# BEYOND GRID-LOCKED VOXELS: NEURAL RE-SPONSE FUNCTIONS FOR CONTINUOUS BRAIN EN-CODING

**Anonymous authors**Paper under double-blind review

000

001

002

004

006

008 009 010

011 012 013

014

015

016

017

018

019

021

025

026

027

028

029

031

032

034

037

039

040

041

042

043

044

046

047

048

051

052

## **ABSTRACT**

Neural encoding models aim to predict fMRI-measured brain responses to natural images. fMRI data is acquired as a 3D volume of voxels, where each voxel has a defined spatial location in the brain. However, conventional encoding models often flatten this volume into a 1D vector and treat voxel responses as independent outputs. This removes spatial context, discards anatomical information, and ties each model to a subject-specific voxel grid. We introduce the NRF Neural Response Function, a framework that models fMRI activity as a continuous function over anatomical space rather than a flat vector of voxels. NRF represents brain activity as a continuous implicit function: given an image and a spatial coordinate (x, y, z) in standardized MNI space, the model predicts the response at that location. This formulation decouples predictions from the training grid, supports querying at arbitrary spatial resolutions, and enables resolution-agnostic analyses. By grounding the model in anatomical space, NRF exploits two key properties of brain responses: (1) local smoothness—neighboring voxels exhibit similar response patterns; modeling responses continuously captures these correlations and improves data efficiency, and (2) cross-subject alignment—MNI coordinates unify data across individuals, allowing a model pretrained on one subject to be fine-tuned on new subjects. In experiments, NRF outperformed baseline models in both intrasubject encoding and cross-subject adaptation. Achieving high performance while reducing the data size needed by orders of magnitude. To our knowledge, NRF is the first anatomically aware encoding model to move beyond flattened voxels, learning a continuous mapping from images to brain responses in 3D space.

## 1 Introduction

A major goal in computational neuroscience is to understand how the human brain maps visual stimuli into neural activity. Neural encoding models aim to address this by predicting neural responsestypically measured by fMRI—from visual stimuli. These models offer powerful tools for analyzing high-dimensional brain data and probing the representations encoded in the visual Downing et al. (2001); Epstein & Kanwisher (1998); Gu et al. (2022); Heeger & Ress (2002); Kanwisher et al. (1997); Naselaris et al. (2011); Huth et al. (2012).

However, the real-world utility of current neural encoding models remains limited. Current neural encoding models represent fMRI responses as a 1D vector in  $\mathbb{R}^n$ , where n is a subject-specific voxel count. Naselaris et al. (2015)St-Yves & Naselaris (2018)Wang et al. (2023)Yamins et al. (2014) This discrete formulation has two critical limitations: 1) Ignoring 3D structure. By flattening fMRI volumes into 1D vectors, conventional models discard spatial information. This removes local context: anatomically adjacent voxels, which are often functionally correlated, are instead treated as independent outputs. 2)Subject specific. Each model is tied to the voxel grid of a single subject, making it non-transferable across individuals. Because voxel counts differ across brains, the output dimensionality of conventional models—and therefore their architecture—is tied to a single subject. As a result knowledge learned from one individual cannot be directly transferred to another, forcing each new subject to require training a separate model from scratch.

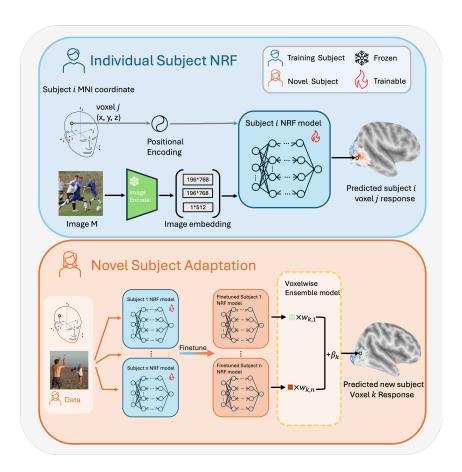


Figure 1: **Overview of NRF. Top:** Individual-subject NRF. Brain responses are modeled as a continuous function of both image features and anatomical coordinates in MNI space. NRF learns how to map to voxel responses while capturing correlations between neighboring voxels through anatomical coordinates. **Bottom:** New subject adaptation. 1) For novel subjects, the learned representation could be transferred by fine-tuning pretrained NRFs with limited data. 2) Predictions from multiple finetuned base models are combined via voxelwise ensembling to capture individual variability. NRF thus moves beyond grid-locked voxel models, offering a continuous, anatomically grounded representation that enables both data-efficient single-subject encoding and flexible cross-subject transfer.

As a result, models fail to exploit the 3D geometry of brain activity, wasting statistical power and requiring more data to learn accurate mappings. These issues are particularly problematic in realistic applications, where data are scarce. Unlike large-scale efforts such as NSD Allen et al. (2022), which collected tens of thousands of trials per subject over more than a year of scanning, most studies can only acquire a few hundred trials per subject due to cost and time constraints. Thus, the inefficiencies of discrete neural encoding models are amplified in real-world low-data regimes. In short, previous encoding models are *grid-locked*: they can only predict responses at the discrete sampling points they were trained on, for a single subject, and at one resolution. However, the human brain is a continuous 3D structure. Within each subject, neighboring voxels also exhibit similar response patterns, reflecting the spatial smoothness of neural activity. Despite individual variability, the visual cortex is highly conserved across people: areas such as the fusiform face area (FFA) and extrastriate body area (EBA) consistently respond to the same stimulus categories, and these regions align well across subjects when mapped into standardized anatomical templates such as MNI space. Ignoring this organization—both local smoothness and cross-subject correspondences — discards valuable structure that could enable more efficient learning and better generalization.

To address these limitations, we propose the **Neural Response Function** (NRF), a coordinate-based neural encoding model that predicts fMRI responses as a continuous function over anatomical space. Given a stimulus M and a spatial coordinate  $\mathbf{x}=(x,y,z)$  in standardized MNI space, NRF outputs the predicted brain response r at that location:

$$\Phi(M, \mathbf{x}) = r, \mathbf{x} \in \mathbb{R}^3, \ r \in \mathbb{R}.$$

This formulation directly addresses the limitations of previous models:

- Exploiting local smoothness. By conditioning predictions on anatomical coordinates, NRF incorporates the 3D spatial structure of fMRI data. This allows nearby voxels—often anatomically connected and functionally correlated—to share information, rather than being treated as independent outputs. As a result, NRF captures local smoothness in brain responses and achieves greater data efficiency.
- Efficient new-subject adaptation. NRF grounds predictions in standardized MNI space, unifying responses across subjects in a shared coordinate system. This enables direct transfer: a model pretrained on one subject can be adapted to a new individual with only minimal fine-tuning. We further introduce a finetune—ensemble strategy that leverages multiple pretrained models to boost adaptation accuracy, reducing the need for extensive subject-specific data collection in real-world settings.
- **Resolution-agnostic predictions.** By defining responses in continuous space, NRF decouples predictions from the voxel grid of the data model that is trained. It can be queried at arbitrary spatial coordinates, independent of the voxel size or sampling scheme used in data collection, and can seamlessly integrate data acquired at different resolutions. This flexibility opens the door to building more general-purpose brain models, moving closer to a functional *digital twin* of the brain.

Through experiments, we show that NRF - the first brain anatomy structure-grounded neural encoding model - offers a new direction for efficient and generalizable neural encoding.

## 2 RELATED WORK

**Neural Encoding models** fMRI encoding models have been extensively studied over the past two decades Mitchell et al. (2008); Huth et al. (2016); Gu et al. (2022); Tang et al. (2023); Kay et al. (2008); Güçlü & Van Gerven (2015); Naselaris et al. (2015). Most existing approaches treat fMRI data as discrete, formulating the problem as a regression task that maps image features to voxelwise responses Naselaris et al. (2011); Han et al. (2019); Wang et al. (2023). In these models, fMRI responses are flattened as a 1-D vector, and each voxel is treated independently, ignoring the anatomical structure of the brain. As a result, they fail to capture the inherent local smoothness of fMRI responses, where neighboring voxels often exhibit correlated activity. To our knowledge, we are the first to develop an anatomically aware continuous encoding model, which leverages the 3D structure of the brain to improve data efficiency and generalization.

Implicit neural representation Implicit neural representations have emerged as a powerful paradigm for modeling continuous signals in computer vision and graphics. Instead of storing data on fixed grids, INRs represent signals such as images Sitzmann et al. (2020) and 3D shapes Park et al. (2019); Mildenhall et al. (2021); Chen & Zhang (2019); Mescheder et al. (2019) as continuous functions parameterized by neural networks. A key advantage of this framework is its ability to capture fine-grained structure and support resolution-agnostic queries. Inspired by this line of work, we adopt a similar coordinate-based formulation for fMRI encoding. Unlike prior voxelwise models that discretize the brain into subject-specific grids, our approach treats brain responses as a continuous function over standardized anatomical coordinates. To our knowledge, NRF is the first attempt to bring the implicit representation framework to computational neuroscience, enabling anatomically aware, resolution-agnostic modeling of fMRI responses.

# 3 METHOD

#### 3.1 Modeling Brain Response Mapping with implicit neural representation

Current encoding models can be summarized in two steps: flatten neural response into 1D vectors, then train an encoding model that takes an image or its embedding as input and directly outputs the predicted response as a flattened vector in  $\mathbb{R}^n$ . Ignoring the 3D spatial information and forcing models to be trained separately for each subject. This leads to poor data efficiency. Our key insight is that brain response should be modeled in its anatomical context. We represent the brain response mapping as a **continuous function** over MNI coordinates, a standardized anatomical space. Formally, given a stimulus image M and a spatial coordinate  $\mathbf{x}=(x,y,z)\in\mathbb{R}^3$ , the Neural Response Function (NRF) outputs the predicted fMRI response  $\hat{r}\in\mathbb{R}$  at that location:

$$\Phi(M, \mathbf{x}) = \hat{r}.$$

Rather than outputting a fixed-length vector tied to a particular subject's voxel grid, NRF predicts the response at any coordinate  $(x,y,z) \in \mathbb{R}^3$ . This shift from 1D discrete outputs to 3D spatial coordinate conditioned continuous predictions makes the model anatomically aware and able to exploit spatial smoothness during training and inference. Because  $\Phi$  is defined over  $\mathbb{R}^3$ , it can be queried at arbitrary spatial resolutions, independent of the voxel grid or sampling scheme used during acquisition. This enables flexible data analysis: fMRI responses can be resampled seamlessly at different resolutions, supporting resolution-agnostic modeling and analysis.

Architecture. We instantiate  $\Phi$  using a two-component design. The first component, G, the image feature extraction block. It extracts multi-scale features from the stimulus image M, capturing both low-level and high-level representations. These features are fused together to obtain a final image embedding G(M). The second component is an implicit neural representation predictor P that conditions on both G(M) and the spatial coordinate  $\mathbf{x}$ . The coordinate is first encoded using Fourier features Tancik et al. (2020):

$$\gamma(\mathbf{x}) = [\cos(2\pi b_1^T \mathbf{x}), \sin(2\pi b_1^T \mathbf{x}), \dots, \cos(2\pi b_m^T \mathbf{x}), \sin(2\pi b_m^T \mathbf{x})]^T,$$

where  $b_j$  are sampled from an isotropic Gaussian. Finally, G(M) and  $\gamma(\mathbf{x})$  are concatenated and passed through an MLP predictor P:

$$\Phi(M, \mathbf{x}) = P(G(M), \gamma(\mathbf{x})).$$

This design makes the mapping explicitly anatomy-aware by conditioning on both image content and spatial location. Details can be found in the appendix A.2.

**Model Training.** The model is trained end-to-end with all components learned together, where the objective of the model is to correctly predict the voxel activation on each input image. Training batches are constructed from 32 randomly selected images, where for each image, we randomly sample 2000 voxels (out of 13000-15000 voxels) along with their corresponding fMRI activations for prediction. The model is trained using the Adam optimizer with a learning rate of 3e-3. For the loss function, we employ the same loss as in Beliy et al. (2019), a convex combination of mean square error and cosine similarity between the predicted response  $\hat{r}$  and ground truth fMRI, r. The fMRI loss is defined as:

$$L(\hat{r}, r) = (1 - \alpha) \|\hat{r}, r\|_2 + \alpha * cos(\angle(\hat{r}, r))$$

Where  $\alpha$  is set to 0.1 during training, which balances absolute error minimization (via MSE) with representational alignment (via cosine similarity).

#### 3.2 Cross Subject Transfer

A major challenge in training visual encoding models is the limited availability of subject-specific data. Collecting fMRI responses for thousands of images requires many hours of scanning, often across multiple sessions, and is infeasible in most clinical or experimental settings. In practice, new

subjects often contribute only a few hundred trials. Discrete neural encoding models underperform in this regime because they are tied to subject-specific voxel grids: each subject requires training a new model from scratch, and knowledge cannot be transferred directly across individuals.

NRF overcomes this limitation by being **voxel-grid agnostic**. Since responses are defined as a continuous function over standardized MNI space, subjects are naturally aligned in a shared anatomical coordinate system. This enables direct transfer: a model trained on one subject can be adapted to another without voxel-wise resampling, on the new subject's coordinates and responses. Unlike classical voxelwise models—which rigidly tie the representation to a subject-specific grid—NRF learns a continuous, anatomically grounded representation that could flexibly generalize across individuals with only minimal data. To exploit this property, we adopt a two-step adaptation strategy:

**Finetuning.** A pretrained NRF is fine-tuned on the new subject's limited data, using their MNI coordinates and measured responses. The two components of NRF  $\Phi$  are the feature extractor G and MLP predictor P. G encodes the visual stimulus into a representation, while P maps this representation and the spatial coordinate to the predicted brain response. Both G and P benefit from adaptation, since individuals vary in both how visual content is processed and how it is mapped to anatomy. Therefore, we perform full end-to-end finetuning of both components on the new subject's data.

**Voxelwise Ensemble.** To further improve performance on new subjects, we perform a voxel-wise ensemble of the predictions from different finetuned models. Similar to the personalized ensemble approach in Gu et al. (2022), this strategy maximizes predictive performance while preserving inter-subject variability to improve model personalization. Specifically, for each voxel v, let  $\{\hat{r}_v^{(1,i)}, \hat{r}_v^{(2,i)}, \dots, \hat{r}_v^{(K,i)}\}$  denote the predictions of K finetuned base models for the ith image. We then learn voxel-specific weights  $w_{v,k}$  (one per base model, for the k th base model) and a bias  $b_v$  by solving a least-squares regression on the limited new subject training data:

$$\min_{\{w_{v,k},b_v\}} \sum_{i=1}^{N} \left( r_v^{(i)} - \sum_{k=1}^{K} w_{v,k} \, \hat{r}_v^{(k,i)} - b_v \right)^2, \tag{1}$$

where  $r_v^{(i)}$  is the measured response of voxel v to image i, and N is the number of adaptation samples. At inference time, the final prediction for voxel v is given by the weighted ensemble:

$$\hat{r}_v = \sum_{k=1}^K w_{v,k} \, \hat{r}_v^{(k)} + b_v. \tag{2}$$

This ensemble leverages common neural structure while accounting for subject-specific variability, yielding higher accuracy in the low-data regime and producing a more personalized model for each subject.

#### 4 EXPERIMENTAL SETUP

# 4.1 DATASETS AND PREPROCESSING

We use the Natural Scene Dataset (NSD)Allen et al. (2022), which includes whole-brain 7T fMRI data from 8 subjects who viewed 10,000 natural scene images from the MS COCO dataset, repeated 1-3 times. The brain activations were computed using the GLMSingle algorithm Prince et al. (2022), and each voxel's response value is normalized per session ( $\mu=0,\sigma^2=1$ ). The brain activation to repeated images within a subject was averaged. The Neural response function(NRF) was trained using 9,000 unique images per subject, with around 1,000 images used for testing model accuracy via  $R^2$ . Since we are focusing on the visual cortex regions, we apply the official nsdgeneral region-of-interest (ROI) mask, which spans visual regions ranging from the early visual cortex to higher visual areas. Our evaluation focuses on Subj01, Subj02, Subj05, and Subj07 because these subjects completed all experiment sessions.

#### 4.2 EVALUATION METRICS.

To quantitatively compare with other models, we assess model performance across two levels.

**Voxel-Level Metrics:** To quantify prediction accuracy, we compute the voxel-wise Pearson correlation (Pearson) and voxel-wise mean square error (MSE) across all testing images.

**Semantic-Level Metrics:** Following prior work Bao et al. (2025), we evaluate the semantic fidelity of predicted responses using MindEye2 Scotti et al. (2024), a pretrained fMRI-to-image decoder. The decoder reconstructs visual stimuli from predicted fMRI responses, which are then compared against the ground-truth stimuli presented during data collection. We employ a suite of image reconstruction metrics: PixCorr and SSIM quantify low-level visual fidelity, while Alex(2) and Alex(5) capture feature similarity at early and deeper layers of AlexNet. To assess semantic alignment, we further compute Incep, CLIP, Eff, and SwAV scores, which measure the correspondence between reconstructed and original images in higher-level representational spaces. Additional details on these metrics are provided in the Appendix.

# 5 RESULTS

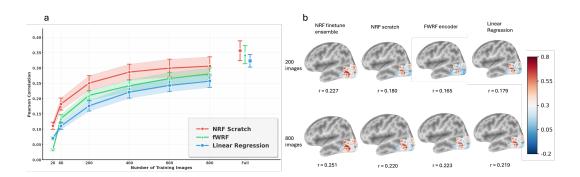


Figure 2: Prediction accuracy (Pearson correlation) in low data regime. a. **Single-subject models.** NRF consistently outperforms baseline models when trained on limited samples from scratch, highlighting the benefit of its continuous mapping. Results are shown for the average median voxel correlation across four subjects, with error bars indicating the standard error of the mean (SEM). b. **Cross-subject transfer.** Voxel-level prediction accuracy visualized on the cortical surface of subject 7. When pretrained base models from other subjects are available, the NRF finetune ensemble further improves performance over NRF scratch and baselines, showing clear gains across visual regions.

#### 5.1 Individual subjects encoding

We first evaluated NRF's neural prediction capability for single-subject data. Training a separate model for each of the 4 subjects and comparing the average neural prediction accuracy across subjects. For comparison, we selected two representative encoding models as baseline comparison. The linear regression model from the BrainDIVE Luo et al. (2023) and the fWRF (Feature-Weighted Receptive Field) encoder St-Yves & Naselaris (2018). Details about the baseline model are in A.3. We also took the result for full data encoding performance from MindSimulator Bao et al. (2025).

We first evaluate NRF under limited-data conditions, since practical applications rarely have access to the tens of thousands of trials collected in large-scale datasets such as NSD. As shown in Figure 2a, NRF achieves significantly higher accuracy than baseline models when trained on small numbers of images. Remarkably, with only 200 training samples, NRF outperforms baselines trained on more than 800 images. We attribute this data efficiency to the anatomical awareness of NRF: by conditioning on spatial coordinates, the model can exploit the smoothness of fMRI responses and learn more effectively from scarce data. This neuroscience-inspired design makes NRF particularly well-suited for realistic, low-data regimes. We then evaluate NRF in the full-data setting, where models were trained on ~9k images. The quantitative evaluation results are shown in Table 1.NRF outperforms baselines on voxelwise prediction metrics and achieves comparable performance on semantic-level evaluations. Figure 3 shows the visualization results of the decoded images. Low-level details and the semantic information of the images are very similar to the ground

truth. These results demonstrate that NRF maintains high voxel-level accuracy while also preserving semantic information, confirming its effectiveness across both limited and full data regimes.

In addition, we observed that fWRF achieves unusually high semantic-level scores, in some cases even surpassing reconstructions from measured fMRI. We attribute this to decoder bias: fWRF outputs, while less neurally accurate, may align more closely with the pretrained decoder's distribution, inflating semantic metrics. These results indicate that semantic-level metrics should be interpreted as a coarse indication of reconstruction quality rather than a strict basis for comparing encoding models.

Method	Voxel-l	Level			Semar	ntic-Level (v	ia decoding	g)		
	Pearson <sup>↑</sup>	MSE↓	PixCorr↑	SSIM↑	Alex(2)↑	Alex(5)↑	IncepT↑	CLIP↑	Eff↓	SwAV↓
Measured fMRI	-	_	0.322	0.431	96.1%	98.6%	95.4%	93.0%	0.619	0.344
Linear Regression	0.323	0.353	0.186	0.271	86.1%	95.0%	90.2%	84.5%	0.750	0.417
fWRF	0.343	0.361	0.303	0.341	96.9%	99.1%	96.2%	91.9%	0.614	0.356
MindSimulator (Trials=1)	0.345	0.403	0.194	0.296	89.0%	96.2%	92.3%	90.3%	0.702	0.399
MindSimulator (Trials=5)	0.355	0.385	0.201	0.298	89.6%	96.8%	93.2%	91.2%	0.688	0.393
NRF (our method)	0.358	0.345	0.261	0.371	91.6%	96.3%	92.1%	89.3%	0.706	0.400

Table 1: Evaluation results of fMRI prediction accuracy for the model trained on the full dataset. Our NRF model achieves state-of-the-art performance on prediction accuracy and comparable results on semantic metrics. All metrics are calculated across 4 subjects.



Figure 3: Visualization comparison between different neural encoding models and NRF. GT = seen during data collection. Measured fMRI = decoded image using measured fMRI. Reconstructions from NRF-predicted responses preserve both low-level visual details and high-level semantic content of the stimuli. Results shown for Subject 1.

#### 5.2 NEW SUBJECT ADAPTATION

More importantly, NRF enables cross-subject transfer, allowing knowledge learned from one subject to be adapted to new subjects—a critical property given that collecting fMRI data for new individuals is both resource-intensive and time-consuming. To evaluate this capability, we tested adaptation with 20, 200, and 800 images, corresponding to approximately 4, 40, and 160 minutes of scanning time. Three subjects were used for pretraining base models, and a fourth subject was held out for adaptation. For the new subject, we applied fine-tuning followed by voxelwise regression ensemble using the limited data. As a baseline, we compared against the "NRF scratch" approach, where a new NRF is trained entirely from the same limited dataset without pretraining. Across all data conditions, fine-tuning + ensemble consistently outperformed NRF scratch, confirming that NRF's

anatomically grounded formulation enables efficient cross-subject transfer, reducing the need for extensive subject-specific data while maintaining high predictive fidelity. The qualitative comparison is shown in Table 2. Predcition accuracy comparison across different methods is shown in Figure 2b. Notably, in the very low-data regime, finetuning + ensemble achieved strong semantic-level decoding performance. This shows that the strategy not only improves voxelwise prediction but also preserves subject variability, enabling predicted responses that more faithfully capture the semantic content of visual stimuli.

Training Images	Method	Voxel-I	Level	Semantic-Level (via decoding)								
		Pearson↑	MSE↓	PixCorr↑	SSIM↑	Alex(2)↑	Alex(5)↑	IncepT <sup>↑</sup>	CLIP↑	Eff↓	SwAV↓	
Full	NRF subject 7 (all data)	0.269	0.348	0.244	0.367	0.880	0.936	0.892	0.846	0.768	0.445	
20	NRF scratch NRF finetune ensemble	0.076 <b>0.114</b>	<b>0.417</b> 0.445	0.060 <b>0.186</b>	0.195 <b>0.366</b>	0.564 <b>0.750</b>	0.597 <b>0.792</b>	0.549 <b>0.732</b>	0.545 <b>0.729</b>	0.962 <b>0.868</b>	0.621 <b>0.515</b>	
200	NRF scratch NRF finetune ensemble	0.180 <b>0.227</b>	0.394 <b>0.390</b>	0.159 <b>0.255</b>	0.284 <b>0.372</b>	0.760 <b>0.908</b>	0.813 <b>0.957</b>	0.774 <b>0.913</b>	0.716 <b>0.873</b>	0.857 <b>0.729</b>	0.515 <b>0.425</b>	
800	NRF scratch NRF finetune ensemble	0.220 <b>0.251</b>	0.376 <b>0.372</b>	0.188 <b>0.269</b>	0.313 <b>0.382</b>	0.856 <b>0.927</b>	0.926 <b>0.970</b>	0.878 <b>0.922</b>	0.834 <b>0.895</b>	0.772 <b>0.700</b>	0.452 <b>0.408</b>	

Table 2: New subject adaptation with limited data (20, 200, 800 images). Here, the result is shown for using the NRF pretrained on subjects 1, 2, and 5 and adapted to subject 7 using finetuning + ensemble. This consistently outperforms training from scratch. With only 200 images, NRF outperforms the model trained with the whole dataset.

#### 5.3 PROBING ANATOMICAL AWARENESS

NRF achieves strong performance, particularly in low-data regimes, by leveraging two key properties of fMRI data: (i) the *local spatial continuity* of voxel responses within a subject, and (ii) the *anatomical alignment* across subjects. To directly test the contribution of these factors, we design controlled perturbation experiments that disrupt either spatial smoothness or cross-subject correspondence. If NRF's gains indeed stem from anatomical awareness, performance should degrade when these structural assumptions are broken.

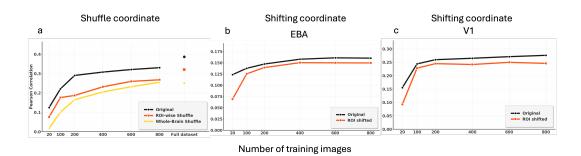


Figure 4: Probing anatomical awareness in NRF. (a) Disrupting spatial smoothness by shuffling coordinate-response pairings reduced accuracy, especially in low-data regimes, confirming that NRF relies on local continuity in brain responses. (b)(c) Breaking cross-subject alignment by shifting MNI coordinates degraded transfer, with the largest effect under limited data, showing that anatomical correspondence is critical for efficient adaptation.

**Local Smoothness.** To test whether NRF's data efficiency stems from exploiting spatial continuity, we disrupted the natural smoothness of fMRI data by shuffling coordinate–response pairings. Voxel responses were randomly reassigned to MNI coordinates, breaking correlations between neighboring voxels. We performed two variants of this perturbation: (i) *global shuffling*, randomizing pairings across the entire visual cortex, and (ii) *ROI-wise shuffling*, randomizing only within each ROI. Traditional voxelwise models should be unaffected, since they treat voxels independently. In contrast, NRF relies on coordinate conditioning, and as expected, its performance dropped sharply in low-data regimes, with global shuffling producing the largest drop. This confirms that NRF's

Method	Voxel-I	Level	Semantic-Level (via decoding)									
	Pearson <sup>†</sup>	MSE↓	PixCorr <sup>↑</sup>	SSIM↑	Alex(2)↑	Alex(5)↑	IncepT↑	CLIP↑	Eff↓	SwAV↓		
NRF finetune ensemble	0.227	0.390	0.255	0.372	0.908	0.957	91.3%	87.3%	0.729	0.425		
NRF finetune average	0.253	0.367	0.167	0.283	0.784	0.852	80.4%	74.9%	0.848	0.514		
NRF finetune base (subj1→subj7)	0.220	0.386	0.246	0.375	0.897	0.952	90.8%	86.9%	0.735	0.431		
NRF finetune base (subj2→subj7)	0.232	0.379	0.243	0.366	0.885	0.937	87.1%	82.4%	0.779	0.457		
NRF finetune base (subj5→subj7)	0.225	0.389	0.226	0.371	0.874	0.938	87.6%	82.4%	0.775	0.452		

Table 3: Ablation on voxelwise regression ensemble. We report the result for adapting from subjects 1,2,5 to subject 7 with 200 images.

improvements are driven by its ability to leverage local smoothness in brain responses. Shown in Figure 4(a).

Cross-Subject Alignment. To test the importance of anatomical correspondence for transfer, we disrupted MNI alignment by shifting voxel coordinates between subjects. Specifically, a model pretrained on Subject 1 was finetuned on Subject 7 using responses from EBA and V1. During finetuning, the MNI coordinates were shifted while remaining within the subject's brain range, breaking the cross-subject anatomy alignment. Compared to finetuning with aligned coordinates, coordinate shifting substantially degraded cross-subject transfer. The effect was most pronounced in low-data regimes: with only a small number of finetuning samples, the misaligned model failed to adapt, whereas alignment enabled effective transfer. With more data, the model gradually compensated for the misalignment, but still required far more samples to match the aligned case. These results demonstrate that NRF's cross-subject generalization depends critically on anatomical alignment. Without it, transfer is possible but far less data-efficient. To avoid artificial overlap after shifting, finetuning was performed using ROI-restricted data rather than the full brain.

#### 5.4 ABLATION STUDY

**Voxelwise ensemble** A key component for new subject adaptation is voxelwise regression ensemble, where each voxel is fit with a linear regression model to optimally combine predictions from multiple fine-tuned base models. This approach improves prediction accuracy while preserving subject-specific variability. Table 3 compares voxelwise regression against single fine-tuned base models and simple averaging. While simple averaging slightly boosts voxelwise prediction accuracy, it hinders subject variability and produces predicted fMRI signals with reduced semantic fidelity, leading to lower decoding performance. In contrast, voxelwise regression leverages complementary information across base models in a flexible, voxel-specific way, achieving both higher voxel-level accuracy and stronger semantic-level decoding results.

Additional ablation results are included in the appendix A.5.

## 6 Conclusion

In this work, we introduced the Neural Response Function (NRF), an anatomically aware neural encoding model that represents fMRI activity as a continuous function over MNI coordinates. Unlike conventional voxelwise models, NRF leverages spatial smoothness and cross-subject alignment to achieve accurate predictions in low-data regimes and to support efficient subject adaptation. Crucially, its continuous formulation moves beyond grid-locked voxels, allowing predictions at arbitrary spatial resolutions and across individuals. In this sense, NRF serves as a resolution-agnostic digital twin of the brain: a unified, flexible representation that integrates data across scales and subjects. These advances offer a new path toward efficient, generalizable, and anatomically grounded neural encoding.

# REFERENCES

Emily J Allen, Ghislain St-Yves, Yihan Wu, Jesse L Breedlove, Jacob S Prince, Logan T Dowdle, Matthias Nau, Brad Caron, Franco Pestilli, Ian Charest, et al. A massive 7t fmri dataset to bridge cognitive neuroscience and artificial intelligence. *Nature neuroscience*, 25(1):116–126, 2022.

- Guangyin Bao, Qi Zhang, Zixuan Gong, Zhuojia Wu, and Duoqian Miao. Mindsimulator: Exploring brain concept localization via synthetic fmri. arXiv preprint arXiv:2503.02351, 2025.
  - Roman Beliy, Guy Gaziv, Assaf Hoogi, Francesca Strappini, Tal Golan, and Michal Irani. From voxels to pixels and back: Self-supervision in natural-image reconstruction from fmri. *Advances in Neural Information Processing Systems*, 32, 2019.
    - Zhiqin Chen and Hao Zhang. Learning implicit fields for generative shape modeling. In *Proceedings* of the IEEE/CVF conference on computer vision and pattern recognition, pp. 5939–5948, 2019.
    - Paul E Downing, Yuhong Jiang, Miles Shuman, and Nancy Kanwisher. A cortical area selective for visual processing of the human body. *Science*, 293(5539):2470–2473, 2001.
    - Russell Epstein and Nancy Kanwisher. A cortical representation of the local visual environment. *Nature*, 392(6676):598–601, 1998.
    - Zijin Gu, Keith Jamison, Mert Sabuncu, and Amy Kuceyeski. Personalized visual encoding model construction with small data. *Communications Biology*, 5(1):1382, 2022.
    - Umut Güçlü and Marcel AJ Van Gerven. Deep neural networks reveal a gradient in the complexity of neural representations across the ventral stream. *Journal of Neuroscience*, 35(27):10005–10014, 2015.
    - Kuan Han, Haiguang Wen, Junxing Shi, Kun-Han Lu, Yizhen Zhang, Di Fu, and Zhongming Liu. Variational autoencoder: An unsupervised model for encoding and decoding fmri activity in visual cortex. *NeuroImage*, 198:125–136, 2019.
    - David J Heeger and David Ress. What does fmri tell us about neuronal activity? *Nature reviews neuroscience*, 3(2):142–151, 2002.
    - Alexander G Huth, Shinji Nishimoto, An T Vu, and Jack L Gallant. A continuous semantic space describes the representation of thousands of object and action categories across the human brain. *Neuron*, 76(6):1210–1224, 2012.
    - Alexander G Huth, Wendy A De Heer, Thomas L Griffiths, Frédéric E Theunissen, and Jack L Gallant. Natural speech reveals the semantic maps that tile human cerebral cortex. *Nature*, 532 (7600):453–458, 2016.
    - Nancy Kanwisher, Josh McDermott, and Marvin M Chun. The fusiform face area: a module in human extrastriate cortex specialized for face perception. *Journal of neuroscience*, 17(11):4302–4311, 1997.
    - Kendrick N Kay, Thomas Naselaris, Ryan J Prenger, and Jack L Gallant. Identifying natural images from human brain activity. *Nature*, 452(7185):352–355, 2008.
    - Andrew Luo, Maggie Henderson, Leila Wehbe, and Michael Tarr. Brain diffusion for visual exploration: Cortical discovery using large scale generative models. *Advances in Neural Information Processing Systems*, 36:75740–75781, 2023.
    - Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3d reconstruction in function space. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 4460–4470, 2019.
    - Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*, 65(1):99–106, 2021.
- Tom M Mitchell, Svetlana V Shinkareva, Andrew Carlson, Kai-Min Chang, Vicente L Malave, Robert A Mason, and Marcel Adam Just. Predicting human brain activity associated with the meanings of nouns. *science*, 320(5880):1191–1195, 2008.
  - Thomas Naselaris, Kendrick N Kay, Shinji Nishimoto, and Jack L Gallant. Encoding and decoding in fmri. *Neuroimage*, 56(2):400–410, 2011.

- Thomas Naselaris, Cheryl A Olman, Dustin E Stansbury, Kamil Ugurbil, and Jack L Gallant. A voxel-wise encoding model for early visual areas decodes mental images of remembered scenes. *Neuroimage*, 105:215–228, 2015.
- Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. Deepsdf: Learning continuous signed distance functions for shape representation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 165–174, 2019.
- Jacob S Prince, Ian Charest, Jan W Kurzawski, John A Pyles, Michael J Tarr, and Kendrick N Kay. Improving the accuracy of single-trial fmri response estimates using glmsingle. *Elife*, 11:e77599, 2022.
- Paul S Scotti, Mihir Tripathy, Cesar Kadir Torrico Villanueva, Reese Kneeland, Tong Chen, Ashutosh Narang, Charan Santhirasegaran, Jonathan Xu, Thomas Naselaris, Kenneth A Norman, et al. Mindeye2: Shared-subject models enable fmri-to-image with 1 hour of data. *arXiv preprint arXiv:2403.11207*, 2024.
- Vincent Sitzmann, Julien Martel, Alexander Bergman, David Lindell, and Gordon Wetzstein. Implicit neural representations with periodic activation functions. *Advances in neural information processing systems*, 33:7462–7473, 2020.
- Ghislain St-Yves and Thomas Naselaris. The feature-weighted receptive field: an interpretable encoding model for complex feature spaces. *NeuroImage*, 180:188–202, 2018.
- Matthew Tancik, Pratul Srinivasan, Ben Mildenhall, Sara Fridovich-Keil, Nithin Raghavan, Utkarsh Singhal, Ravi Ramamoorthi, Jonathan Barron, and Ren Ng. Fourier features let networks learn high frequency functions in low dimensional domains. *Advances in neural information processing systems*, 33:7537–7547, 2020.
- Jerry Tang, Meng Du, Vy Vo, Vasudev Lal, and Alexander Huth. Brain encoding models based on multimodal transformers can transfer across language and vision. *Advances in Neural Information Processing Systems*, 36:29654–29666, 2023.
- Aria Y Wang, Kendrick Kay, Thomas Naselaris, Michael J Tarr, and Leila Wehbe. Better models of human high-level visual cortex emerge from natural language supervision with a large and diverse dataset. *Nature Machine Intelligence*, 5(12):1415–1426, 2023.
- Daniel LK Yamins, Ha Hong, Charles F Cadieu, Ethan A Solomon, Darren Seibert, and James J DiCarlo. Performance-optimized hierarchical models predict neural responses in higher visual cortex. *Proceedings of the national academy of sciences*, 111(23):8619–8624, 2014.

# A APPENDIX

## A.1 USE OF LLMS

LLMs were used only for polishing grammar and writing clarity.

#### A.2 ADDITIONAL DETAILS ON NRF

Image Feature Extraction Block. We leverage the pretrained OpenAI CLIP ViT-B/16 to obtain multiscale image features. Representations are extracted from the 3rd and 6th transformer layers, each yielding features of shape  $(196\times768)$ , along with the final CLIP embedding of shape  $(1\times512)$ . The two intermediate feature maps are each processed by separate two-layer projection modules with identical architecture: the first layer reduces the dimensionality from  $(196\times768)$  to  $(196\times256)$ , and the second compresses this to a  $(1\times256)$  vector. These two compact embeddings are then concatenated with the  $(1\times512)$  CLIP embedding to form the fused multiscale image representation G(M).

**MLP Predictor.** The predictor is a coordinate-conditioned MLP that takes as input both the Fourier positional encoding of the MNI coordinate and the fused image embedding G(M). It outputs a single scalar—the predicted fMRI response at that voxel location. We use an 8-layer MLP with hidden dimension 4096, applying ReLU activations after each layer except the final output. Further ablations on model architecture are included in the appendix A.5.

#### A.3 ADDITIONAL DETAILS ON BASELINE MODELS

**FWRF encoding model** The encoder uses AlexNet as the base feature extractor, processing 227×227 pixel input images normalized to [0,1]. Feature selection retains the top 256 features per layer based on variance across the training data. The receptive field model employs a 3×3 spatial grid with aperture size 0.8, covering RF sizes from 0.15 to 0.25 across 2 logarithmically-spaced scales, yielding 18 total RF candidates per voxel. Ridge regression optimization uses regularization parameters  $\lambda \in [10^4, 10^5]$  with adaptive holdout validation.

**Linear Regression Encoding.** For linear regression baseline we used the same encoding model as Luo et al. (2023). Specifically, we extract the  $(1 \times 512)$  CLIP embedding from OpenAI CLIP ViT-B/16 and directly map it to the voxel dimension (e.g., 15,724 voxels) using a linear layer. The model is trained for 150 epochs with the AdamW optimizer, with a learning rate that decays linearly from  $3 \times 10^{-4}$ . During inference, we select the checkpoint that achieves the lowest validation MSE.

#### A.4 ADDITIONAL DETAILS ON EVALUATION METRICS

We used the evaluation metrics for decoded image evaluation from MindEye2 Scotti et al. (2024) directly. PixCorr measures the pixel-wise correlation between the ground-truth image and the reconstruction. SSIM refers to the Structural Similarity Index, which evaluates perceptual similarity between ground-truth and reconstructed images. Alex(2), Alex(5), Incep, and CLIP are two-way identification metrics (chance = 50%) based on feature similarity. Specifically, Alex(2) uses features from the 2nd layer of AlexNet, Alex(5) from the 5th layer of AlexNet, Incep from the final pooling layer of InceptionV3, and CLIP from the final layer of CLIP ViT-L/14. In two-way identification, the task is to decide whether the voxel embedding is closer to its paired image embedding or to a randomly selected image embedding, reported as percent correct. Eff and SwAV denote representational similarity metrics, computed as the average correlation distance between voxel embeddings and features extracted from EfficientNet-B1 and SwAV-ResNet50, respectively.

## A.5 ADDITIONAL ABLATION RESULTS

**Finetuning strategy** During finetuning, there are two options: the projection model and the image feature merger. We explored finetuning both/projector-only and the feature extraction block only. We finetuned subject 1 NRF with 800 images from subject 8 data. The results are shown in the Table A.5. The result suggests that finetuning should be done on both components to get the maximum performance boost.

Finetune strategy	Voxel-I	Level		Semantic-Level (via decoding)											
	Pearson <sup>↑</sup>	Pearson↑ MSE↓ Pix		SSIM↑	Alex(2)↑	Alex(5)↑	IncepT↑	CLIP↑	Eff↓	SwAV↓					
Both	0.234	0.361	0.254	0.377	0.910	0.958	90.9%	86.7%	0.730	0.424					
Image extractor only	0.104	0.394	0.102	0.274	0.613	0.636	59.2%	55.6%	0.958	0.596					
Projector only	0.233			0.361	0.894	0.945	88.8%	83.8%	0.766	0.446					

Table 4: Ablation of different finetuning strategies. Models pre-trained on subject 1 finetuned on 800 images from subject 7.

**Number of layers** To investigate the effect of model architecture on neural response prediction accuracy, we conducted an ablation study by varying the number of layers and the hidden dimension of the MLP projector. As shown in Table 5, we computed 4, 8, 16-layer configurations under subject-specific settings on subject1. The results show that 8-layer model results in the best performance, which is also the setting we utilized for our experiments.

Layers	Voxel-I	Level		Semantic-Level (via decoding)										
	Pearson <sup>↑</sup>	MSE↓	PixCorr↑	SSIM↑	Alex(2)↑	Alex(5)↑	IncepT↑	CLIP↑	Eff↓	SwAV↓				
4	0.358	0.348	0.258	0.351	0.914	0.961	89.7%	85.1%	0.748	0.437				
8	0.360	0.350	0.324	0.387	0.956	0.983	94.1%	89.9%	0.680	0.396				
16	0.349	0.353	0.331	0.385	0.956	0.983	93.99%	89.9%	0.678	0.395				

Table 5: Ablation of the number of layers of the MLP projector. Model trained on subject 1 data.

**Hidden dimension** We also conducted an ablation study by varying the hidden dimension of the MLP projector. As shown in Table, we computed for hidden dimension = 2048, 4096, 8192 configurations under subject-specific settings on subject1. The results show that model with hidden dimension = 4096 results in the best performance, which is also the setting we utilized for our experiments.

Hidden dimension	Voxel-I	Level		Semantic-Level (via decoding)										
	Pearson <sup>†</sup>	Pearson↑ MSE↓ I		SSIM↑	Alex(2)↑	Alex(5)↑	IncepT↑	CLIP↑	Eff↓	SwAV↓				
2048	0.358	0.348	0.253	0.353	0.903	0.948	87.4%	83.2%	0.777	0.451				
4096	0.360	0.350	0.324	0.387	0.956	0.983	94.1%	89.9%	0.680	0.396				
8192	0.353	0.353	0.323	0.386	0.955	0.981	94.1%	89.9%	0.683	0.396				

Table 6: Ablation on hidden layer dimension of the MLP projector. Model trained on subject 1 data.

## A.6 ADDITIONAL SUBJECT ADAPTATION RESULT

Here, we present additional results of new subject adaptation for subjects 1, 2, and 5 in Table 7, Table 8, and Table 9, respectively. The results show that our method consistently yields superior performance compared to the scratch method.

Training Images	Method	Voxel-I	Level	Semantic-Level (via decoding)									
		Pearson↑	MSE↓	PixCorr <sup>↑</sup>	SSIM↑	Alex(2)↑	Alex(5)↑	IncepT↑	CLIP↑	Eff↓	SwAV↓		
20	NRF scratch	0.116	0.411	0.023	0.163	0.548	0.552	56.8%	53.9%	0.969	0.660		
	NRF finetune ensemble	0.184	0.463	0.139	0.308	0.744	0.809	72.5%	68.9%	0.885	0.540		
200	NRF scratch	0.261	0.377	0.132	0.242	0.750	0.811	73.6%	70.3%	0.892	0.555		
	NRF finetune ensemble	0.306	0.379	0.266	0.375	0.917	0.958	88.2%	85.9%	0.758	0.437		
800	NRF scratch	0.314	0.369	0.244	0.307	0.915	0.962	90.6%	86.1%	0.742	0.432		
	NRF finetune ensemble	0.342	0.361	0.316	0.382	0.945	0.980	93.3%	88.7%	0.698	0.404		

Table 7: New subject adaptation with limited data (20, 200, 800 images). NRF pretrained on subjects 2,5,7 are used as base models to adapt to subject 1.

# **B** ETHIC STATEMENT

Our research adheres to the ICLR Code of Ethics. All experiments in this paper are conducted using open-source datasets, and no potential ethical concerns are associated with this work.

# C REPRODUCIBILITY STATEMENT

All preprocessed data, code, and model parameters used in our research will be made publicly available upon publication. Detailed protocols for data preprocessing, model training, and evaluation have been provided in our manuscript, enabling independent reproduction.

Training Images	Method	Voxel-I	Level	Semantic-Level (via decoding)								
		Pearson↑	MSE↓	PixCorr↑	SSIM↑	Alex(2)↑	Alex(5)↑	IncepT↑	CLIP↑	Eff↓	SwAV↓	
20	NRF scratch	0.124	0.462	0.023	0.344	0.499	0.498	49.6%	49.6%	0.971	0.636	
	NRF finetune ensemble	0.168	0.457	0.140	0.328	0.780	0.848	74.0%	67.9%	0.874	0.532	
200	NRF scratch	0.266	0.386	0.076	0.323	0.617	0.674	62.7%	57.8%	0.944	0.605	
	NRF finetune ensemble	0.317	0.375	0.255	0.365	0.916	0.966	90.3%	85.5%	0.735	0.427	
800	NRF scratch	0.323	0.372	0.192	0.299	0.885	0.951	87.4%	82.5%	0.778	0.452	
	NRF finetune ensemble	0.356	0.354	0.274	0.374	0.935	0.976	92.4%	88.0%	0.711	0.413	

Table 8: New subject adaptation with limited data (20, 200, 800 images). NRF pretrained on subjects 1,5,7 are used as base models to adapt to subject 2.

Training Images	Method	Voxel-I	Level	Semantic-Level (via decoding)								
		Pearson↑	MSE↓	PixCorr↑	SSIM↑	Alex(2)↑	Alex(5)↑	IncepT↑	CLIP↑	Eff↓	SwAV↓	
20	NRF scratch	0.128	0.433	0.126	0.213	0.608	0.633	56.6%	56.8%	0.953	0.592	
	NRF finetune ensemble	0.184	0.470	0.161	0.366	0.695	0.744	70.2%	67.3%	0.889	0.533	
200	NRF scratch	0.293	0.415	0.113	0.238	0.687	0.719	68.0%	64.3%	0.911	0.573	
	NRF finetune ensemble	0.353	0.373	0.242	0.374	0.882	0.936	88.4%	84.8%	0.765	0.445	
800	NRF scratch	0.365	0.370	0.211	0.326	0.883	0.945	89.3%	87.0%	0.734	0.426	
	NRF finetune ensemble	0.400	0.355	0.259	0.383	0.916	0.969	92.8%	90.9%	0.682	0.401	

Table 9: New subject adaptation with limited data (20, 200, 800 images). NRF pretrained on subjects 1,2,7 are used as base models to adapt to subject 5.