

# ALIGNING BRAIN FUNCTIONS BOOSTS THE DECODING OF VIDEOS IN NOVEL SUBJECTS

Anonymous authors

Paper under double-blind review

## ABSTRACT

1 Deep learning is leading to major advances in the realm of brain decoding from  
 2 functional Magnetic Resonance Imaging (fMRI). However, the large inter-subject  
 3 variability in brain characteristics has limited most studies to train models on one  
 4 subject at a time. Consequently, this approach hampers the training of deep learning  
 5 models, which typically requires very large datasets. Here, we propose to  
 6 boost brain decoding by aligning brain responses to videos across subjects. Compared  
 7 to the anatomically-aligned baseline, our method improves out-of-subject  
 8 decoding performance by up to 75%. Moreover, it also outperforms classical  
 9 single-subject approaches when fewer than 100 minutes of data is available for the  
 10 tested subject. Furthermore, we propose a new multi-subject alignment method,  
 11 which obtains comparable results to that of classical single-subject approaches  
 12 while easing out-of-subject generalization. Finally, we show that this method  
 13 aligns neural representations in accordance with brain anatomy. Overall, this study  
 14 lays foundations to leverage extensive neuroimaging datasets and enhance the  
 15 decoding of individuals with a limited amount of brain recordings.

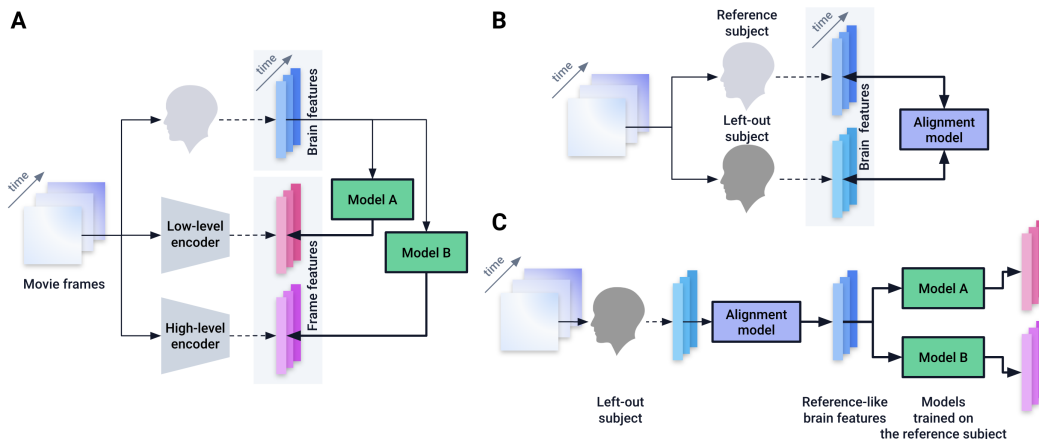


Figure 1: **General outline of video decoding from BOLD fMRI signal in left-out subjects**

**A.** For every image associated with a brain volume, one computes its low-level and high-level latent representations using pre-trained models. Subsequently, regression models can be fitted to map brain features onto each of these latent representations. **B.** BOLD signal acquired in two subjects watching the same movie can be used to derive an alignment model which associates voxels of the two subjects based on functional similarity. **C.** Once this alignment model is trained, it can be used to transform brain features of the left-out subject into brain features that resemble that of the reference subject. In particular, this allows one to use models that have been trained on a lot of data coming from a reference subject data, and apply it on a left-out subject for whom less data was collected.

## 16 1 INTRODUCTION

17 **Decoding the brain** Deep learning is greatly accelerating the possibility of decoding mental repre-  
18 sentations from brain activity. Originally restricted to linear models (Mitchell et al., 2004; Harrison  
19 & Tong, 2009; Haynes & Rees, 2006), the decoding of brain activity can now be carried out with  
20 deep learning techniques. In particular, using functional Magnetic Resonance Imaging (fMRI) sig-  
21 nals, significant progress has been made in the decoding of images (Ozcelik & VanRullen, 2023;  
22 Chen et al., 2023a; Scotti et al., 2023; Takagi & Nishimoto, 2023; Gu et al., 2023; Ferrante et al.,  
23 2023; Mai & Zhang, 2023), speech (Tang et al., 2023), and videos (Kupersmidt et al., 2022; Wen  
24 et al., 2018; Wang et al., 2022; Chen et al., 2023b; Lahner et al., 2023; Phillips et al., 2022).

25 **Challenge** However, brain representations are highly variable across subjects, which makes it  
26 challenging to train the same model on multiple subjects. Therefore, with few noteworthy exceptions  
27 (Haxby et al., 2020; Ho et al., 2023), studies typically train a decoder on a single subject at a time.  
28 With this constraint in mind, major effort has been put towards building fMRI datasets collecting  
29 a lot of data in a limited number of participants (Allen et al., 2022; Wen et al., 2017; LeBel et al.,  
30 2023; Pinho et al., 2018). Nonetheless, the necessity to train and test models on a single subject  
31 constitutes a major impediment to using notoriously data-hungry deep learning approaches.

32 **Functional alignment** Several methods can align the functions – as opposed to the anatomy – of  
33 multiple brains, and thus offer a potential solution to inter-subject variability: differentiable wrap-  
34 pings of the cortical surface (Robinson et al., 2014), rotations between brain voxels in the functional  
35 space (Haxby et al., 2011), shared response models (Chen et al., 2015; Richard et al., 2020), per-  
36 mutations of voxels minimizing an optimal transport cost (Bazeille et al., 2019), or combinations  
37 of these approaches (Feilong et al., 2022). However, it is not clear which of these methods offers  
38 the best performance and generalization capabilities (Bazeille et al., 2021). Besides, several studies  
39 rely on deep learning models trained in a self-supervised fashion to obtain a useful embedding of  
40 brain activity, in hope that this embedding could be meaningful across subjects (Thomas et al., 2022;  
41 Chen et al., 2023a). However, it is currently unknown whether any of these methods improve the  
42 decoding of naturalistic stimuli such as videos, and how such hypothetical gain would vary with the  
43 amount of fMRI recording available in a given a subject.

44 **Approach** To address this issue, we leverage fMRI recordings of multiple subjects to boost the  
45 decoding of videos in a single left-out subject. This requires fitting two models: an alignment model  
46 and a decoder. The alignment aims at making brain responses of a left-out subject most similar to  
47 those of a reference subject. Here, we leverage optimal transport to compute this transformation  
48 using functional and anatomical data from both subjects. The decoder consists of a linear regression  
49 trained to predict the latent representations of movie frames from the corresponding BOLD signals.

50 We evaluate video decoding in different setups. In particular, we assess (1) whether training a de-  
51 coder with several subjects improves performance, (2) whether decoders generalize to subjects on  
52 which they were not trained and (3) the extent to which functional alignment improves aforemen-  
53 tioned setups.

54 **Contributions** We first confirm the feasibility of decoding, from 3T fMRI, the semantics of videos  
55 watched by the subjects. Our study further makes three novel contributions:

- 56 1. functional alignment across subjects boosts video decoding performance when left-out sub-  
57 jects have a limited amount of data
- 58 2. training a decoder on multiple aligned subjects reaches the same performance as training a  
59 single model per subject
- 60 3. the resulting alignments, computed from movie watching data, yield anatomically-coherent  
61 maps.

62 From a representation learning perspective, this is one more piece of evidence that representations  
63 learnt by deep learning models can help model and decode brain signal, even with stimuli as complex  
64 as naturalistic videos. Our results also show that, in high-data regimes, naturalistic movie-watching  
65 yields functional features which can help discriminate between parts of the cortex much beyond the  
66 visual system.

## 67 2 METHODS

68 Our goal is to decode visual stimuli seen by subjects from their brain activity. To this end, we train  
69 a linear model to predict latent representations – shortened as *latents* – of these visual stimuli from  
70 BOLD fMRI signals recorded in subjects watching naturalistic videos.

71 In the considered data, brains are typically imaged at a rate of one scan every 2 seconds. During this  
72 period, a subject sees 60 video frames on average. For simplicity, we consider the restricted issue of  
73 decoding only the first video frame seen by subjects at each brain scan. Formally, for a given subject,  
74 let  $\mathbf{X} \in \mathbb{R}^{n,v}$  be the BOLD response collected in  $v$  voxels over  $n$  brain scans and  $\mathbf{Y} \in \mathbb{R}^{n,m}$  the  
75  $m$ -dimensional latent representation of each selected video frame for all  $n$  brain scans.

## 76 2.1 BRAIN ALIGNMENT

77 **Anatomical alignment** As a baseline, we consider the alignment method implemented in  
78 Freesurfer (Fischl, 2012), which relies on anatomical information to project each subject onto a  
79 surface template of the cortex (in our case *fsaverage5*). Consequently, brain data from all subjects  
80 lie on a mesh of size  $v = 10\,242$  vertices per hemisphere.

81 **Functional alignment** On top of the aforementioned anatomical alignment, we apply a recent  
82 method from Thual et al. (2022) denoted as Fused Unbalanced Gromov-Wasserstein (FUGW)<sup>1</sup>.  
83 As illustrated in Figure 1.B, this method consists in using functional data to train an alignment that  
84 transforms brain responses of a given left-out subject into the brain responses of a reference subject.  
85 This approach can be seen as a soft permutation of voxels<sup>2</sup> of the left-out subject which maximizes  
86 the functional similarity to voxels of the reference subject.

87 Formally, for a left-out subject, let  $\mathbf{D}^{\text{out}} \in \mathbb{R}^{v,v}$  be the matrix of anatomical distances between  
88 vertices on the cortex, and  $\mathbf{w}^{\text{out}} \in \mathbb{R}_+^v$  a probability distribution on vertices.  $\mathbf{w}^{\text{out}}$  can be interpreted  
89 as the relative importance of vertices ; without prior knowledge, we use the uniform distribution.  
90 Reciprocally, we define  $\mathbf{D}^{\text{ref}}$  and  $\mathbf{w}^{\text{ref}}$  for a reference subject. Note that, in the general case,  $v$   
91 can be different from one subject to the other, although we simplify notations here.

92 We derive a transport plan  $\mathbf{P} \in \mathbb{R}^{v,v}$  to match the vertices of the two subjects based on functional  
93 similarity, while preserving anatomical organisation. For this, we simultaneously optimize multiple  
94 constraints, formulated in the loss function  $\mathcal{L}(\mathbf{P})$  described in Equation 1:

$$\begin{aligned} \mathcal{L}(\mathbf{P}) \triangleq & (1 - \alpha) \underbrace{\sum_{0 \leq i,j < n} \|\mathbf{X}_i^{\text{out}} - \mathbf{X}_j^{\text{ref}}\|_2^2 \mathbf{P}_{i,j}}_{\text{Wasserstein loss}} + \alpha \underbrace{\sum_{0 \leq i,k,j,l < n} |\mathbf{D}_{i,k}^{\text{out}} - \mathbf{D}_{j,l}^{\text{ref}}|^2 \mathbf{P}_{i,j} \mathbf{P}_{k,l}}_{\text{Gromov-Wasserstein loss}} \\ & + \rho \left( \underbrace{\text{KL}(\mathbf{P}_{\#1} \otimes \mathbf{P}_{\#1} | \mathbf{w}^{\text{out}} \otimes \mathbf{w}^{\text{out}})}_{\text{Marginal constraints}} + \underbrace{\text{KL}(\mathbf{P}_{\#2} \otimes \mathbf{P}_{\#2} | \mathbf{w}^{\text{ref}} \otimes \mathbf{w}^{\text{ref}})}_{\text{Entropy}} \right) + \varepsilon \text{H}(\mathbf{P}) \quad (1) \end{aligned}$$

95

96 with  $\mathbf{P}_{\#1} \triangleq (\sum_j \mathbf{P}_{i,j})_{0 \leq i < v}$  and  $\mathbf{P}_{\#2} \triangleq (\sum_i \mathbf{P}_{i,j})_{0 \leq j < v}$  the first and second marginal distri-  
97 butions of  $\mathbf{P}$ ,  $\otimes$  the Kronecker product between two matrices, and  $\text{KL}(\cdot, \cdot)$  the Kullback-Leibler  
98 divergence.  $\alpha$ ,  $\rho$  and  $\varepsilon$  are hyper-parameters setting the relative importance of each constraint.

99 Following Thual et al. (2022), we minimize  $\mathcal{L}(\mathbf{P})$  with 10 iterations of a block coordinate descent  
100 algorithm (Séjourné et al., 2021), each running 1 000 Sinkhorn iterations (Cuturi, 2013). Subse-  
101 quently, we define  $\phi_{\text{out} \rightarrow \text{ref}}: \mathbf{X} \mapsto (\mathbf{P}^T \mathbf{X}^T) \oslash \mathbf{P}_{\#2} \in \mathbb{R}^{n,v}$  where  $\oslash$  is the element-wise di-  
102 vision, a function which transports any matrix of brain features from the left-out subject to the  
103 reference subject. To simplify notations, for any  $\mathbf{X}$  defined on the left-out subject, we define  
104  $\mathbf{X}^{\text{out} \rightarrow \text{ref}} \triangleq \phi_{\text{out} \rightarrow \text{ref}}(\mathbf{X})$ .

<sup>1</sup><https://alexisthual.github.io/fugw><sup>2</sup>We use the words *voxel* (volumetric pixel) or *vertex* (point on a mesh) indifferently.

## 105 2.2 DECODING

106 **Brain input** There is a time *lag* between the moment a stimulus is played and the moment it elicits  
 107 a maximal BOLD response in the brain (Glover, 1999). Moreover, since the effect induced by this  
 108 stimulus might span over multiple consecutive brain volumes, we set a *window size* describing the  
 109 number of brain volumes to aggregate together. To account for these effects, we use a standard  
 110 Finite Impulse Response (FIR) approach. FIR consists in fitting the decoder on a time-lagged,  
 111 multi-volume version of the BOLD response. Different *aggregation functions* can be used, such as  
 112 stacking or averaging. Figure S2 describes these concepts visually.

113 **Video output** The matrix of latent features  $\mathbf{Y}$  is obtained by using a pre-trained image encoder on  
 114 each video frame and concatenating all obtained vectors in  $\mathbf{Y}$ . Similarly to Ozcelik & VanRullen  
 115 (2023), and as illustrated in Figure 1.A, we seek to predict CLIP  $257 \times 768$  (high-level) and VD-  
 116 VAE (low-level) latent representations. We use visual – as opposed to textual – CLIP representations  
 117 (Radford et al., 2021). For comparison, we also reproduce our approach on latent representations  
 118 from CLIP CLS (high-level) and AutoKL (low-level), which happen to be much smaller<sup>3</sup> and might  
 119 be computationally easier to fit.

120 **Model** Fitting the decoder consists in deriving  $\mathbf{W} \in \mathbb{R}^{v,m}$ ,  $\mathbf{b} \in \mathbb{R}^m$  the solution of a Ridge  
 121 regression problem – i.e. a linear regression with L2 regularization – predicting  $\mathbf{Y}$  from  $\mathbf{X}$ .

122 **Evaluation** We evaluate the performance of the decoder with retrieval metrics. Let us denote  $\mathbf{X}$   
 123 and  $\mathbf{Y}$  the brain and latent features used to train the decoder,  $\mathbf{X}_{\text{test}}$  and  $\mathbf{Y}_{\text{test}}$  those to test the decoder,  
 124 and  $\hat{\mathbf{Y}} \triangleq \mathbf{W}\mathbf{X}_{\text{test}} + \mathbf{b}$  the predicted latents. We ensure that the train and test data are disjoint.

125 We randomly draw a retrieval set  $K$  of 499 frames without replacement from the test data. For each  
 126 pair  $\hat{\mathbf{y}}, \mathbf{y}$  of predicted and ground truth latents, one derives their cosine similarity score  $s(\hat{\mathbf{y}}, \mathbf{y})$ , as  
 127 well as similarity scores to all latents  $\mathbf{y}_{\text{neg}}$  of the retrieval set  $s(\hat{\mathbf{y}}, \mathbf{y}_{\text{neg}})$ . Let us denote  $r(\hat{\mathbf{y}}, \mathbf{y})$  the  
 128 rank of  $\mathbf{y}$ , which we define as the number of elements of  $K$  whose similarity score to  $\hat{\mathbf{y}}$  is larger  
 129 than  $s(\hat{\mathbf{y}}, \mathbf{y})$ . In order for the rank to not depend on the size of  $K$ , we define the *relative rank* as  
 130  $\frac{r(\hat{\mathbf{y}}, \mathbf{y})}{|K|}$ . Eventually, one derives the median relative rank  $\text{MR}(\hat{\mathbf{Y}}, K)$ :

$$r(\hat{\mathbf{y}}, \mathbf{y}) \triangleq |\{\mathbf{y}_{\text{neg}} \in K \mid s(\hat{\mathbf{y}}, \mathbf{y}_{\text{neg}}) > s(\hat{\mathbf{y}}, \mathbf{y})\}|$$

$$\text{MR}(\hat{\mathbf{Y}}, K) \triangleq \text{median}\left(\left\{\frac{r(\hat{\mathbf{y}}, \mathbf{y})}{|K|}, \forall (\hat{\mathbf{y}}, \mathbf{y})\right\}\right)$$

## 131 2.3 DECODING AND ALIGNMENT SETUPS

132 **Within- vs out-of-subject** The *within-subject* setup consists in training a decoder with data  $\mathbf{X}_{\text{train}}^{S_1}$ ,  
 133  $\mathbf{Y}_{\text{train}}^{S_1}$  from a given subject, and testing it on left-out data  $\mathbf{X}_{\text{test}}^{S_1}$ ,  $\mathbf{Y}_{\text{test}}^{S_1}$  acquired in the same subject.  
 134 The *out-of-subject* setup consists in training a decoder with data from a given subject, and testing it  
 135 on data  $\mathbf{X}_{\text{test}}^{S_2}$ ,  $\mathbf{Y}_{\text{test}}^{S_2}$  acquired in a left-out subject.

136 **Single- vs multi-subject** The *single-subject* setup consists in training a decoder predicting  $\mathbf{Y}$  from  
 137  $\mathbf{X}$  for each subject. The *multi-subject* setup consists in training a single decoder using data from  
 138 multiple subjects. In this case, data from several subjects is stacked together, resulting in a matrix  
 139  $\mathbf{X}_{\text{multi}} \in \mathbb{R}^{n_1 + \dots + n_p, v}$  and  $\mathbf{Y}_{\text{multi}} \in \mathbb{R}^{n_1 + \dots + n_p, m}$ , where  $p$  is the number of subjects.

140 **Un-aligned vs aligned** In multi-subject and out-of-subject setups, data coming from different sub-  
 141 jects can be *aligned* to a *reference* subject. Let us assume that  $S_1$  is the reference subject. In the  
 142 case of multi-subject, all subjects are aligned to  $S_1$  and the decoder is trained on a concatenation of  
 143  $\mathbf{X}^{S_1}, \mathbf{X}^{S_2 \rightarrow S_1}, \dots, \mathbf{X}^{S_p \rightarrow S_1}$  (see notations introduced at the end of section 2.1) and  $\mathbf{Y}^{S_1}, \dots, \mathbf{Y}^{S_p}$ ,  
 144 where  $p$  is the number of subjects. In the case of out-of-subject, it corresponds to aligning  $S_2$  onto  
 145  $S_1$ , such that a decoder trained on  $S_1$  will be tested on  $\mathbf{X}_{\text{test}}^{S_2 \rightarrow S_1}, \mathbf{Y}_{\text{test}}^{S_2}$ .

<sup>3</sup>Dimensions for CLIP CLS: 768 ; CLIP  $257 \times 768$  :  $257 \times 768 = 197\,376$  ; AutoKL:  $4 \times 32 \times 32 = 4\,096$   
 ; VD-VAE:  $2 \times 2^4 + 4 \times 2^8 + 8 \times 2^{10} + 16 \times 2^{12} + 2^{14} = 91\,168$

146 The aforementioned setups are described visually in Figure 3.A.

147 **Evaluation under different data regimes** Note that alignment and decoding models need not be  
 148 fitted using the same amount of data. In particular, we are interested in evaluating out-of-subject  
 149 performance in setups where a lot of data is available for a *reference* subject, and little data is  
 150 available for a *left-out* subject: this would typically be the case in clinical setups where little data is  
 151 available in patients. In this case, we evaluate whether it is possible to use this small amount of data  
 152 to align the left-out subject onto the reference subject, and have the left-out subject benefit from a  
 153 decoder previously trained on a lot of data.

## 154 2.4 DATASET

155 We analyze the dataset from Wen et al. (2017). This dataset comprises 3 human subjects who each  
 156 watched 688 minutes of video in an MRI scanner. The videos consists of 18 train segments of 8  
 157 minutes each and 5 test segments of 8 minutes each. Each train segment was presented twice. Each  
 158 test segment was presented 10 times. Each segment consists of a sequence of roughly 10-second  
 159 video clips.

160 The fMRI data was acquired at 3 Tesla (3T), 3.5mm isotropic spatial resolution and 2-second tem-  
 161 poral resolution. It was minimally pre-processed with the same pre-processing pipeline than that  
 162 of the Human Connectome Project (Glasser et al., 2013). In particular, data from each subject are  
 163 projected onto a common volumetric anatomical template.

164 Comparably to prior work on this dataset (Wen et al., 2018; Kupersmidt et al., 2022; Wang et al.,  
 165 2022), we use runs related to the first 18 video segments - 288 minutes - as training data, and runs  
 166 related to the last 5 video segments as test data.

## 167 2.5 PREPROCESSING

168 We implement minimal additional preprocessing steps for each subject separately. For this, we (1)  
 169 project all volumetric data onto the FreeSurfer average surface template *fsaverage5* (Fischl, 2012),  
 170 then (2) regress out cosine drifts in each vertex and each run and finally (3) center and scale each  
 171 vertex time-course in each run. Figure S1 gives a visual explanation as to why the last two steps are  
 172 needed. The first two steps are implemented with nilearn (Abraham et al., 2014)<sup>4</sup> and the last one  
 173 with scikit-learn (Pedregosa et al., 2011).

174 Additionally, for a given subject, we try out two different setups: a first one where runs showing the  
 175 same video are averaged, and a second one where they are stacked.

## 176 2.6 HYPER-PARAMETERS SELECTION

177 To train decoders, we use the same regularization coefficient  $\alpha_{\text{ridge}}$  across latent types and choose it  
 178 by running a cross-validated grid search on folds of the training data. We find that results are robust  
 179 to using different values and stick to  $\alpha_{\text{ridge}} = 50\,000$ . Similarly, values for lag, window size and  
 180 aggregation function are determined through a cross-validated grid search.

181 Finally, for functional alignment, we stick to default parameters shipped with version 0.1.0 of  
 182 FUGW. Namely,  $\alpha$ , which balances between Wasserstein and Gromov-Wasserstein losses – i.e. how  
 183 important functional data is compared to anatomical data – is set to 0.5. Empirically, we see that this  
 184 value yields values for the Wasserstein loss which are bigger than that of the Gromov-Wasserstein  
 185 loss, meaning that functional data drives these alignments.  $\varepsilon$ , which controls for entropic regulariza-  
 186 tion – i.e. how blurry computed alignments will be – is set to  $10^{-4}$ . Empirically, this value yields  
 187 very anatomically sharp alignments.  $\rho$ , which sets the importance of marginal constraints – i.e. to  
 188 what extent more or less mass can be transported to / from each voxel – is set to 1. Empirically, this  
 189 value leads to all voxels being transported / matched.

<sup>4</sup><https://nilearn.github.io>

### 190 3 RESULTS

#### 191 3.1 WITHIN-SUBJECT PREDICTION OF VISUAL REPRESENTATIONS FROM BOLD SIGNAL 192 AND RETRIEVAL OF VISUAL INPUTS

193 We report the retrieval predictions of video decoding results in Table 1. For all three subjects of the  
194 Wen et al. (2017) dataset, and for all four types of latent representations considered, a Ridge regres-  
195 sion fitted within-subject achieves significantly above-chance performance. Besides, performance  
196 varies across subjects, although well-performing subjects reach good performance on all types of  
197 latents.

198 Results reported in Table 1 were obtained for a lag of 2 brain volumes (i.e. 4 seconds since TR =  
199 2 seconds) and a window size of 2 brain volumes which were averaged together (see definitions in  
200 section 2.5). These parameters were chosen after running a grid search for lag values ranging from  
201 1 to 5, a window size ranging from 1 to 3, and 2 possible aggregation functions for brain volumes  
202 belonging to the same window (namely averaging and stacking). Figure S4 shows results using  
203 the averaging aggregation function for different values of lag and window size, averaged across  
204 subjects. Eventually, these results were obtained by stacking all runs of the training dataset, as  
205 opposed to averaging repetitions of the same video clip. The two approaches yielded very similar  
206 metrics. We expand on this matter in section 3.3.

207 Finally, Figure 2 shows retrieved images for Subject 2. Qualitatively, we observe that retrieved  
208 images often fit the theme of images shown to subjects (with categories like indoor sports, human  
209 faces, animals, etc.), but also regularly exhibit failure cases. It is also possible to use predicted  
210 latents to reconstruct seen video clips at a low frame-per-second rate (see Figure S3), which we do  
211 not attempt in this study.

Table 1: **Within-subject metrics for all subjects and all latent types on the test set** Reported metrics are relative median rank  $\downarrow$  (MR) of retrieval on a set of 500 samples, top-5 accuracy %  $\uparrow$  (Acc) of retrieval on a set of 500 samples. These results were averaged across 50 retrieval sets, hence results are reported with a standard error of the mean (SEM) smaller than 0.01. The *Dummy* model systematically predicts the mean latent representation of the training set.

|       | CLIP $257 \times 768$ |      | VD-VAE |     | CLIP CLS |      | AutoKL |     |
|-------|-----------------------|------|--------|-----|----------|------|--------|-----|
|       | MR                    | Acc  | MR     | Acc | MR       | Acc  | MR     | Acc |
| Dummy | 50.0                  | 1.0  | 50.0   | 1.0 | 50.0     | 1.0  | 50.0   | 1.0 |
| S1    | 9.4                   | 13.8 | 29.9   | 3.0 | 15.1     | 8.4  | 24.9   | 3.9 |
| S2    | 6.8                   | 16.4 | 30.2   | 3.5 | 10.6     | 10.5 | 21.8   | 3.8 |
| S3    | 7.8                   | 13.6 | 28.5   | 3.1 | 11.0     | 9.9  | 26.0   | 3.3 |

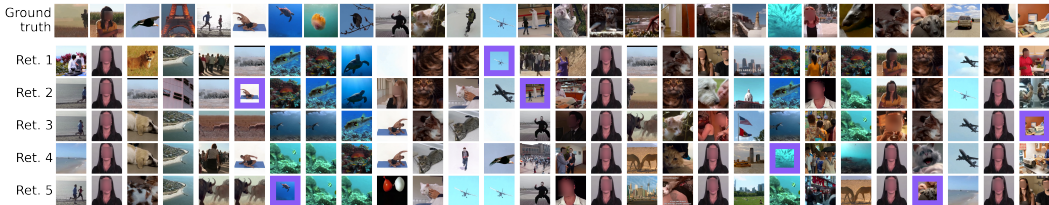


Figure 2: **Image retrievals using predicted latent representations of CLIP  $257 \times 768$  latents** We use a model fitted on Subject 2 (S2) and predict the latent representation of unseen videos (test set). Ground truth images featured within the first 5 retrieved images are indicated with a bold purple border. In a given column, images which appear similar across rows are actually different frames of the same video clip. Images featuring human faces were blurred.

#### 212 3.2 OUT-OF-SUBJECT DECODING AND MULTI-SUBJECT TRAINING

213 As illustrated in Figure 3, models trained on one subject do not generalise well to other subjects.  
214 However, we demonstrate that functional alignment can successfully be used as a transfer learning  
215 strategy to generalize a pre-trained model to left-out subjects. In particular, we show that left-out

216 subjects need not have the same amount of available data than training subjects to benefit from their  
 217 model: with just 30 minutes of data, left-out subjects can reach performance which would have  
 218 needed roughly 100 minutes of data in a within-subject setting. Besides, compared to the out-of-  
 219 subject baseline, we obtain 25 to 75 percents improvement in relative median rank across latent  
 220 types. Note that, in this study, we chose the best performing subject (S2) as the reference subject.  
 221 Finally, we show that a single model trained on all functionally aligned subjects can reach slightly  
 222 better results than models trained on all un-aligned subjects. In every subject, this multi-subject  
 223 aligned model performs comparably to their associated within-subject model. Supplementary Fig-  
 224 ures S5 and S6 show that these results hold for all types for latents.  
 225 Other interesting setups are reported in Figures S8, S9, S10, S11. In particular, they show that a  
 226 multi-subject aligned model (e.g. trained on S1 and S2) has better performance on aligned left-out  
 227 subjects (e.g. S3) than a single-subject model (e.g. trained on S2 only).

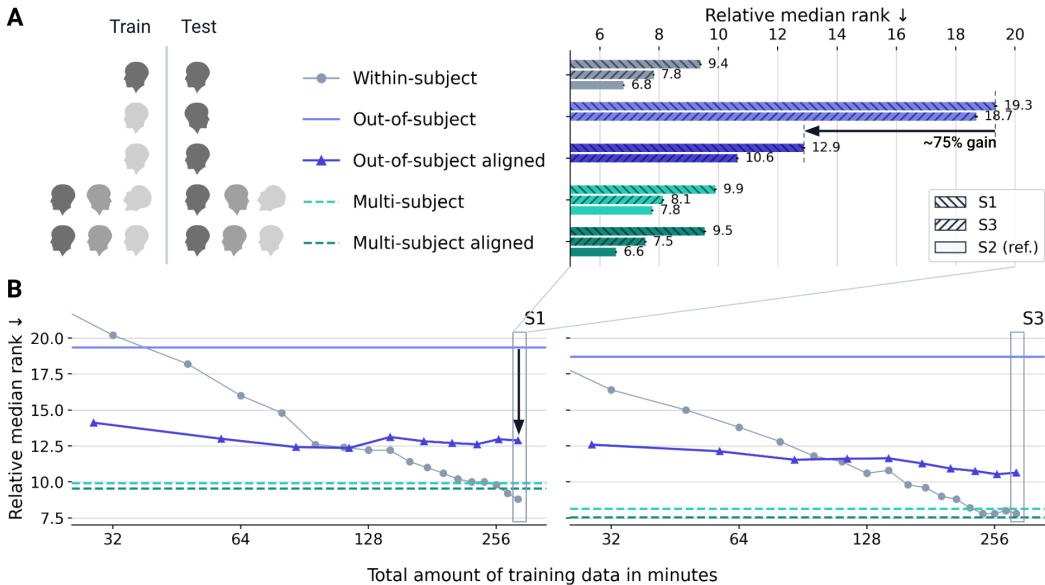


Figure 3: **Effects of functional alignment on multi-subject and out-of-subject setups**  
 We report relative median rank  $\downarrow$  in all setups described in section 2.3 for CLIP  $257 \times 768$ . In all *aligned* cases, S1 and S3 were aligned onto S2. In all *out-of-subject* cases, we test S1 and S3 onto a decoder trained on S2. In all *multi-subject* cases, the decoder was trained on all data from all 3 subjects. **A.** In this panel, all models (alignment and decoding) were trained on all available training data. Results for other latent types are available in Figure S5. **B.** In left-out S1 and S3, decoding performance is much better when using functional alignment to S2 (solid dark purple) than when using anatomical alignment only (solid pale purple). Performance increases slightly as the amount of data used to align subjects grows, but does not always reach levels which can be achieved with a single-subject model fitted in left-out subjects (solid pale gray dots) when a lot of training data is available. Training a model on multiple subjects yields good performance in all 3 subjects (dashed pale teal) which can be further improved by using functional alignment (dashed dark teal). Results for other latent types are available in Figure S6.

228 To better understand how brain features are transformed by functional alignment, we show in Figure  
 229 4 how vertices from S1 are permuted to fit those of S2. Note that both subjects' data lie on fsaver-  
 230 age5. To this end, we colorize vertices in S1 using the MMP 1.0 atlas (Glasser et al., 2016) and  
 231 use  $\phi_{S1 \rightarrow S2}$  to transport each of the three RGB channels of this colorization. We see that, even in  
 232 low data regimes, FUGW scrambles most of the brain but can leverage signal to recover the cortical  
 233 organization of the occipital lobe. Higher regimes yield anatomically-consistent matches in a much  
 234 higher number of cortical areas such as the temporal and parietal lobes, and more surprisingly in the  
 235 primary motor cortex as well, while the prefrontal cortex and temporo-parietal junction (TPJ) still  
 236 seem challenging to map.

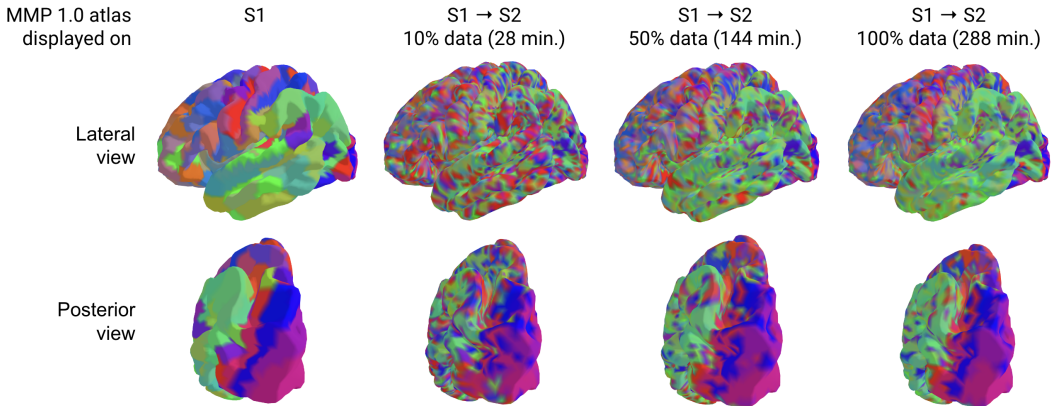


Figure 4: **Visualizing functional alignments in the left hemisphere** Vertices of the source subject (left) are permuted by FUGW. The result of this permutation is visualized on the target subject (columns 2, 3, and 4). Fitting FUGW with different increasing amounts of data gradually unfolds the cortical organisation of multiple areas, even non-visual ones. Note that all 3 models have been fitted using the same number of iterations.

237 3.3 INFLUENCE OF TRAINING SET SIZE AND TEST SET REPETITIONS

238 Recent publications in brain-decoding using non-invasive brain imagery show impressive results.  
 239 However, we stress that these results are obtained in setups which are very advantageous when it  
 240 comes to both dataset size and signal-to-noise ratio. To better assess the importance of these two  
 241 factors, we report in Figure 5 performance metrics for subject models trained and tested with various  
 242 amounts of data and various amounts of noise.

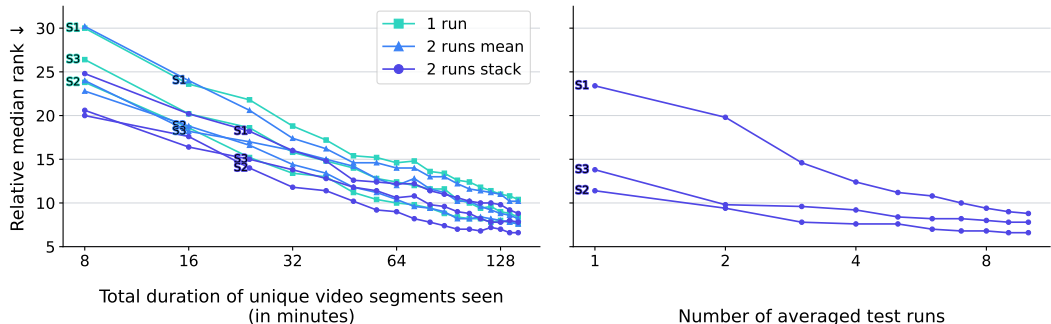


Figure 5: **Effect of training set size and test set noise on retrieval metrics** Relative median rank  $\downarrow$  on a fixed test set gets better as more training data is used to fit the model (left). Interestingly, averaging brain volumes of 2 similar runs does not bring improvements compared to using just 1 run. Instead, stacking runs does yield significant improvements. Note that training sets using 2 runs have twice as much data as those using 1 run. Finally, these metrics are highly affected by the noise level of the test set (right): averaging more runs in the test set yields better metrics despite using the same decoder.

243 Firstly, using a fixed test set in which brain features were averaged across runs, we find that expo-  
 244 nentially more training data is required per subject to achieve better performance. This finding is  
 245 similar to that of systematic scaling studies on similar topics (Tang et al., 2023). More interestingly,  
 246 in this given signal-to-noise ratio setup, it seems that more diverse training data should bring com-  
 247 parable or better performance than repeating already seen content, while potentially covering more  
 248 semantic domains.  
 249 Secondly, reported performance metrics only hold in favorable signal-to-noise setups. Indeed, the  
 250 test set associated with the Wen 2017 dataset comes with 10 runs for each video segment, which,  
 251 when averaged together, greatly reduce the noise level. However, as reported here, when tested  
 252 in real-life signal-to-noise conditions (i.e. only one run per video clip), our models’ performance  
 253 degrades: it is approximately twice as bad for each subject when using CLIP latents.



## 254 4 DISCUSSION

255 **Impact** The present work confirms the feasibility of using BOLD fMRI signal acquired in a nat-  
 256 uralistic setup to decode high level visual features (Nishimoto et al., 2011). It further demonstrates  
 257 that it is possible to leverage fMRI signal from naturalistic movie watching to derive meaningful  
 258 functional alignments between subjects, which in turn can be used to transfer decoding models to  
 259 novel subjects.

260 In particular, our study shows that decoding brain data from a left-out subject can be substantially  
 261 improved by aligning this left-out subject to a large reference dataset on which a decoder was trained.  
 262 Our method thus paves the way towards using models used on large amounts of individual data to  
 263 decode signal acquired in smaller neuro-imaging studies, which typically record one hour of fMRI  
 264 for each subject (Madan, 2022).

265 Besides, this study reports decoding accuracy in setups where subjects are showed test stimuli for  
 266 the first time only, hence yielding insights on how these models would perform in real-time decod-  
 267 ing. While performance improves with the number of repetitions at test time, reasonable decoding  
 268 performance of semantics can be achieved in two out of three subjects with just one repetition.

269 Lastly, by systematically quantifying decoding accuracy as a function of the amount of training data,  
 270 the present work brings insightful recommendations as to what stimuli should be played in future  
 271 fMRI datasets collecting large amounts of data in a limited number of subjects. In the current setup  
 272 (naturalistic movie watching at 3T), more diverse semantic content is more valuable than repeated  
 273 content for fitting decoding models.

274 **Limitations** This work is a first step towards training accurate semantic decoders which generalize  
 275 across individuals, but subsequent work remains necessary to ensure the generality of our findings.

276 Firstly, although reported gains in out-of-subject setups are significant, the small number of partic-  
 277 ipants present in the dataset under study calls for replications on other – and potentially larger –  
 278 cohorts. However, to our knowledge, no other dataset presented similar features to that of Wen et al.  
 279 (2017) – i.e. high quantity of data per subject and large variety of video stimuli. The recent Courtois  
 280 Neuromod dataset<sup>5</sup> might be useful in this regard.

281 Secondly, our approach currently requires left-out subjects to watch the same videos as reference  
 282 subjects. It is yet unclear whether functional alignment could bring improvements without this  
 283 constraint. However, multi-subject decoding can probably help partially address this issue: since it  
 284 is possible to train a decoder on multiple subjects and because not all of them have to watch the same  
 285 movies, it is possible that a lot of different movies could be used as “anchor” for left-out individuals.

286 Thirdly, unlike other approaches (Défossez et al., 2022), our approach relies on pre-trained encoders,  
 287 and cannot align all subjects at once. Consequently, overall performance highly depends on the  
 288 quality of other models and of data acquired in reference individuals.

289 Finally, while restricting this study to linear models makes sense to establish baselines and ensure  
 290 reproducibility, non-linear models have proved to be very efficient. A natural improvement on this  
 291 work could include these architectures.

292 **Ethical implications** Out-of-subject generalization is an important test for decoding models, but it  
 293 raises legitimate concerns. In this regard, this study highlights that signal-to-noise ratio still currently  
 294 makes it challenging to very accurately decode semantics in a real-time setup, and that a non-trivial  
 295 amount of data is needed per individual for these models to work. Moreover, we stress that, while  
 296 decoding perceived stimuli is making great progress, imagined stimuli are still very challenging  
 297 (Horikawa & Kamitani, 2017). Nonetheless, it is important for advances in this domain to be pub-  
 298 licly documented. We thus advocate that open and peer-reviewed research is the best way forward  
 299 to safely explore the implications of inter-subject modeling, and more generally brain decoding.

300 **Conclusion** Overall, these results provide a significant step towards real-time, subject-agnostic  
 301 visual decoding of semantics using fMRI.

<sup>5</sup><https://www.cneuromod.ca>

## 302 REFERENCES

- 303 Alexandre Abraham, Fabian Pedregosa, Michael Eickenberg, Philippe Gervais, Andreas Mueller,  
304 Jean Kossaifi, Alexandre Gramfort, Bertrand Thirion, and Gael Varoquaux. Machine learning for  
305 neuroimaging with scikit-learn. *Frontiers in Neuroinformatics*, 8, 2014. ISSN 1662-5196. URL  
306 <https://www.frontiersin.org/article/10.3389/fninf.2014.00014>.
- 307 Emily J. Allen, Ghislain St-Yves, Yihan Wu, Jesse L. Breedlove, Jacob S. Prince, Logan T. Dow-  
308 dle, Matthias Nau, Brad Caron, Franco Pestilli, Ian Charest, J. Benjamin Hutchinson, Thomas  
309 Naselaris, and Kendrick Kay. A massive 7T fMRI dataset to bridge cognitive neuroscience  
310 and artificial intelligence. *Nature Neuroscience*, 25(1):116–126, January 2022. ISSN 1546-  
311 1726. doi: 10.1038/s41593-021-00962-x. URL <https://www.nature.com/articles/s41593-021-00962-x>. Number: 1 Publisher: Nature Publishing Group.
- 313 T. Bazeille, H. Richard, H. Janati, and B. Thirion. Local Optimal Transport for Functional Brain  
314 Template Estimation. In Albert C. S. Chung, James C. Gee, Paul A. Yushkevich, and Siqi Bao  
315 (eds.), *Information Processing in Medical Imaging*, Lecture Notes in Computer Science, pp. 237–  
316 248, Cham, 2019. Springer International Publishing. ISBN 978-3-030-20351-1. doi: 10.1007/  
317 978-3-030-20351-1\_18.
- 318 Thomas Bazeille, Elizabeth DuPre, Hugo Richard, Jean-Baptiste Poline, and Bertrand Thirion.  
319 An empirical evaluation of functional alignment using inter-subject decoding. *NeuroImage*,  
320 245:118683, December 2021. ISSN 1053-8119. doi: 10.1016/j.neuroimage.  
321 2021.118683. URL <https://www.sciencedirect.com/science/article/pii/S1053811921009563>.
- 323 Po-Hsuan (Cameron) Chen, Janice Chen, Yaara Yeshurun, Uri Hasson, James Haxby,  
324 and Peter J Ramadge. A Reduced-Dimension fMRI Shared Response Model. In  
325 *Advances in Neural Information Processing Systems*, volume 28. Curran Associates,  
326 Inc., 2015. URL [https://papers.nips.cc/paper\\_files/paper/2015/hash/b3967a0e938dc2a6340e258630febd5a-Abstract.html](https://papers.nips.cc/paper_files/paper/2015/hash/b3967a0e938dc2a6340e258630febd5a-Abstract.html).
- 328 Zijiao Chen, Jiaxin Qing, Tiange Xiang, Wan Lin Yue, and Juan Helen Zhou. Seeing Beyond the  
329 Brain: Conditional Diffusion Model with Sparse Masked Modeling for Vision Decoding, March  
330 2023a. URL <http://arxiv.org/abs/2211.06956>. arXiv:2211.06956 [cs].
- 331 Zijiao Chen, Jiaxin Qing, and Juan Helen Zhou. Cinematic Mindscapes: High-quality Video Recon-  
332 struction from Brain Activity, May 2023b. URL <http://arxiv.org/abs/2305.11675>.  
333 arXiv:2305.11675 [cs].
- 334 Marco Cuturi. Sinkhorn Distances: Lightspeed Computation of Optimal Transportation Distances.  
335 *arXiv*, June 2013. doi: 10.48550/arXiv.1306.0895.
- 336 Alexandre Défossez, Charlotte Caucheteux, Jérémy Rapin, Ori Kabeli, and Jean-Rémi King. De-  
337 coding speech from non-invasive brain recordings, August 2022. URL <http://arxiv.org/abs/2208.12266>. arXiv:2208.12266 [cs, eess, q-bio].
- 339 Ma Feilong, Samuel A. Nastase, Guo Jiahui, Yaroslav O. Halchenko, M. Ida Gobbi, and James V.  
340 Haxby. The Individualized Neural Tuning Model: Precise and generalizable cartography of func-  
341 tional architecture in individual brains, May 2022. URL <https://www.biorxiv.org/content/10.1101/2022.05.15.492022v1>. Pages: 2022.05.15.492022 Section: New  
342 Results.  
343
- 344 Matteo Ferrante, Furkan Ozcelik, Tommaso Boccato, Rufin VanRullen, and Nicola Toschi. Brain  
345 Captioning: Decoding human brain activity into images and text, May 2023. URL <http://arxiv.org/abs/2305.11560>. arXiv:2305.11560 [cs].
- 347 Bruce Fischl. FreeSurfer. *NeuroImage*, 62(2):774–781, August 2012. ISSN 1053-8119. doi: 10.  
348 1016/j.neuroimage.2012.01.021. URL <https://www.sciencedirect.com/science/article/pii/S1053811912000389>.

- 350 Matthew F. Glasser, Stamatiios N. Sotiropoulos, J. Anthony Wilson, Timothy S. Coalson, Bruce  
351 Fischl, Jesper L. Andersson, Junqian Xu, Saad Jbabdi, Matthew Webster, Jonathan R. Polimeni,  
352 David C. Van Essen, Mark Jenkinson, and WU-Minn HCP Consortium. The minimal preprocess-  
353 ing pipelines for the Human Connectome Project. *NeuroImage*, 80:105–124, October 2013. ISSN  
354 1095-9572. doi: 10.1016/j.neuroimage.2013.04.127.
- 355 Matthew F. Glasser, Timothy S. Coalson, Emma C. Robinson, Carl D. Hacker, John Harwell, Essa  
356 Yacoub, Kamil Ugurbil, Jesper Andersson, Christian F. Beckmann, Mark Jenkinson, Stephen M.  
357 Smith, and David C. Van Essen. A multi-modal parcellation of human cerebral cortex. *Nature*,  
358 536(7615):171–178, August 2016. ISSN 1476-4687. doi: 10.1038/nature18933. URL  
359 <https://www.nature.com/articles/nature18933>. Number: 7615 Publisher: Na-  
360 ture Publishing Group.
- 361 G. H. Glover. Deconvolution of impulse response in event-related BOLD fMRI. *NeuroImage*, 9(4):  
362 416–429, April 1999. ISSN 1053-8119. doi: 10.1006/nimg.1998.0419.
- 363 Zijin Gu, Keith Jamison, Amy Kuceyeski, and Mert Sabuncu. Decoding natural image stimuli from  
364 fMRI data with a surface-based convolutional network, March 2023. URL <http://arxiv.org/abs/2212.02409>. arXiv:2212.02409 [cs, q-bio].
- 365  
366 Stephenie A Harrison and Frank Tong. Decoding reveals the contents of visual working memory in  
367 early visual areas. *Nature*, 458(7238):632–635, 2009.
- 368 James V. Haxby, J. Swaroop Guntupalli, Andrew C. Connolly, Yaroslav O. Halchenko, Bryan R.  
369 Conroy, M. Ida Gobbini, Michael Hanke, and Peter J. Ramadge. A common, high-dimensional  
370 model of the representational space in human ventral temporal cortex. *Neuron*, 72(2):404–416,  
371 October 2011. ISSN 1097-4199. doi: 10.1016/j.neuron.2011.08.026.
- 372 James V Haxby, J Swaroop Guntupalli, Samuel A Nastase, and Ma Feilong. Hyperalignment:  
373 Modeling shared information encoded in idiosyncratic cortical topographies. *eLife*, 9:e56601,  
374 June 2020. ISSN 2050-084X. doi: 10.7554/eLife.56601. URL [https://doi.org/10.](https://doi.org/10.7554/eLife.56601)  
375 [7554/eLife.56601](https://doi.org/10.7554/eLife.56601). Publisher: eLife Sciences Publications, Ltd.
- 376 John-Dylan Haynes and Geraint Rees. Decoding mental states from brain activity in humans. *Nature*  
377 *reviews neuroscience*, 7(7):523–534, 2006.
- 378 Jun Kai Ho, Tomoyasu Horikawa, Kei Majima, Fan Cheng, and Yukiyasu Kamitani. Inter-individual  
379 deep image reconstruction via hierarchical neural code conversion. *NeuroImage*, 271:120007,  
380 May 2023. ISSN 1053-8119. doi: 10.1016/j.neuroimage.2023.120007. URL [https://www.](https://www.sciencedirect.com/science/article/pii/S1053811923001532)  
381 [sciencedirect.com/science/article/pii/S1053811923001532](https://www.sciencedirect.com/science/article/pii/S1053811923001532).
- 382 Tomoyasu Horikawa and Yukiyasu Kamitani. Generic decoding of seen and imagined objects using  
383 hierarchical visual features. *Nat. Commun.*, 8(15037):1–15, May 2017. ISSN 2041-1723. doi:  
384 10.1038/ncomms15037.
- 385 Ganit Kupersmidt, Roman Belyi, Guy Gaziv, and Michal Irani. A Penny for Your (visual)  
386 Thoughts: Self-Supervised Reconstruction of Natural Movies from Brain Activity, June 2022.  
387 URL <http://arxiv.org/abs/2206.03544>. arXiv:2206.03544 [cs].
- 388 Benjamin Lahner, Kshitij Dwivedi, Polina Iamshchinina, Monika Graumann, Alex Lascelles,  
389 Gemma Roig, Alessandro Thomas Gifford, Bowen Pan, SouYoung Jin, N. Apurva Ratan Murty,  
390 Kendrick Kay, Aude Oliva, and Radoslaw Cichy. BOLD Moments: modeling short visual events  
391 through a video fMRI dataset and metadata, March 2023. URL [https://www.biorxiv.](https://www.biorxiv.org/content/10.1101/2023.03.12.530887v1)  
392 [org/content/10.1101/2023.03.12.530887v1](https://www.biorxiv.org/content/10.1101/2023.03.12.530887v1). Pages: 2023.03.12.530887 Section:  
393 New Results.
- 394 Amanda LeBel, Lauren Wagner, Shailee Jain, Aneesh Adhikari-Desai, Bhavin Gupta, Allyson  
395 Morgenthal, Jerry Tang, Lixiang Xu, and Alexander G. Huth. A natural language fMRI  
396 dataset for voxelwise encoding models. *Scientific Data*, 10(1):555, August 2023. ISSN 2052-  
397 4463. doi: 10.1038/s41597-023-02437-z. URL [https://www.nature.com/articles/](https://www.nature.com/articles/s41597-023-02437-z)  
398 [s41597-023-02437-z](https://www.nature.com/articles/s41597-023-02437-z). Number: 1 Publisher: Nature Publishing Group.

- 399 Christopher R. Madan. Scan Once, Analyse Many: Using Large Open-Access Neuroimaging  
400 Datasets to Understand the Brain. *Neuroinformatics*, 20(1):109–137, January 2022. ISSN 1559-  
401 0089. doi: 10.1007/s12021-021-09519-6.
- 402 Weijian Mai and Zhijun Zhang. UniBrain: Unify Image Reconstruction and Captioning All in One  
403 Diffusion Model from Human Brain Activity, August 2023. URL [http://arxiv.org/abs/  
404 2308.07428](http://arxiv.org/abs/2308.07428). arXiv:2308.07428 [cs].
- 405 Tom M Mitchell, Rebecca Hutchinson, Radu S Niculescu, Francisco Pereira, Xuerui Wang, Marcel  
406 Just, and Sharlene Newman. Learning to decode cognitive states from brain images. *Machine  
407 learning*, 57:145–175, 2004.
- 408 Shinji Nishimoto, An T. Vu, Thomas Naselaris, Yuval Benjamini, Bin Yu, and Jack L. Gallant.  
409 Reconstructing visual experiences from brain activity evoked by natural movies. *Current Biology*,  
410 21(19):1641–1646, October 2011. ISSN 0960-9822. doi: 10.1016/j.cub.2011.08.031. URL  
411 <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3326357/>.
- 412 Furkan Ozcelik and Rufin VanRullen. Natural scene reconstruction from fMRI signals using  
413 generative latent diffusion, June 2023. URL <http://arxiv.org/abs/2303.05334>.  
414 arXiv:2303.05334 [cs, q-bio].
- 415 Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier  
416 Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, Jake Vanderplas,  
417 Alexandre Passos, David Cournapeau, Matthieu Brucher, Matthieu Perrot, and Édouard Duch-  
418 esnay. Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research*,  
419 12(85):2825–2830, 2011. ISSN 1533-7928. URL [http://jmlr.org/papers/v12/  
420 pedregosalla.html](http://jmlr.org/papers/v12/pedregosalla.html).
- 421 Erin M. Phillips, Kirsten D. Gillette, Daniel D. Dilks, and Gregory S. Berns. Through  
422 a Dog’s Eyes: fMRI Decoding of Naturalistic Videos from the Dog Cortex. *JoVE*  
423 (*Journal of Visualized Experiments*), (187):e64442, September 2022. ISSN 1940-  
424 087X. doi: 10.3791/64442. URL [https://www.jove.com/fr/v/64442/  
425 through-dog-s-eyes-fmri-decoding-naturalistic-videos-from-dog](https://www.jove.com/fr/v/64442/through-dog-s-eyes-fmri-decoding-naturalistic-videos-from-dog).
- 426 Ana Luísa Pinho, Alexis Amadon, Torsten Ruest, Murielle Fabre, Elvis Dohmatob, Isabelle  
427 Denghien, Chantal Ginisty, Séverine Becuwe-Desmidt, Séverine Roger, Laurence Laurier,  
428 Véronique Joly-Testault, Gaëlle Médiouni-Cloarec, Christine Doublé, Bernadette Martins,  
429 Philippe Pinel, Evelyn Eger, Gaël Varoquaux, Christophe Pallier, Stanislas Dehaene, Lucie Hertz-  
430 Pannier, and Bertrand Thirion. Individual Brain Charting, a high-resolution fMRI dataset for cog-  
431 nitive mapping. *Scientific Data*, 5(1):180105, June 2018. ISSN 2052-4463. doi: 10.1038/sdata.  
432 2018.105. URL <https://www.nature.com/articles/sdata2018105>. Number: 1  
433 Publisher: Nature Publishing Group.
- 434 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agar-  
435 wal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya  
436 Sutskever. Learning Transferable Visual Models From Natural Language Supervision, February  
437 2021. URL <http://arxiv.org/abs/2103.00020>. arXiv:2103.00020 [cs].
- 438 Hugo Richard, Luigi Gresele, Aapo Hyvärinen, Bertrand Thirion, Alexandre Gramfort, and Pierre  
439 Ablin. Modeling Shared Responses in Neuroimaging Studies through MultiView ICA, December  
440 2020. URL <http://arxiv.org/abs/2006.06635>. arXiv:2006.06635 [cs, stat].
- 441 Emma C. Robinson, Saad Jbabdi, Matthew F. Glasser, Jesper Andersson, Gregory C. Burgess,  
442 Michael P. Harms, Stephen M. Smith, David C. Van Essen, and Mark Jenkinson. MSM: a new  
443 flexible framework for Multimodal Surface Matching. *NeuroImage*, 100:414–426, October 2014.  
444 ISSN 1095-9572. doi: 10.1016/j.neuroimage.2014.05.069.
- 445 Paul S. Scotti, Atmadeep Banerjee, Jimmie Goode, Stepan Shabalín, Alex Nguyen, Ethan Cohen,  
446 Aidan J. Dempster, Nathalie Verlinde, Elad Yundler, David Weisberg, Kenneth A. Norman, and  
447 Tanishq Mathew Abraham. Reconstructing the Mind’s Eye: fMRI-to-Image with Contrastive  
448 Learning and Diffusion Priors, May 2023. URL <http://arxiv.org/abs/2305.18274>.  
449 arXiv:2305.18274 [cs, q-bio].

- 450 Thibault Séjourné, François-Xavier Vialard, and Gabriel Peyré. The Unbalanced Gromov Wasser-  
451 stein Distance: Conic Formulation and Relaxation. *arXiv:2009.04266 [math, stat]*, June 2021.  
452 URL <http://arxiv.org/abs/2009.04266>. arXiv: 2009.04266.
- 453 Yu Takagi and Shinji Nishimoto. High-resolution image reconstruction with latent diffusion models  
454 from human brain activity, March 2023. URL [https://www.biorxiv.org/content/  
455 10.1101/2022.11.18.517004v3](https://www.biorxiv.org/content/10.1101/2022.11.18.517004v3). Pages: 2022.11.18.517004 Section: New Results.
- 456 Jerry Tang, Amanda LeBel, Shailee Jain, and Alexander G. Huth. Semantic reconstruction of con-  
457 tinuous language from non-invasive brain recordings. *Nature Neuroscience*, 26(5):858–866, May  
458 2023. ISSN 1546-1726. doi: 10.1038/s41593-023-01304-9. URL [https://www.nature.  
459 com/articles/s41593-023-01304-9](https://www.nature.com/articles/s41593-023-01304-9). Number: 5 Publisher: Nature Publishing Group.
- 460 Armin Thomas, Christopher Ré, and Russell Poldrack. Self-Supervised Learn-  
461 ing of Brain Dynamics from Broad Neuroimaging Data. *Advances in Neu-  
462 ral Information Processing Systems*, 35:21255–21269, December 2022. URL  
463 [https://proceedings.neurips.cc/paper\\_files/paper/2022/hash/  
464 8600a9df1a087a9a66900cc8c948c3f0-Abstract-Conference.html](https://proceedings.neurips.cc/paper_files/paper/2022/hash/8600a9df1a087a9a66900cc8c948c3f0-Abstract-Conference.html).
- 465 Alexis Thuau, Quang Huy Tran, Tatiana Zemsanova, Nicolas Courty, Rémi Flamary, Stanislas  
466 Dehaene, and Bertrand Thirion. Aligning individual brains with fused unbalanced Gromov  
467 Wasserstein. *Advances in Neural Information Processing Systems*, 35:21792–21804, Decem-  
468 ber 2022. URL [https://proceedings.neurips.cc/paper\\_files/paper/2022/  
469 hash/8906cac4ca58dcaf17e97a0486ad57ca-Abstract-Conference.html](https://proceedings.neurips.cc/paper_files/paper/2022/hash/8906cac4ca58dcaf17e97a0486ad57ca-Abstract-Conference.html).
- 470 Chong Wang, Hongmei Yan, Wei Huang, Jiayi Li, Yuting Wang, Yun-Shuang Fan, Wei Sheng,  
471 Tao Liu, Rong Li, and Huaifu Chen. Reconstructing rapid natural vision with fMRI-conditional  
472 video generative adversarial network. *Cerebral Cortex*, 32(20):4502–4511, October 2022. ISSN  
473 1047-3211. doi: 10.1093/cercor/bhab498. URL [https://doi.org/10.1093/cercor/  
474 bhab498](https://doi.org/10.1093/cercor/bhab498).
- 475 Haiguang Wen, Junxing Shi, Yizhen Zhang, Kun-Han Lu, Jiayue Cao, and Zhongming Liu. Data for  
476 Neural Encoding and Decoding with Deep Learning for Dynamic Natural Vision Tests, September  
477 2017. URL <https://purr.purdue.edu/publications/2809/1>.
- 478 Haiguang Wen, Junxing Shi, Yizhen Zhang, Kun-Han Lu, Jiayue Cao, and Zhongming Liu. Neural  
479 Encoding and Decoding with Deep Learning for Dynamic Natural Vision. *Cerebral Cortex (New  
480 York, N.Y.: 1991)*, 28(12):4136–4160, December 2018. ISSN 1460-2199. doi: 10.1093/cercor/  
481 bhx268.