 827 (in-subject) datasets. Each method is trained and evaluated in each 828 subject independently, with results for different retrieval configura-829 tions: 2-way, 4-way, 10-way, and the Top-1 and Top-5 accuracy of 830 200-way. The methods compared include NICE (Song et al., 2023), 831 EEGNetV4 (Lawhern et al., 2018), B.D. (Benchetrit et al., 2024), 832 MindEyeV2 (Scotti et al., 2024b), UMBRAE (Xia et al., 2024a), ATM-S (Li et al., 2024b), CongnitionCapturer (Zhang et al., 2024), 833 MindBridge (Wang et al., 2024b), WAVE (Wang et al., 2024d) and 834 025 MB2C (Wei et al., 2024).

1 1	· · · · · · · · · · · · · · · · · · ·					
6			THIN	GS-EEG2	(in-subject	)
0	Model	2-way	4-way	10-way	200-way	200-way
		Top-1	Top-1	Top-1	Top-1	Top-5
	NICE	93.23	83.93	69.22	21.67	51.34
	EEGNetV4	91.42	80.21	63.37	16.84	42.58
	B.D.	89.03	56.77	56.77	13.29	35.50
	CogCap	93.15	82.85	69.35	22.05	51.60
	MB2C	78.40	62.25	43.75	8.85	25.20
	Wave	88.65	73.70	54.70	13.00	34.50
	MindBridge	89.10	75.35	57.45	15.00	39.15
	ATM-S	94.70	86.73	74.00	26.85	57.21
	MindEyeV2	92.50	82.80	66.10	23.80	50.25
	FLORA-unimodal	<u>93.75</u>	85.70	<u>69.65</u>	21.85	51.65

cross-subject challenges. FLORA's robust performance can
be attributed to its efficient encoder architecture across different data scales and the flexibility to incorporate or remove
components like the MoE module. These results confirm
FLORA's capability to deliver strong performance across
various experimental paradigms while effectively leveraging
larger-scale datasets when available.

847

855

876

877

878

879

856 Table 9. In-subject Retrieval on THINGS-MEG. We report the 857 retrieval performance of different methods on the THINGS-MEG 858 dataset in the in-subject setting. Each method is trained and eval-859 uated in each subject independently, with results for different 860 retrieval configurations: 2-way, 4-way, 10-way, and the Top-1 861 and Top-5 accuracy of 200-way. The methods compared include 862 NICE (Song et al., 2023), EEGNetV4 (Lawhern et al., 2018), B.D., 863 MindEyeV2 (Scotti et al., 2024b) and ATM-S (Li et al., 2024b).

64			THIN	GS-MEG	(in-subject)	)
65	Model	2-way	4-way	10-way	200-way	200-way
05		Top-1	Top-1	Top-1	Top-1	Top-5
66	NICE (unave.)	68.65	43.63	24.88	2.63	9.25
67	EEGNetV4 (unave.)	69.25	49.63	29.00	4.37	13.25
57	B.D. (unave.)	56.75	35.00	16.12	0.75	5.37
80	MindEyeV2 (unave.)	68.12	45.13	24.62	2.25	9.37
59	ATM-S (unave.)	80.13	60.75	42.75	7.38	21.75
0	FLORA-unimodal(unave.)	76.25	<u>59.12</u>	37.25	<u>6.25</u>	17.25
/	NICE	79.75	60.75	40.00	6.00	20.50
L	EEGNetV4	80.37	62.87	41.00	6.75	21.50
2	B.D.	64.25	43.32	21.50	2.00	8.50
2	MindEyeV2	76.25	56.88	37.62	6.00	18.75
	ATM-S	82.87	68.00	46.13	8.12	27.12
ŀ	FLORA-unimodal	<u>81.75</u>	<u>66.00</u>	48.13	10.25	26.25
5						

The in-subject evaluation on the THINGS-MEG dataset demonstrates that FLORA-unimodal achieves state-of-theart performance in the averaged condition while maintaining robust performance in the unaveraged setting (Tab.9). This superior performance, particularly with high-quality averaged MEG signals, validates FLORA's strong decoding capabilities even when operating on smaller-scale datasets.

This consistent performance can be attributed to two fundamental aspects of our architecture: first, the model's adaptable design that effectively accommodates various experimental paradigms, and second, its ability to efficiently utilize the enhanced signal-to-noise ratios in averaged MEG data without relying on the advantages of large-scale datasets. These findings further substantiate FLORA's effectiveness as a versatile neural decoding framework, demonstrating robust performance across both jointsubject and in-subject experimental paradigms.

*Table 10.* **In-subject Retrieval on THINGS-fMRI.** We report the retrieval performance of different methods on the THINGS-fMRI dataset. Each method is trained and evaluated in each subject independently, with results for different retrieval configurations: 2-way, 4-way, 10-way, and the Top-1 and Top-5 accuracy of 200-way. The methods compared include NICE (Song et al., 2023), EEGNetV4 (Lawhern et al., 2018), B.D. (Benchetrit et al., 2024), MindEyeV2 (Scotti et al., 2024b), UMBRAE (Xia et al., 2024a), ATM-S (Li et al., 2024b), CongnitionCapturer (Zhang et al., 2024), MindBridge (Wang et al., 2024b), WAVE (Wang et al., 2024d) and MB2C (Wei et al., 2024).

	THINGS-fMRI (in-subject)									
Model	2-way	4-way	10-way	100-way	100-way					
	Top-1	Top-1	Top-1	Top-1	Top-5					
NICE	90.00	82.67	62.67	25.00	56.00					
EEGNetV4	89.00	76.00	56.00	20.33	50.67					
B.D.	80.33	69.67	33.33	9.67	24.00					
CogCap	93.67	81.33	64.00	27.33	64.00					
MB2C	93.00	80.33	61.33	23.67	57.00					
Neuro.V2L	86.67	70.67	50.67	15.00	39.67					
Wave	91.00	79.67	65.67	29.00	61.00					
ATM-S	95.00	81.00	66.67	25.67	60.33					
UMBRAE	89.00	77.00	60.00	22.00	50.67					
MindEyeV2	95.00	86.67	76.67	32.67	69.00					
FLORA-unimodal	<u>92.33</u>	87.00	<u>67.00</u>	24.33	<u>61.67</u>					

The in-subject evaluation on the THINGS-fMRI dataset demonstrates FLORA-unimodal's strong adaptability in decoding capabilities, even when operating on limited singlesubject data. As shown in Tab.10, our model maintains competitive performance, validating its effectiveness across varying data scales and modalities. These results substantiate that FLORA, while primarily optimized for cross-subject scenarios with larger-scale datasets, retains robust decoding capabilities even when constrained to limited single-subject fMRI data. This finding further underscores the versatility of our architectural design across diverse experimental paradigms.

Examining in-subject retrieval performanceTab. 8, 9, and 10 demonstrate FLORA's consistent performance across EEG, MEG, and fMRI modalities in in-subject settings, showing

880 that our model, while optimized for cross-subject general-881 ization, maintains competitive performance when trained on 882 single-subject data. Comparing with Tab.1, notably, while 883 traditional approaches often suffer significant performance 884 degradation in cross-subject scenarios, FLORA maintains or 885 improves its performance when moving from in-subject to 886 joint-subject settings, proposing a solution to joint-subject 887 large neural models.

## 889 G Additional images results

888

890

909 910

911

912

913 914 915

917

918

919 920 921

922

923

929

930

931

932 933 934

We present a qualitative comparison of the highest-, 891 moderate-, and lowest-fidelity generated images in Fig. 7. 892 EEG data were randomly sampled from all subjects view-893 894 ing 200 images, and corresponding EEG-derived embeddings were extracted to condition the image generation 895 process. By computing the cosine similarity between the 896 CLIP embeddings of the generated and original images, 897 898 we identified six examples for each fidelity category. In the highest-fidelity group, the generated images exhibit 899 900 strong semantic correspondence with the original images while effectively preserving low-level visual attributes. In 901 the moderate-fidelity group, the generated images main-902 903 tain semantic coherence with the original stimuli, yet exhibit partial alterations in low-level visual details. In the 904 lowest-fidelity group, both semantic integrity and low-level 905 906 visual structures are substantially distorted, leading to a pronounced divergence from the original images. See Fig. 8 907 and Fig. 9 for MEG and fMRI reconstruction results. 908



*Figure 8.* Examples of cross-modality MEG-guided crosssubject visual reconstruction from FLORA. From top to bottom, we exhibit the best, median, and worst 6 generated images, respectively. We show the images subjects seen and the generated images by our two-stage image generator.



*Figure 7.* Examples of cross-modality EEG-guided crosssubject visual reconstruction from FLORA. From top to bottom, we exhibit the best, median, and worst 6 generated images, respectively. We show the images subjects seen and the generated images by our two-stage image generator.



*Figure 9.* Examples of cross-modality fMRI-guided crosssubject visual reconstruction from FLORA. From top to bottom, we exhibit the best, median, and worst 6 generated images, respectively. We show the images subjects seen and the generated images by our two-stage image generator.