

## Self-supervised Natural Image Reconstruction and Large-scale Semantic Classification from Brain Activity



Guy Gaziv<sup>a,1,\*</sup>, Roman Beliy<sup>a,1</sup>, Niv Granot<sup>a,1</sup>, Assaf Hoogi<sup>a</sup>, Francesca Strappini<sup>b</sup>, Tal Golan<sup>c</sup>, Michal Irani<sup>a</sup>

<sup>a</sup> Dept. of Computer Science and Applied Math, Weizmann Institute of Science, Rehovot, Israel

<sup>b</sup> Dept. of Neurobiology, Weizmann Institute of Science, Rehovot, Israel

<sup>c</sup> Zuckerman Institute, Columbia University, New York, NY USA

### ARTICLE INFO

**Keywords:**

Self-Supervised learning, Decoding, Encoding, fMRI, Image reconstruction, Classification vision

### ABSTRACT

Reconstructing natural images and decoding their semantic category from fMRI brain recordings is challenging. Acquiring sufficient pairs of images and their corresponding fMRI responses, which span the huge space of natural images, is prohibitive. We present a novel *self-supervised* approach that goes well beyond the scarce paired data, for achieving both: (i) state-of-the art fMRI-to-image reconstruction, and (ii) first-ever large-scale semantic classification from fMRI responses. By imposing cycle consistency between a pair of deep neural networks (from image-to-fMRI & from fMRI-to-image), we train our image reconstruction network on a large number of “unpaired” natural images (images without fMRI recordings) from many novel semantic categories. This enables to adapt our reconstruction network to a very rich semantic coverage without requiring any explicit semantic supervision. Specifically, we find that combining our self-supervised training with *high-level perceptual losses*, gives rise to new reconstruction & classification capabilities. In particular, this perceptual training enables to classify well fMRIs of never-before-seen semantic classes, *without requiring any class labels during training*. This gives rise to: (i) Unprecedented image-reconstruction from fMRI of never-before-seen images (evaluated by image metrics and human testing), and (ii) Large-scale semantic classification of categories that were never-before-seen during network training. *Such large-scale (1000-way) semantic classification from fMRI recordings has never been demonstrated before.* Finally, we provide evidence for the biological consistency of our learned model.

### 1. Introduction

Natural images span a vastly rich visual and semantic space that humans are experts at processing and recognizing. An intriguing inverse problem is to decode images seen by a person and their semantic categories directly from brain activity (Fig 1a) (Miyawaki et al., 2008; Thirion et al., 2006). Modeling the stimulus as a function of cortical activation (decoding) complements the conventional neuroscientific paradigm of modeling cortical activation as a function of the stimulus (encoding) (Thirion et al., 2006). It provides us with an additional benchmark of our understanding of the cortical visual code and is thus a cornerstone in cognitive neuroscience (Cox and Savoy, 2003; Haxby et al., 2001; Haynes and Rees, 2006; Hebart and Baker, 2018; Kamitani and Tong, 2005; 2006; Naselaris et al., 2011; O’Craven and Kanwisher, 2000). Specifically, reconstruction and classification of seen images provide a window into top-down perceptual processes whose probing was previously limited to subjective report. These include visual re-

call (Naselaris et al., 2015), visual imagery in wakefulness (Cichy et al., 2012; Naselaris et al., 2015; Pearson et al., 2015; Shen et al., 2019a; 2019b) and sleep (Horikawa and Kamitani, 2017b; Horikawa et al., 2013). Last, advancements in fMRI-based decoding may be useful for other neuroimaging modalities that provide opportunities for effective brain-machine interfaces (Milekovic et al., 2018; Naci et al., 2012; Tripathy et al., 2021) or for diagnosing disorders of consciousness (Monti et al., 2010; Owen et al., 2006).

In the image reconstruction task, one attempts to decode natural images which were observed by a human subject from the induced brain activity captured by functional magnetic resonance imaging (fMRI). To learn the mapping between fMRI and image representation, typical fMRI datasets provide many pairs of images and their corresponding fMRI responses, referenced to as “paired” data. The goal is to learn an fMRI-to-image decoder which generalizes well to reconstructing images from novel “test-fMRI”, fMRI response induced by novel images from totally different semantic categories than those in the training data (referred to

\* Corresponding author.

E-mail addresses: [guy.gaziv@weizmann.ac.il](mailto:guy.gaziv@weizmann.ac.il) (G. Gaziv), [michal.irani@weizmann.ac.il](mailto:michal.irani@weizmann.ac.il) (M. Irani).

<sup>1</sup> The authors contributed equally to this work.

as “test-images”). Moreover, a complementary challenge to reconstructing the underlying image is also to decode its semantic category. *However, the shortage of “paired” training data limits the generalization power of today’s fMRI decoders.* The number of obtainable image-fMRI pairs is bounded by the limited time a human can spend in an MRI scanner. This also results in a limited number of semantic categories associated with fMRI data. Accordingly, most datasets provide only up to a few thousands of (Image,fMRI) pairs, which are usually obtained from a small number of image categories. Such limited data cannot span the huge space of natural images and their semantic categories, nor the space of their fMRI recordings. Moreover, the poor spatio-temporal resolution of fMRI signals, as well as their low Signal-to-Noise Ratio (SNR), reduce the reliability of the already scarce paired training data (Glasser et al., 2013; Glover, 2011; Lage-Castellanos et al., 2019).

Reconstructing natural images from fMRI was approached by a number of methods, which can broadly be classified into three families: (i) Linear regression between fMRI data and handcrafted image-features (e.g., Gabor wavelets) (Kay et al., 2008; Naselaris et al., 2009; Nishimoto et al., 2011), (ii) Linear regression between fMRI data and deep (CNN-based) image-features (e.g., using pretrained AlexNet) (Guclu et al., 2015; Shen et al., 2019b; Wen et al., 2018d; Zhang et al., 2018a), or latent spaces of pretrained generative models (Han et al., 2019; Mozafari et al., 0000; Qiao et al., 2020; Ren et al., 2021), and (iii) End-to-end Deep Learning (Lin et al., 2019; Seeliger et al., 2018; Shen et al., 2019a; St-Yves and Naselaris, 2019).

To our best knowledge, methods (Shen et al., 2019b) and St-Yves and Naselaris (2019) are the current state-of-the-art in this field. All these methods inherently rely on the available “paired” data to train their decoder (pairs of images and their corresponding fMRI responses). When trained on limited data, purely supervised models are prone to overfitting, which leads to poor generalization to new test-data (fMRI response evoked by new images). To overcome this problem, recent methods (Han et al., 2019; Lin et al., 2019; Mozafari et al., 0000; Ren et al., 2021; Seeliger et al., 2018; Shen et al., 2019b; St-Yves and Naselaris, 2019) introduced natural-image priors in the form of Generative Adversarial Networks (GANs) or Variational Auto-Encoders (VAEs). These methods enabled leap advancement in reconstruction quality from fMRI, and tend to produce natural-looking images. Nevertheless, despite their pleasant natural appearance, their reconstructed images are often *unfaithful* to the actual images underlying the test-fMRI (see Fig 8a).

Prior work on semantic classification of fMRI recordings induced by natural-images, can be characterized as two families: (i) Classifying new images from previously seen categories, and (ii) Classifying new images from novel never-before-seen categories. In the first, the categories to-be-decoded (of the test data) are represented in the training data (Eickenberg et al., 2017; Haxby et al., 2001; Konkle and Caramazza, 2013; Qiao et al., 2019; Wen et al., 2018b). This widely-explored family is limited to decode only the few and typically coarse classes (e.g., faces, scenes, inanimate objects, and birds), which are represented in the limited “paired” data used for decoder training. The second family, which was introduced in Horikawa and Kamitani (2017a), addresses the much more challenging case, where the test-categories are *novel*, namely, not directly represented in the training data. Under this setting, decoding novel, diverse, and fine-grained semantic categories (e.g., ImageNet (Deng et al., 2009)) remains a difficult task because of the narrow semantic coverage spanned by the limited paired training data.

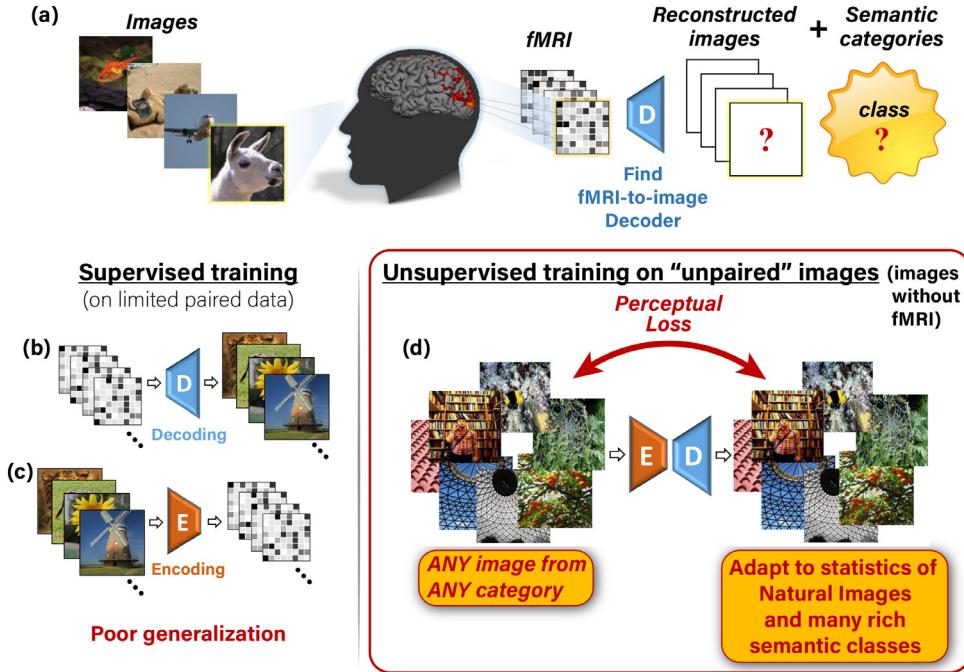
To cope with this data limitation, recent approaches (Eickenberg et al., 2017; Horikawa and Kamitani, 2017a; Qiao et al., 2019; Wen et al., 2018b; 2018d) harnessed a *pretrained* and semantically separable embedding. In this approach, voxel responses are linearly mapped to *higher-level* feature representation of an image classification network. Once mapped, categorization is achieved by either (i) forward-propagating the decoded representation to the classification layer (Wen et al., 2018d), or by (ii) nearest-neighbor classification against a gallery of category representatives, which are the mean feature representations of many natural images from that cate-

gory (Horikawa and Kamitani, 2017a). While these methods benefited from the wide-coverage of their semantic representation (which stems from pretrained image features independently of the fMRI data), their decoding method remained **supervised** in essence. This is because their training of the mapping from fMRI-to-feature representation relies solely on the limited “paired” training data. Consequently, they are prone to poor generalization and are still limited by the poor category coverage of the “paired” training data.

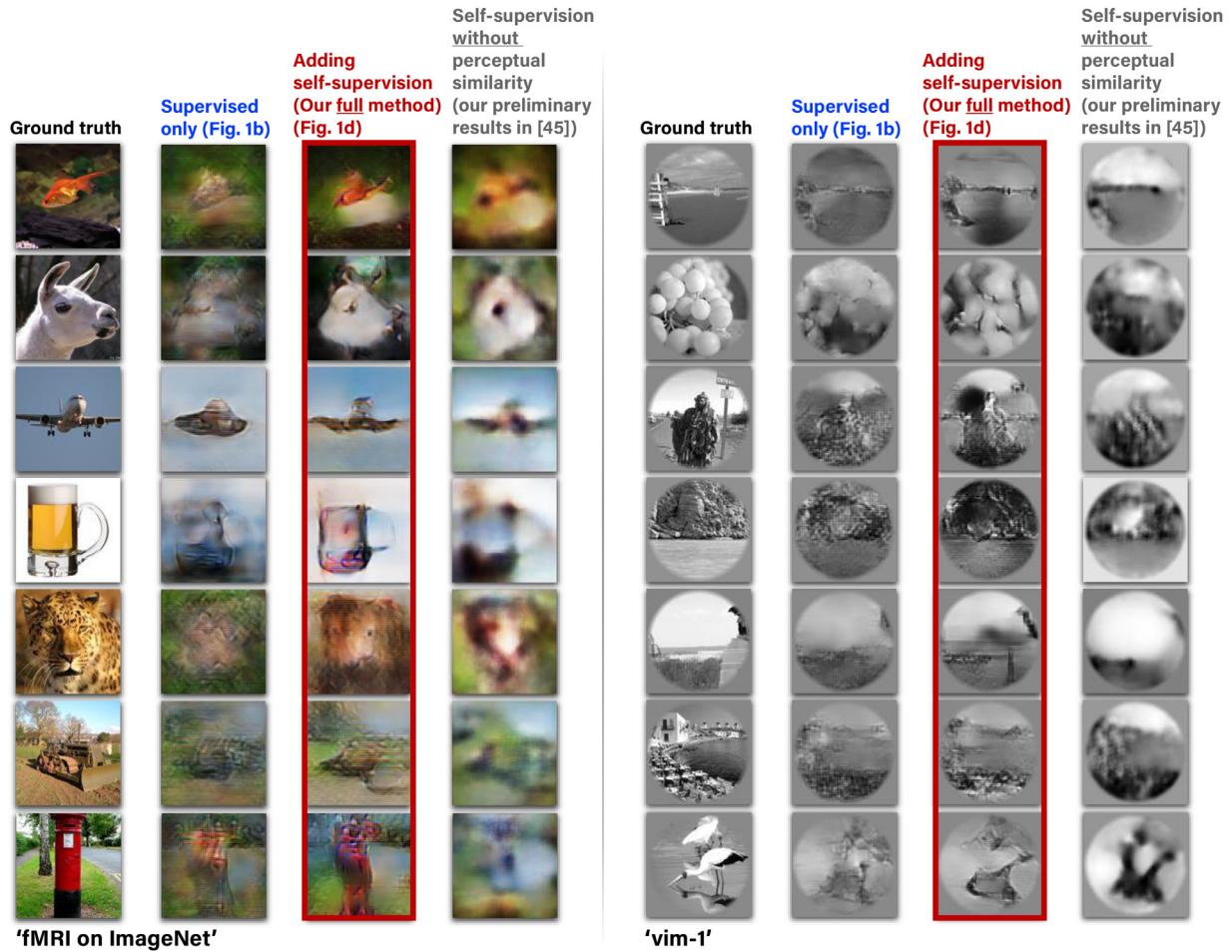
We present a new approach to overcome the limitations mentioned above and inherent lack of training data, simultaneously for both tasks – image reconstruction as well as for *large-scale semantic classification*. We achieve this by introducing *self-supervised training on unpaired images* (*images without fMRI recordings*). Our approach is illustrated in Fig 1. We train two types of deep neural networks: an Encoder *E*, to map natural images to their corresponding fMRI response, and a Decoder *D*, to map fMRI recordings to their corresponding images. Composing those two networks, *E-D*, yields a combined network whose input and output are the same image (Fig 1d). *This allows for unsupervised training on unpaired images* (i.e., images without fMRI recordings, e.g., 50,000 natural images from 1000 semantic categories in our experiments). Such self-supervision adapts the Decoder to the statistics of novel natural images and their novel categories (including categories not represented in the “paired” training data or even in the “unpaired” natural images). Importantly, we combine our self-supervised approach with semantics learning **without requiring any explicit class labels in the training**. This is achieved via: (i) Using perceptual similarity reconstruction criteria. These criteria require for the reconstructed image to be similar to the original image not only at a pixel level, but also that the two images be ***perceptually similar***, in their higher-level deep “semantic” representations. (ii) Designing an Encoder architecture that benefits from a higher-level “semantic” representation of a backbone network that was trained for image classification. Fig 2 exemplifies the power of adding unsupervised training on unpaired images together with the perceptual criteria. Furthermore, beyond the dramatic improvement in the image reconstruction task, we show that our combined self-supervised perceptual approach gives rise also to new capabilities in semantic category decoding (against 1000 classes), despite the scarce fMRI-based training data.

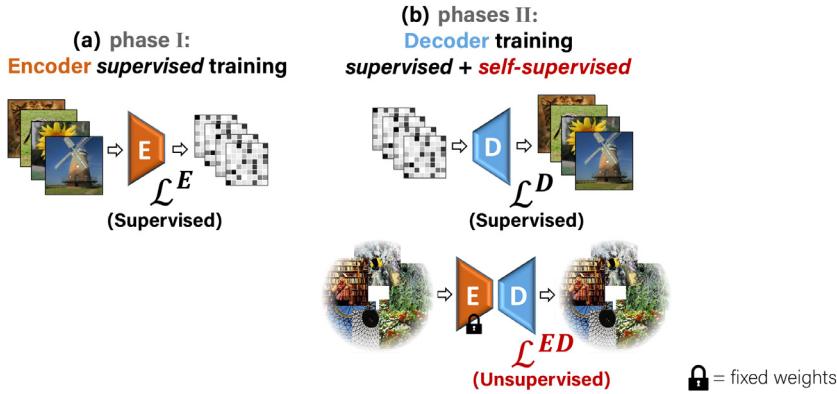
Our self-supervision leverages the cycle-consistency nature of our two networks when combined. Notably, our self-supervised training on unpaired natural images differs from an auto-encoder: When training under this configuration the Encoder weights are kept fixed, thus constraining the input images to go through the true fMRI space. Furthermore, unsupervised training on unpaired natural images was also recently proposed in St-Yves and Naselaris (2019), where they used these images to produce additional surrogate fMRI-data to train their model. However, this was never addressed in the context of semantic decoding from fMRI data. Moreover, their image reconstruction criteria were limited to pixel level alone and did not include any type of perceptual or visual-features based reconstruction criteria. To the best of our knowledge, we are the first to present classification of semantic categories that are never-before-seen during training at a large-scale (1000-way – i.e., detecting the correct class out of more than 1000 rich classes).

A preliminary version of our self-supervised approach and partial results (in the context of image reconstruction only) were previously presented in a conference proceeding (Belyi et al., 2019). However, here we present our complete and advanced reconstruction algorithm that has undergone major extensions, enabling a leap improvement in the image reconstruction quality over our previous method (Belyi et al., 2019) (as demonstrated in Fig 2). Moreover, our new extended framework now enables also impressive *semantic classification capabilities* from fMRI data. As an aid for readers familiar with our previous report (Belyi et al., 2019), we list the four major extensions introduced in the present paper: (i) We introduced two significant improvements to our algorithm: Adding high-level perceptual criteria (Zhang et al., 2018b) on the reconstructed images (in contrast with optimizing Mean-Square-Error loss,



**Fig. 1. Our self-supervised approach.** (a) The task: reconstructing images and classifying their semantic category from evoked brain activity, recorded via fMRI. (b,c) Supervised training for decoding (b) and encoding (c) using limited training pairs. This gives rise to poor generalization. (d) Illustration of our added self-supervision, which enables training on “unpaired images” (any natural image with no fMRI recording). This self-supervision allows adapting the decoder to the statistics of natural images and many rich semantic classes.





**Fig. 3. Training phases.** (a) The first training phase: Supervised training of the Encoder with {Image, fMRI} pairs. (b) Second phase: Training the Decoder with two types of data simultaneously: {fMRI, Image} pairs (supervised examples), and unpaired natural images (self-supervision). The pretrained Encoder from the first training phase is kept fixed in the second phase.

and on low-level features alone), and increasing the expressiveness of the Encoder architecture to include higher-level “semantic” representations by using multiple levels of the pretrained VGG network (in contrast with a single readout layer before). Our present algorithm provides state-of-the-art reconstructions compared to leading existing methods to-date, as evident through image-metric-based comparisons as well as extensive human behavioral evaluations. (ii) We extended the self-supervised approach to allow also for semantic classification of reconstructed images. This classification relies on and demonstrates the fidelity of the reconstructions. (iii) We analyze the contributions of different visual cortex areas to the resulting image reconstructions, and (iv) We evaluate the biological consistency of the learned models.

Importantly, adding the perceptual similarity was powerful enough to shadow and marginalize the previously dominant effect of training on unpaired fMRI data in Belyi et al. (2019). Accordingly, training on unpaired fMRI is omitted in the present framework (despite being a legitimate feature enabled by our self-supervised approach).

Our contributions are therefore several-fold:

- A self-supervised approach for simultaneous image reconstruction and semantic category decoding, which can handle the inherent lack of image-fMRI training data.
- Unprecedented state-of-the-art image-reconstruction quality from fMRI of never-before-seen images (from never-before-seen semantic categories).
- Large-scale semantic classification (1000+ rich classes) of *never-before-seen semantic categories*. To the best of our knowledge, such large-scale semantic classification capabilities from fMRI data has never been demonstrated before.
- We provide analyses showing predominance of early visual areas in image reconstruction quality, and demonstrate biologically plausible receptive field formation in our learned models (although such constraints were not imposed in the training, but rather automatically learned by our system).

## 2. Overview of the approach

This section provides a high-level overview of our approach and its key components. The technical details and the experimental datasets are described in Methods.

### 2.1. Self-supervised image reconstruction from brain activity

The essence of our approach is to enrich the scarce paired image-fMRI training data with easily accessible natural images for which *there are no fMRI recordings*. This type of training is enabled by imposing cycle-consistency on the “unpaired images”, using two networks, which learn two inverse mappings: from images to fMRI (Encoding), and vice versa – from fMRI to images (Decoding).

Our training consists of Encoder training followed by Decoder training, which we define in the two phases illustrated in Fig 3. In the

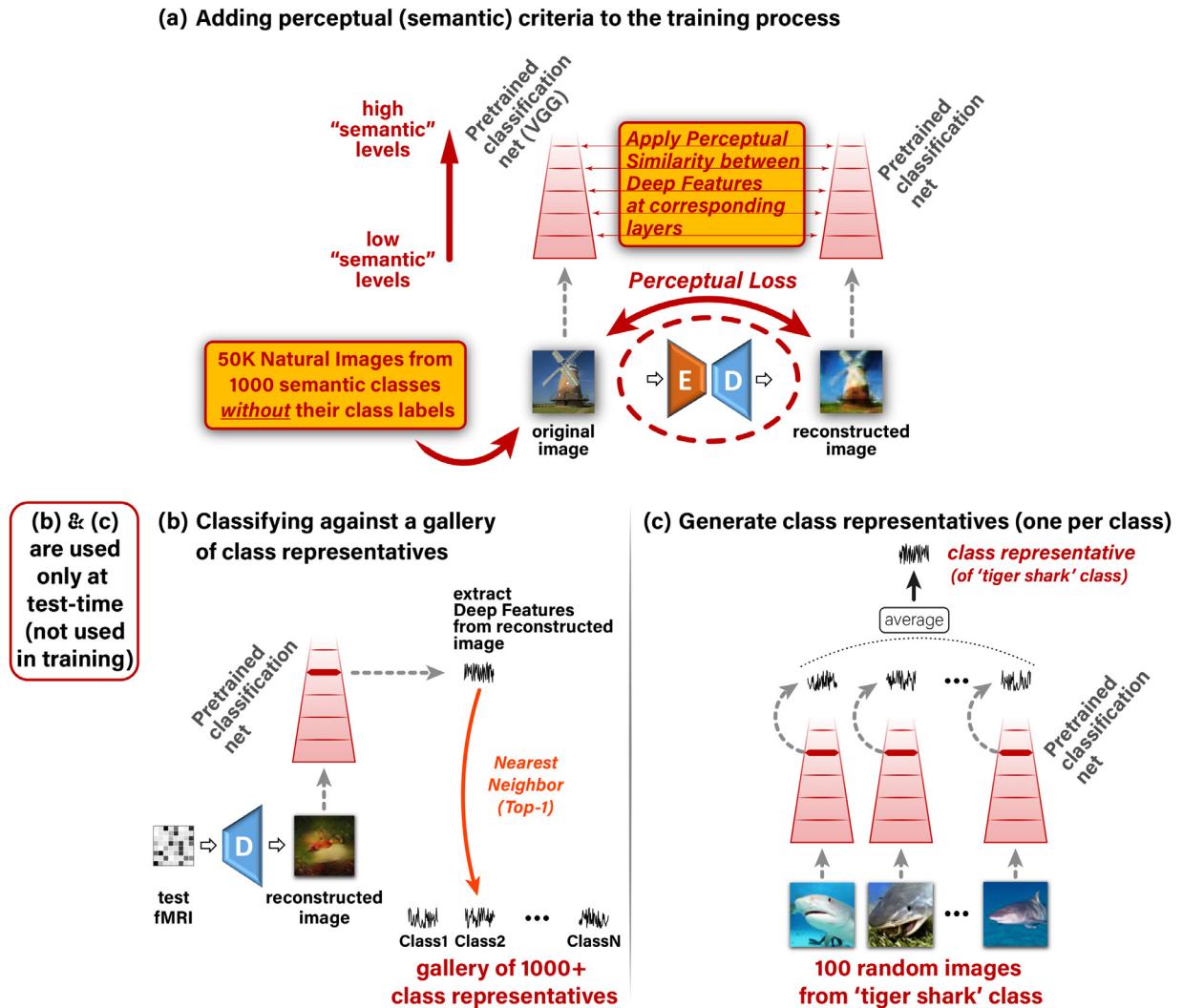
first phase, we apply supervised training of the Encoder  $E$  alone. We train it to predict the fMRI responses of input images using the image-fMRI training pairs (Fig 3a). In the second phase, we use the pretrained Encoder (from the first phase) and train the Decoder  $D$ , keeping the weights of  $E$  fixed (Fig 3b).  $D$  is trained using both the paired and the unpaired data, simultaneously. Here, each training batch consists of two types of training data: (i) image-fMRI pairs from the training set (Fig 1b), and (ii) unpaired natural images (with no fMRI, Fig 1d).

Training on *unpaired natural images* (without fMRI) allows to augment the training with data from a much richer semantic space than the one spanned by the paired training data alone. Specifically, we draw the unpaired images from a large *external database* of 49K images from 980 ImageNet (“ILSVRC”) classes, which are mutually exclusive not only to the test-images contained in the fMRI (paired) dataset, but also to their underlying test classes (test categories). In principle, for optimal networks  $E$  and  $D$ , the combined  $E - D$  network should yield an output image which is identical to its input image. This should hold for any natural image (regardless if an fMRI was ever recorded for it). Importantly, our image reconstruction losses (in both  $\mathcal{L}_{ED}$  and  $\mathcal{L}_D$ , Fig 3) require for the reconstructed image to be similar to the original image not only at a pixel level. We further require these two images to be *perceptually similar*, at a semantic level.

Perceptual Similarity is a metric that highly correlates with *human* image-similarity perception, and involves a broad range of visual feature representation levels (Zhang et al., 2018b). Our Perceptual Similarity loss is illustrated in Fig 4. We apply a pre-trained VGG classification network (a network which was trained for the task of object recognition from images (Simonyan and Zisserman, 2014)), on the reconstructed and the original images. We then impose similarity between their corresponding deep-image-features, extracted from multiple deep layers of VGG. *Using this metric as our reconstruction criterion enables to learn low-level to high-level “semantic” information from the broad semantic space, which is spanned by the external database of ‘unpaired’ images and their classes; It, thus, encourages the Decoder to output images that are not only accurate at low-level visual features, but are also semantically meaningful.* Importantly, the *class labels* of the unpaired images (or the paired images) are never used in our training process, hence may be unknown. Adding the perceptual loss in combination with training on additional unpaired images gives rise to a leap improvement in the fMRI-to-image reconstruction quality compared to any previous method (including our previous method (Belyi et al., 2019) – see Fig 2, as well as compared with other methods – see Fig 8). It provides a dramatic improvement in detail level and perceptual interpretability of the reconstructed images. Our new approach further enables large-scale semantic classification of fMRI into rich novel semantic categories.

### 2.2. Self-supervised image classification to semantic categories

Our new self-supervised perceptual approach extends well beyond the task of image reconstruction. It further allows for large-scale



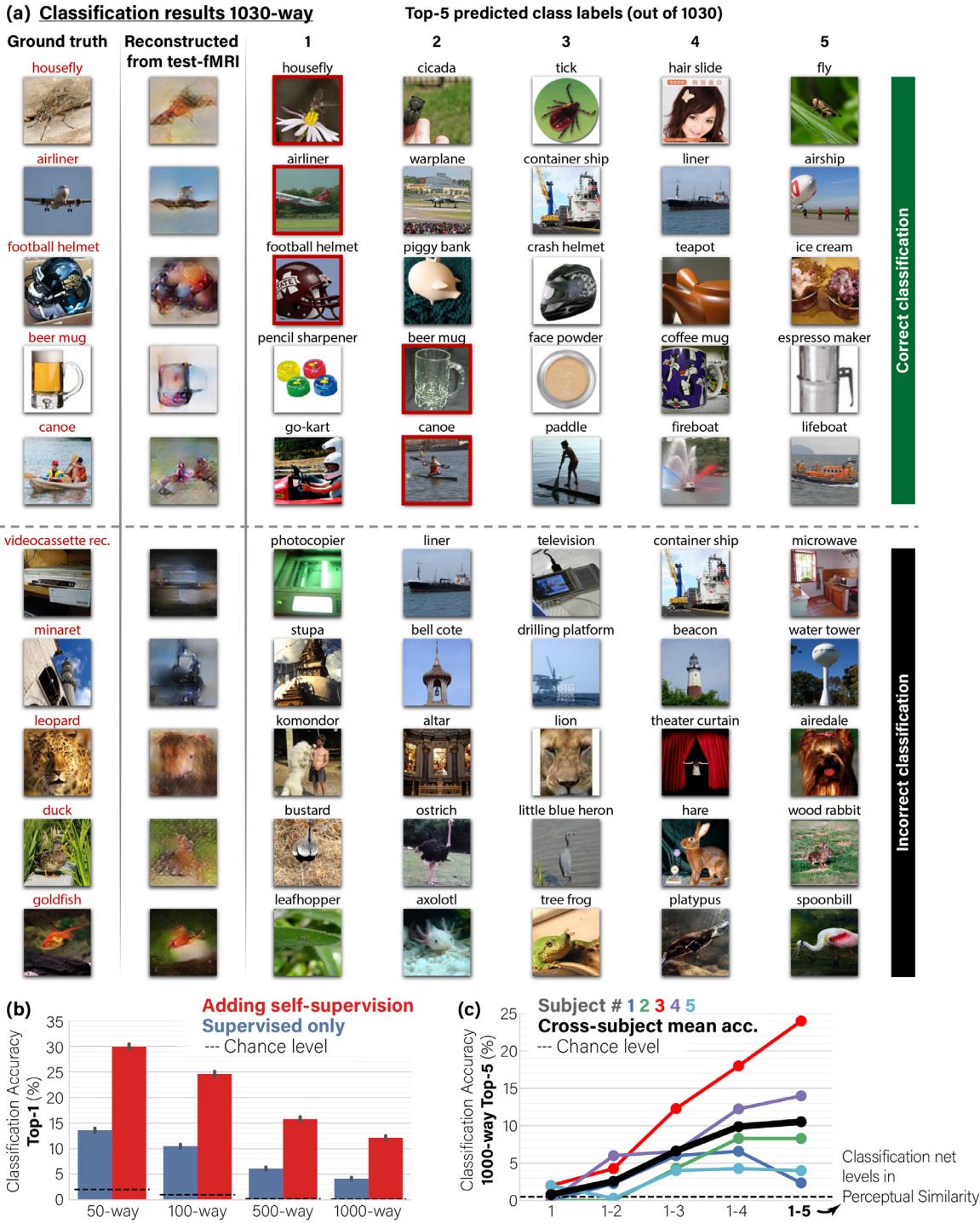
**Fig. 4. Adding high-level perceptual criteria improves reconstruction accuracy and enables large-scale semantic classification.** (a) Imposing Perceptual Similarity on the reconstructed image at the Decoder's output, as applied when training on unpaired natural images (*without* fMRI and *without* any class labels) from many novel semantic classes. This adapts the Decoder to a significantly broader semantic space despite not having any explicit semantic supervision. (b) To classify a reconstructed image to its novel semantic class we extract Deep Features using a pretrained classification network, and follow a nearest-neighbor class-centroid approach against a large-scale gallery of 1000+ ImageNet classes. (c) We define class representatives as the mean-embedding of many same-class images (Horikawa and Kamitani, 2017a).

semantic classification of fMRI data. We present classification of fMRI data against a gallery of more than 1000 rich image classes, in a challenging 1000-way classification task (see Fig 5). Such large-scale semantic classification of fMRI data (to 1000-way, with promising results) has never been demonstrated before.

Our classification approach is based on our self-supervised perceptual reconstruction method described above. We use our perceptually trained Decoder to reconstruct the test-images from their test-fMRI. We then classify the reconstructed images against 1000+ rich ImageNet semantic classes. Fig 4b,c shows our classification approach. To classify a reconstructed image to its novel semantic class, we match a “Deep-Feature signature” extracted from the reconstructed image, against “class-representative Deep-Feature signatures” (one per class), in a gallery of 1000+ semantic categories, which also include the 50 *novel* test classes. More specifically: (i) We extract Deep Features from the reconstructed image at an intermediate level of a pretrained classification network (Fig 4b), (ii) Following Horikawa and Kamitani (2017a), for each class in the gallery, we compute a single “class representative” using 100 randomly sampled images from that class. The class representative is defined as the average Deep Features (centroid) of those 100 randomly sampled images from that class (see Fig 4c). (iii) We com-

pute the correlation between the Deep Features extracted from the reconstructed image and each of the 1000+ class representatives. These yield 1000+ “semantic similarity” scores (Fig 4b). Ranking the gallery classes according to these similarity scores (for each reconstructed image) provides the basis for semantic classification at any desired ‘Top-X’ accuracy level (Fig 5). Specifically, the classification is marked ‘correct’ when the ground truth category obtains the best similarity score among all 1000+ classes (‘Top-1’), or when it is among the top 5 most similar classes (‘Top-5’). The location of the ground-truth class within the sorted list of 1000+ classes further provides a “rank score” for evaluating our classification accuracy (see Table 2).

This classification procedure greatly benefits from our self-supervised perceptual approach, which enables to train on *additional unpaired images from arbitrarily many novel semantic categories* (Fig 1b). This allows to adapt the Decoder to a much richer (practically unlimited) semantic coverage in a **completely category-free way**, namely without any explicit semantic supervision during training. The key component which induces semantically meaningful decoding is the Perceptual Similarity metric, which involves higher-level “semantic” criteria. This type of non-specific semantic supervision enables our method to generalize well to new never-before-seen semantic classes – classes which are



**Fig. 5. Self-supervision allows classification to rich and novel semantic categories.** (a) Visual classification results showing the Top-5 predictions out of 1030 classes for reconstructed test-image. We show examples where the ground-truth class (marked in red) is ranked among the Top-5 (**correct classification**) or excluded from it (**incorrect classification**). For visualization purposes only, each class is represented by the nearest-neighbor image from the 100 randomly sampled images of the particular class. Note that "incorrect" predicted classes are often reasonable (e.g., "Leopard" wrongly predicted as "Lion"; "Duck" wrongly predicted as "Ostrich"). (b) Top-1 Classification accuracy in an  $n$ -way classification task. Adding unsupervised training on unpaired data (Fig 1d,e) dramatically outperforms the baseline of the supervised approach (Fig 1b). (c) Ablation study of the Classification accuracy as a function of the Perceptual Similarity criterion for decoder training: Applying "partial" perceptual similarity using only the outputs of the first VGG16 block (low "semantic" layers), and up to all its 5 blocks (high "semantic" layers). Applying full Perceptual Similarity on higher-level VGG features substantially improves classification performance. Panels a, b show results for subject three. 95% Confidence Intervals shown on charts.

**Table 1**

**Summary of fMRI datasets used in analyses.** Repeat count refers to the number of fMRI recordings per presented stimulus. Voxel count refers to approximated number of voxels used in analysis.

fMRI Dataset	N train images (K repeats)	N test images (K repeats)	N voxels
fMRI on ImageNet (Horikawa and Kamitani, 2017a)	1200 (1)	50 (35)	4500
vim-1 (Kay et al., 2008)	1750 (2)	120 (13)	8500

**Table 2**

**Classification rank (out of 1030) for ‘fMRI on ImageNet’ (Horikawa and Kamitani, 2017a)** Self-supervision allows classification to rich and novel semantic categories. Median rank of the ground truth class among 1030 class representatives (Lower is better). Significant differences between the two methods are marked with asterisks (Mann-Whitney test,  $N = 50$ ,  $p < .05$ ). Adding self-supervision leads to significant improvement in classification rank for the four (out of five) subjects with the highest fMRI median noise-ceiling. SD and SE are standard deviation and error, respectively.

fMRI data noise-ceiling median (SD)	Supervised only median rank (SE)	Adding self-supervision median rank (SE)
Subject 1	0.56 (0.28)	136.0 (48.3)
Subject 2	0.57 (0.30)	* 85.0 (31.1)
Subject 3	0.73 (0.28)	* 67.0 (32.6)
Subject 4	0.68 (0.30)	* 72.0 (25.5)
Subject 5	0.58 (0.29)	* 71.0 (34.9)

neither contained in the paired training data, nor in the unpaired external images used at train time.

### 3. Methods

#### 3.1. Self-supervised Encoder/Decoder training

The Encoder (E) is trained first on predicting fMRI responses from input images. Once trained, the encoder is fixed, and the Decoder (D) is trained on reconstructing images from fMRI responses. The motivation for training the Encoder and Decoder in separate phases (with a fixed Encoder during Decoder training) is to ensure that the Encoder’s output does not diverge from predicting the fMRI responses by the unsupervised training objectives 1 d. Additionally, we start by supervised training of the Encoder in order to allow it to converge at the first phase, and then serve as strong guidance for the more severely ill-posed decoding task (Naselaris et al., 2011), which is the focus of the next phase. We next describe each phase in more detail.

##### 3.1.1. Encoder supervised training (Phase I)

The supervised training of the Encoder is illustrated in Fig 3a. Let  $\hat{r} = E(s)$  denote the encoded fMRI response from image,  $s$ , by Encoder  $E$ . We define an fMRI loss by a convex combination of mean square error and cosine proximity with respect to the ground truth fMRI,  $r$ . The **fMRI loss** is defined as:

$$\mathcal{L}_r(\hat{r}, r) = \alpha \cdot \text{MSE}(\hat{r}, r) - (1 - \alpha) \cos(\angle(\hat{r}, r)), \quad (1)$$

where  $\alpha$  is a hyperparameter set empirically ( $\alpha = 0.9$ ). We use this loss for training the Encoder  $E$ .

Notably, in the considered fMRI datasets, the subjects who participated in the experiments were instructed to fixate at the center of the images. Nevertheless, eye movements were not recorded during the scans, thus the fixation performance is not known. To account for the center-fixation uncertainty, we introduced small random shifts (+/- a few pixels) of the input images during Encoder training (Shen et al., 2019a). This resulted in a substantial improvement in the Encoder performance and subsequently in the image reconstruction quality. Upon completion of Encoder training, we transition to training the Decoder together with the fixed Encoder.

##### 3.1.2. Decoder training (Phase II)

The training loss of our Decoder consists of two main losses illustrated in Fig 3b:

$$\mathcal{L}^D + \mathcal{L}^{ED}. \quad (2)$$

$\mathcal{L}^D$  is a supervised loss on training pairs of image-fMRI.  $\mathcal{L}^{ED}$  (Encoder-Decoder) is an unsupervised loss on unpaired images (without corresponding fMRI responses). Both components of the loss are normalized to have the same order of magnitude (all in the range [0, 1], with equal weights), to guarantee that the total loss is not dominated by any individual component. We found our reconstruction results to be relatively insensitive to the exact balancing between the two-loss components. We next detail each component of the loss.

$\mathcal{L}^D$ : **Decoder Supervised Training** is illustrated in Fig 1b. Given {fMRI, Image} training pairs  $\{(r, s)\}$ , the supervised loss  $\mathcal{L}^D$  is imposed on the decoded image,  $\hat{s} = D(r)$ , and is defined via the image reconstruction objective,  $\mathcal{L}_s$ , as

$$\mathcal{L}^D = \mathcal{L}_s(\hat{s}, s).$$

$\mathcal{L}_s$  consists of losses on image RGB values,  $\mathcal{L}_{RGB}$ , as well as losses on Deep Image Features extracted from the image using a pretrained VGG16 network (Simonyan and Zisserman, 2014) (a deep network that was trained for the task of object recognition from images). Using features of pretrained deep neural network classifiers was recently introduced to the task of reconstructing images from fMRI responses in (Shen et al., 2019a; 2019b), although our implementation is somewhat different (for a comparison of the approaches, see the Discussion section). We denote the deep features extracted from an image,  $s$ , by  $\varphi(s)$ , on which we apply a **Perceptual Similarity** criterion,  $\mathcal{L}_{perceptual}$ . This last component gave a significant performance leap. Unlike our preliminary work (Belyi et al., 2019), where we imposed only a Mean-Square-Error loss on the low level features alone (hence failed to capture or exploit any “semantic” appearance or interpretation), here we impose Perceptual similarity (Zhang et al., 2018b) using the outputs of all the five feature-extractor blocks of VGG (from low to high VGG layers, i.e., lower-to-higher “semantic” levels), denoted as  $\varphi_{vgg\_blocks}^b(s)$  for a particular block  $b$ . This metric is implemented by cosine similarity between channel-normalized ground-truth and predicted features at each block output. The complete criterion is then a sum of the block-wise contributions. The Image loss for a reconstructed image  $\hat{s}$  thus reads:

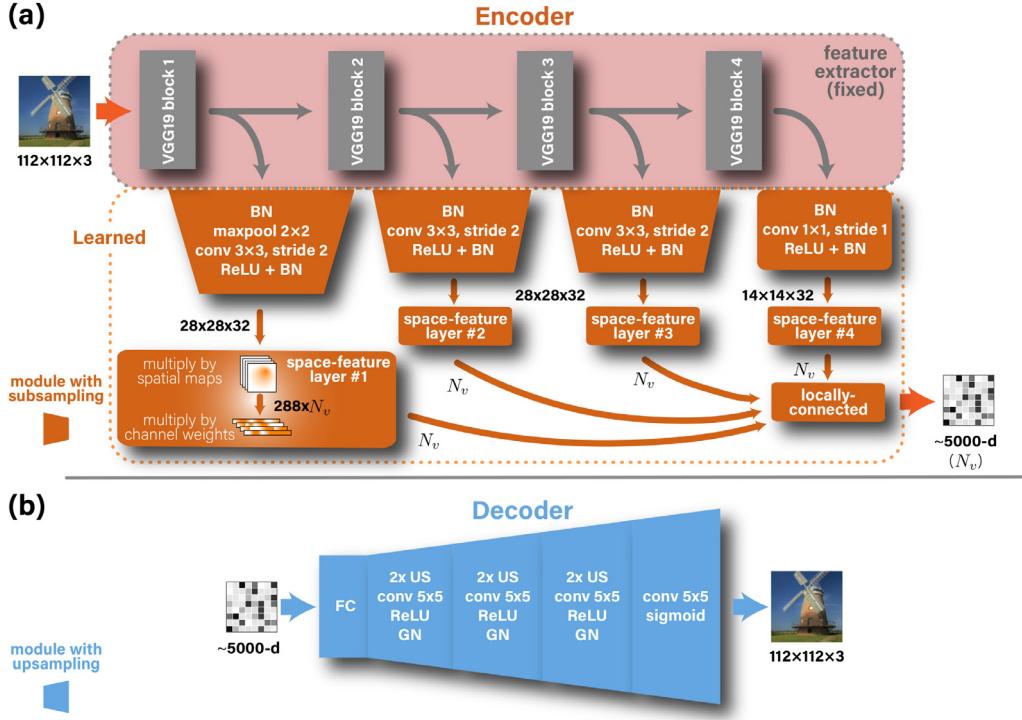
$$\mathcal{L}_s(\hat{s}, s) = \mathcal{L}_{RGB}(\hat{s}, s) + \mathcal{L}_{perceptual}(\hat{s}, s) + \mathcal{R}(\hat{s}) \quad (3)$$

where:

$$\begin{cases} \mathcal{L}_{RGB}(\hat{s}, s) \propto \|\hat{s} - s\|_1 \\ \mathcal{L}_{perceptual}(\hat{s}, s) \propto -\sum_{b=1}^5 \cos(\angle(\varphi_{vgg\_blocks}^b(\hat{s}), \varphi_{vgg\_blocks}^b(s))) \end{cases} \quad (4)$$

The last term,  $\mathcal{R}(\hat{s})$ , corresponds to total variation (TV) regularization of the reconstructed (decoded) image,  $\hat{s} = D(r)$ . The Image loss in Eq. 3 is also used for the self-supervised Encoder-Decoder training on unpaired images (images without fMRI), as explained next.

$\mathcal{L}^{ED}$ : **Encoder-Decoder training on unpaired Natural Images** is illustrated in Fig 1d. This objective enables to train on any desired unpaired image (images for which no fMRI was ever recorded), well beyond the 1200 or 1750 images included in the fMRI dataset (‘fMRI on ImageNet’ or ‘vim-1’ respectively). In particular, we used ~50K additional natural images from ImageNet’s 1000-class data (Deng et al., 2009) (excluding the test classes). This allows adaptation to the statistics



**Fig. 6. Encoder & Decoder Architectures.** BN, GN, US, and ReLU stand for batch normalization, group normalization, up-sampling, and rectified linear unit, respectively. We designed a custom space-feature locally-connected layer (see text).

of many more novel semantic categories, thus learning the common higher-level feature representation of various novel classes. To train on images without corresponding fMRI responses, we map images to themselves through our Encoder-Decoder transformation,

$$s \mapsto \hat{s}_{ED} = D(E(s)).$$

The unsupervised component  $\mathcal{L}^{ED}$  of the loss in Eq. 2 on unpaired images,  $s$ , reads:

$$\mathcal{L}^{ED} = \mathcal{L}_s(\hat{s}_{ED}, s),$$

where  $\mathcal{L}_s$  is the Image loss defined in Eq 3. In other words,  $\mathcal{L}^{ED}$  imposes cycle-consistency on any natural image, but at a perceptual level (not only at the pixel level). Note that this high-level perceptual consistency does not require using any class labels in the training.

### 3.2. Deep architecture

We focused on  $112 \times 112$  RGB or grayscale image reconstruction (depending on the dataset), although our method works well also on other resolutions.

**Architectures** of the Encoder and the Decoder are illustrated in Fig 6. The Encoder comprises four parallel branches of representation, built on top of features extracted from blocks 1–4 of VGG19. This enables to benefit from the hierarchy of “semantic” levels of the pretrained VGG network. The outputs of the four resulting branches (with their various resolutions) are then fed into branch-specific learned convolutional modules, which are designed to reduce the representation’s dimensions to more compact representations of  $28 \times 28 \times 32$  or  $14 \times 14 \times 32$  (Height×Width×ConvolutionChannels). These modules consist of batch normalization,  $3 \times 3$  convolution with 32 channels, ReLU,  $\times 2$  subsampling, and batch normalization. The first branch is preceded by an additional  $\times 2$  maximum pooling while the fourth branch is not subsampled. Inspired by the feature-weighted receptive field (St-Yves and Naselaris, 2017) for complex feature spaces, we designed a locally-connected layer which acts on the spatial and channel dimensions separately. This

induced separability enables a dramatic decrease in the number of parameters required to regress the voxel activations. In this spatial-feature locally-connected layer, for each spatial coordinate we stack along the channel dimension the values of the immediate 9 neighboring coordinates. This results in tensors of shapes  $26 \times 26 \times 288$  for branches 1–3 and  $12 \times 12 \times 288$  for branch 4 (after eliminating the boundaries), and effectively constrains spatial association among adjacent pixels within  $3 \times 3$  patches. Next, the tensors are multiplied by voxel-specific spatial maps, which learn the voxels’ receptive fields. This outputs space-reduced tensors of shape  $288 \times N_v$ , where  $N_v$  is the number of voxels. To encourage the smoothness of the spatial maps, we penalize for the within-map total variation. The next layer is a cross-channel locally-connected layer, which weighs the contribution of each feature/channel per voxel and outputs a vector of size  $N_v$ . This vector is produced by every branch separately. Finally, the outputs of the 4 branches are concatenated along a new dimension and followed by a locally-connected layer ( $4 \times N_v$  parameters) that weighs the contribution from each branch and output the final fMRI activations. We initialized all weights using Glorot normal initializer, except for the last layer, which was 1-initialized (and restricted to remain non-negative).

The Decoder architecture uses a locally-connected layer to transform and reshape the input vector-form fMRI input into 64 feature maps with spatial resolution  $14 \times 14$ . This representation is then followed by three blocks, and each consists of: (i)  $\times 2$  up-sampling, (ii)  $5 \times 5$  convolution with unity stride, 64 channels, and ReLU activation, and (iii) group normalization (16 groups). To yield the output image we finally performed an additional convolution, similar to the preceding ones, but with three channels to represent colors, and a sigmoid activation to keep the output values in the 0–1 range. We used Glorot-normal (Glorot and Bengio, 2010) to initialize the weights.

### 3.3. Experimental datasets

We tested our self-supervised approach on two publicly available (and very different) benchmark fMRI datasets that we summarize in

**Table 1.** These datasets provide elicited fMRI recordings of human subjects paired with their corresponding underlying natural images. In both datasets, subjects were instructed to fixate at the center of the images. The same architectures and hyperparameters were used for both datasets. In ‘fMRI on ImageNet’ dataset (Horikawa and Kamitani, 2017a), which was used also for image classification, 1250 images were drawn from 200 selected ImageNet categories. 150 categories (classes) were used as training data (8 images per category – altogether 1200 training images). The 50 remaining image categories were designated as the novel test categories, represented by 50 test images (1 image from each test category) to be recovered from their fMRI recordings (“test-fMRI”). Importantly, we averaged across same-stimulus repeated fMRI recordings, when available (Table 1). This gave rise to a higher signal-to-noise ratio fMRI samples, and thus to best reconstruction and classification performance (see Supplementary-Material for single-trial reconstruction, identification, and classification results). This further enabled direct comparison with previous works (Shen et al., 2019a; 2019b; St-Yves and Naselaris, 2019).

**External (unpaired) images database.** For unsupervised training on unpaired images (Encoder-Decoder objective, Fig 1d) we used additional 49K natural images from 980 classes of ImageNet (“ILSVRC”) train-data (Deng et al., 2009). We verified that the images and categories in our additional unpaired external dataset do not overlap with the test-images and test-categories in the ‘fMRI on ImageNet’ dataset (the inference target). Since the 50 test classes of ‘fMRI on ImageNet’ (Horikawa and Kamitani, 2017a) partially overlap with the 1000 original ILSVRC classes, we particularly discarded the 20 overlapping classes. At test-time, we tested classification against 1030 classes (980 + 50).

## 4. Results

To test the feasibility of our approach we experimented with two publicly available benchmark fMRI datasets: (i) **fMRI on ImageNet** (Horikawa and Kamitani, 2017a), and (ii) **vim-1** (Kay et al., 2008).

### 4.1. Image reconstruction from fMRI

Fig 2 shows our results with the proposed method, including the combined supervised and self-supervised training with perceptual criteria. These results (in red frames – 3rd column) are contrasted with the results obtainable when using supervised training only (e.g., the 1200 paired training examples of the “fMRI on ImageNet” dataset – 2nd column), as well as when not using perceptual criteria (4th column). All the displayed images were reconstructed from the dataset’s test-fMRI cohort (fMRI of new images from never-before-seen semantic categories). The red-framed images show many faithfully-reconstructed shapes, textures, and colors, which depict recognizable scenes and objects. In contrast, using the supervised objective alone led to reconstructions that were considerably less recognizable (2nd column). Similarly, excluding the perceptual criteria as well as the extended Encoder architecture led to reconstructions that were less perceptually understandable (4th column). The reconstructions of the entire test cohort (50 images in the “fMRI on ImageNet” dataset) as well as an assessment of the separated contribution of the perceptual criteria and the extended Encoder architecture can be found in the Supplementary-Material.

To verify that our method can successfully be applied to different subjects, Fig 7 shows the reconstructions for all five subjects in the “fMRI on ImageNet” dataset. Note that using fMRI data of different subjects give rise to varying quality of reconstruction as driven by the varied SNR in the subjects’ fMRI data (which is substantially affected by the subject’s ability to maintain fixation at the center pixel in all repeated scans of the same image). Nevertheless, clear and common identifying markers of the ground truth image appear across all subjects. The red frame indicates the results for the best subject (Subject 3 in the dataset)

which is the subject of focus in the remaining parts of this paper (unless mentioned otherwise).

#### 4.1.1. Comparison with state-of-the-art image reconstruction methods

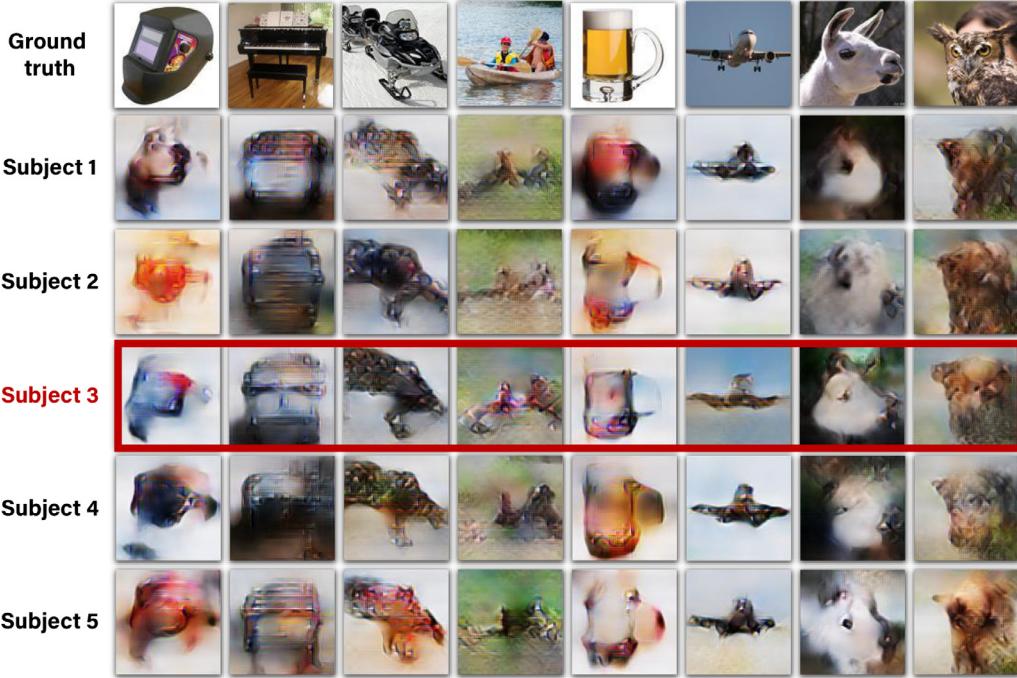
We compared our reconstruction results against the two leading methods: Shen et al. (Shen et al., 2019b)) and St-Yves et al. (St-Yves and Naselaris, 2019) – each on its relevant dataset. Fig 8a,b compares the results of our method with those two methods (both of which are deep-learning GAN-based methods). Visual comparison of Shen et al. (2019b); St-Yves and Naselaris (2019) with our method (Fig 8a,b) highlights that despite their natural-like visual appearance enforced by the GAN, the reconstructed images of Shen et al. (2019b); St-Yves and Naselaris (2019) are often not faithful to the underlying ground truth image. We further report quantitative comparisons, both by image-metric-based evaluation and by human visual evaluation for the top-SNR subject from each dataset (Subject 3 in ‘fMRI on ImageNet’ dataset (Horikawa and Kamitani, 2017a); Subject 1 in ‘vim-1’ dataset (Kay et al., 2008)). Our quantitative evaluations are based on an n-way identification task (Ren et al., 2021; Seeliger et al., 2018; Shen et al., 2019b; Zhang et al., 2018a). Namely, each reconstructed image is compared against  $n$  candidate images (the ground truth image, and  $(n - 1)$  other randomly selected images), and the goal is to identify its ground truth. We considered two identification methods under this task:

(i) **Image-metric-based** identification performed using the Perceptual Similarity metric (Zhang et al., 2018b) (between the reconstructed image and the candidate image). Each reconstructed image is compared against  $n$  images, and the nearest-neighbor candidate image under this metric was determined to be the identified “correct” image. Panels 8 c,d show the correct-identification rate (for each method separately) for n-way identification tasks for  $n = 2, 5, 10, 50$ . We evaluate our method and two variants of the method provided by Shen et al. (2019b) on the “fMRI on ImageNet” benchmark dataset (Fig 8c). In 2-way identification task our method scored an accuracy of 98.5% ( $SEM^1 = 0.6\%$ ,  $N = 50$ ). Our accuracy remains quite high for all  $n$ , outperforming both variants of Shen et al. (2019b) by a margin of 8–46% across all task difficulty levels ( $n = 2, 5, 10, 50$ ). We repeated the analysis for “vim-1” fMRI dataset (Fig 8d), where our method scored an accuracy of 87.5% ( $SEM = 1.7\%$ ,  $N = 120$  2-way task), outperforming the method from St-Yves and Naselaris (2019) by a large margin of 28–40% across the same difficulty levels. Particularly in the challenging 50-way task our method achieved striking leaps: outperforming the baseline by more than x2 prediction accuracy in “fMRI on ImageNet”, and more than x10 better prediction accuracy in “vim-1”.<sup>2</sup> Importantly, the statistical power of these findings generalizes beyond the specific 50 test examples (image-fMRI) in “fMRI on ImageNet”, or the 120 test examples in “vim-1” (Wilcoxon test,  $N = 50, 120$ ,  $p < 10^{-5}$ ).

(ii) **Human-based** identification. Panels 8 e,f show reconstruction evaluation results when repeating the same quantitative comparison approach, but this time outsourcing the n-way identification task ( $n = 2, 5, 10$ ) to random human raters. We used Mechanical Turk to launch surveys to new 45 raters for each evaluated method. Our method scored 90.4% ( $SEM = 0.6\%$ ,  $N = 45$ ) and 85.0% ( $SEM = 1.5\%$ ,  $N = 45$ ) in a 2-way identification task on “fMRI on ImageNet” and “vim-1” respectively; Scaling the task difficulty up to 10-way, our method scored 65.0% ( $SEM = 3.0\%$ ,  $N = 45$ ) and 44.3% ( $SEM = 2.5\%$ ,  $N = 45$ ). Overall our method significantly outperformed the previous methods, on both datasets, and across difficulty levels by a margin of at least 8.7% (Mann-Whitney test,  $N = 45$ ,  $p < .001$ ).

Notably, the identification accuracy of each reconstructed image when using the image-metric for evaluation, mostly depended on the reconstruction quality of the specific image, and was robust to randomizing the selection of the (non-ground-truth) candidate images. Furthermore, the choice of candidate images in the human-based evaluation was fixed

<sup>2</sup> Standard Error of the Mean.



**Fig. 7.** Reconstructions for all five subjects in “fMRI on ImageNet” (Horikawa and Kamitani, 2017a). Reconstructed images when using the full method, which includes training on unpaired data (d,e). Reconstruction quality varies across subjects, depending on noise-ceiling/SNR of subjects ” data (voxel median noise-ceiling for subjects 1–5:0.56, 0.57, 0.73, 0.68, 0.58). Subject 3 (in the dataset), which is framed above in red, is the subject of focus in the remaining parts of this paper unless remarked otherwise.

across the raters and the methods we compared. Therefore, while the two types of evaluations (image-metric and human-based) consider a seemingly similar n-way identification task, *they are not directly comparable*. Additionally, note that they suggest different statistical generalization insights – generalization beyond the specific set of test examples, when using the image-metric (Wilcoxon test,  $N = 50, 120, p < 10^{-5}$ ), and generalization beyond the specific pool of human raters, in the human-based evaluations (Mann-Whitney test,  $N = 45, p < .001$ ). Overall, our method significantly outperforms state-of-the-art methods by a large margin in both image-metric-based and human-based evaluations.

#### 4.2. Decoding rich novel semantic categories of reconstructed images

The benefit of our self-supervised approach extends beyond the task of image reconstruction. It introduces significant gains in the task of semantic classification of the reconstructed images to their novel semantic categories (i.e., categories/classes never seen during training – classes neither represented in the “paired” training set, nor in the “unpaired” external images). For the classification task, we consider the “fMRI on ImageNet” dataset (Horikawa and Kamitani, 2017a), whose train-set contains 1200 images from 150 ImageNet classes. Its test-set contains 50 images from *disjoint* ImageNet classes. The unpaired external images used for the self-supervised training of our network, are drawn from 980 ImageNet classes. Notably, the unpaired images have no fMRI recordings whatsoever. The 50 test classes (to be recovered from the 50 test-fMRIs), are neither included in the 150 ‘paired’ train-classes, nor in the 980 classes of the ‘unpaired’ external images. *These are totally novel classes* (details in Sec 3.3). When checking the classification results of the 50 test fMRI, we test the classification of their reconstructed images against the gallery of 1030 ImageNet classes<sup>3</sup> As mentioned earlier (Sec

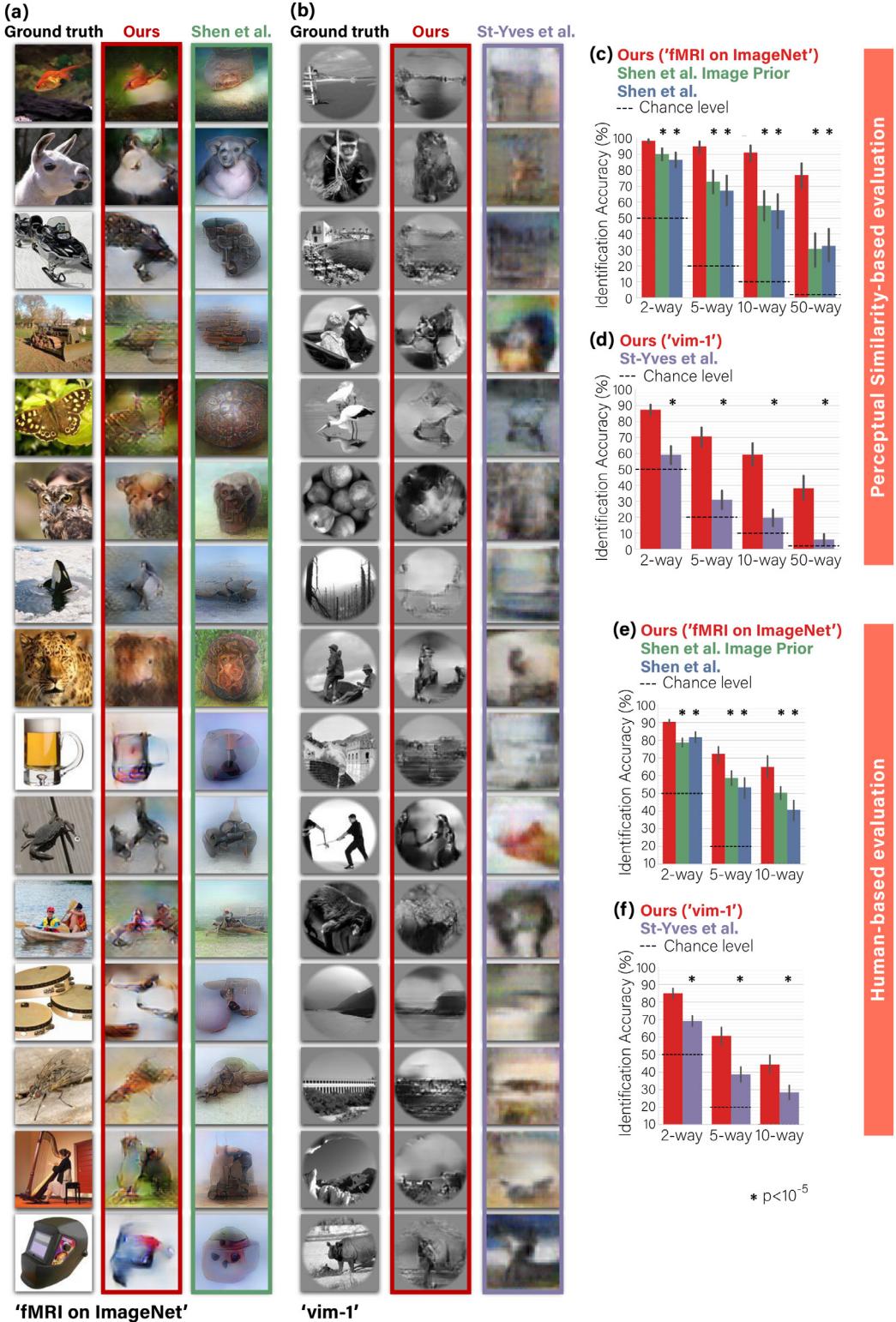
2.2), such classification requires no explicit class-labels in the training (see Fig 4). It is achieved by comparing the Deep-Features of the reconstructed image, with the Deep-Feature class-representative vector – one representative vector for each of the 1030 classes in our gallery (or any other gallery of novel classes).

Fig 5 a exemplifies novel-class classification results by our method. We present each reconstructed image alongside the five top predicted (“Top-5”) classes among the 1030 classes. For visualization purpose only, each of the Top-5 classes are visually exemplified by the nearest-neighbor image (most similar Deep-Features) among 100 randomly selected images from that class label. The first five rows show correct-classification cases at Top-5 accuracy level. In these cases our method successfully includes the ground truth class among the nearest five classes (marked by a red frame).

Interestingly, many of the non-ground-truth classes which are assigned among the Top-5, are also reasonable, frequently representing semantic and visual content, which is reminiscent of the ground truth class as well. For example, the ‘housefly’ is found similar also to other flies and comparable shape insects, and the ‘beer mug’ is also found similar to other types of mugs. The bottom rows show incorrect-classification cases, where the ground truth class is not found among the Top-5 classes. Nonetheless, even in these allegedly “failure” cases, many of the Top-5 classes (and even Top-1) are considerably relevant both semantically and visually. For example, the reconstructed ‘video’ image was associated with the ‘television’ class; the ‘duck’ was wrongly classified as an ‘ostrich’ ; the ‘leopard’ was wrongly classified as a ‘lion’; the ‘minaret’ was associated with several other tower-like classes, etc.

<sup>3</sup> ImageNet consists of 15K semantic classes, from which only 1000 classes participate in the ImageNet classification challenge (ILSVRC). In ‘fMRI on ImageNet’ (Horikawa and Kamitani, 2017a) which we use, only 20 out of the 50

test classes are included among the ILSVRC classes. The remaining 30 classes are taken from the larger collection of 15K ImageNet classes. Since our gallery is based on the 1000 ILSVRC classes, at train-time we omit the test classes, resulting in 980 train-classes (= 1000 – 20). At classification test-time, we add the 50 test-labels to the gallery, resulting in 1030 class labels (= 980 + 50).



**Fig. 8. Comparison of image-reconstruction with state-of-the-art methods.** (a), (b) Visual comparison with (Shen et al., 2019b; St-Yves and Naselaris, 2019) – each compared on its relevant dataset. Our method reconstructs shapes, details and global layout in images better than the leading methods. (c), (d) Quantitative comparisons of identification accuracy (per method) in n-way identification task according to Perceptual Similarity metric (see text for details). (e), (f) n-way identification responses of human raters via Mechanical Turk. Our self-supervised approach significantly outperforms all baseline methods on two datasets and across n-way difficulty levels by both types of experiments – image-metric-based and behavioral human-based (Wilcoxon,  $N = 50, 120$  for panels (c), (d); Mann-Whitney,  $N = 45$  for panels (e), (f)). 95% Confidence Intervals by bootstrap shown on charts.

We evaluate the performance of our classification results using two types of quantitative evaluations: (i) the “Classification Rank” – the average rank of the correct class among all classes (Table 2), and (ii) the more familiar “n-way classification” accuracy (evaluated for  $n = 50; 100; 500; 1000$ , Fig 5b).

These 2 types of numerical evaluation are detailed below.

(i) **Classification Rank.** As explained in Sec 2.2, the location of the ground-truth class within the list of 1030 classes (sorted according to our similarity-based classification measure - see Sec 2.2), provides a “rank score” for evaluating our classification accuracy per image (where rank=1 out of 1030 means perfect score). Moreover, this rank criterion further allows to assess the *generalization power* of our method beyond the limited test-set (the 50 test images and their 50 specific test classes available in the benchmark dataset). Table 2 summarizes novel-class classification rank results for all five subjects in ‘fMRI on ImageNet’. To demonstrate the power of our self-supervised approach, we compare its classification performance with a baseline of the purely supervised approach. This baseline uses reconstructions that were produced by a Decoder trained using the scarce paired data alone (Fig 1b). This comparison shows a leap improvement in median classification rank in favor of our self-supervised approach in all five subjects. Notably, for Subjects 2–5 (excluding Subjects 1 who has the lowest median noise-ceiling), the advantage of our self-supervised approach generalizes beyond the specifically chosen 50 test images and classes of the considered dataset.

(ii) **n-way Classification:** In addition to the Ranking-score within the 1030 gallery classes (for each test-fMRI), we present another alternative way of evaluating the classification results – through classification accuracy in n-way classification experiments (for  $n = 50, 100, 500, 1000$ ), using our automated Deep-Features class-similarity criterion. Fig 5b shows Top-1 classification accuracy across a range of classification task difficulties (shown for Subject 3; the remaining subjects can be found in the Supplementary-Material). The tasks differ in the number of candidate classes (n-way) from which prediction is made (i.e., the percent of cases that the Top-1 predicted class out of  $n$  class labels is indeed the ground-truth class). Our full method scores 29.9% Top-1 accuracy ( $SEM = 0.3\%$ ,  $N = 25000$ ) in 50-way classification task. Even when scaling to 1000-way (as in ImageNet classification), our method scores 12.1% Top-1 accuracy ( $SEM = 0.2\%$ ,  $N = 25000$ ), which exceeds chance level accuracy by more than 100-fold. Contrasting this performance with the baseline of the supervised approach shows a striking leap improvement in classification-accuracy in favor of our self-supervised approach: between x2 and x3 accuracy improvement, in all the n-way experiments ( $n = 50, 100, 500, 1000$ ).

We further performed an *ablation study* of the Perceptual Similarity for the task of semantic classification. Fig 5c shows 1000-way Top-5 classification accuracy by our self-supervised method, where the reconstructions used are produced by ablated versions of the Perceptual Similarity (Zhang et al., 2018b). Specifically, we limit the Perceptual Similarity criterion, which is used in Decoder training, to a varying range of Deep VGG layers, starting from using only the outputs of the first block (low “semantic” features) of VGG16, and up to aggregating outputs from all five blocks (high “semantic” features) of the pretrained network as in the full method. We find that the classification accuracy of the reconstructed images shows an increasing trend with the number of higher-level features, which are used as reconstruction criteria. This highlights the significance of the Perceptual Similarity reconstruction criterion, which includes higher-level features, for semantic classification. Note that the increasing trend appears to a various degree for different subjects, depending on experiment noise and subject-specific noise-ceiling (e.g., Subject 1 having the lowest noise-ceiling); Nevertheless, the trend is well illustrated by their cross-subject mean accuracy. Furthermore, we find that our classification approach, which is applied onto the reconstructed images, outperforms a direct-fMRI classification approach, where the test-fMRI samples are directly classified (see Supplementary-Material).

Our classification approach is inspired by Horikawa and Kamitani (2017a). Both methods use deep-feature embeddings to search for the nearest class centroid in a gallery of novel classes (our embedding is extracted from the reconstructed image (Fig 4b), whereas that of Horikawa and Kamitani (2017a) employs an intermediate deep image-embedding decoded from fMRI). Notably, (Horikawa and Kamitani, 2017a) presented the classification of novel categories in a 2-way task (i.e., discriminating between the correct category and a single random category). Here we scale up this classification task to 1000-way (i.e., finding the correct category among 1000 rich categories). *To the best of our knowledge, we are the first to demonstrate such large-scale semantic classification capabilities from fMRI data.*

#### 4.3. Predominance of early visual areas in reconstruction

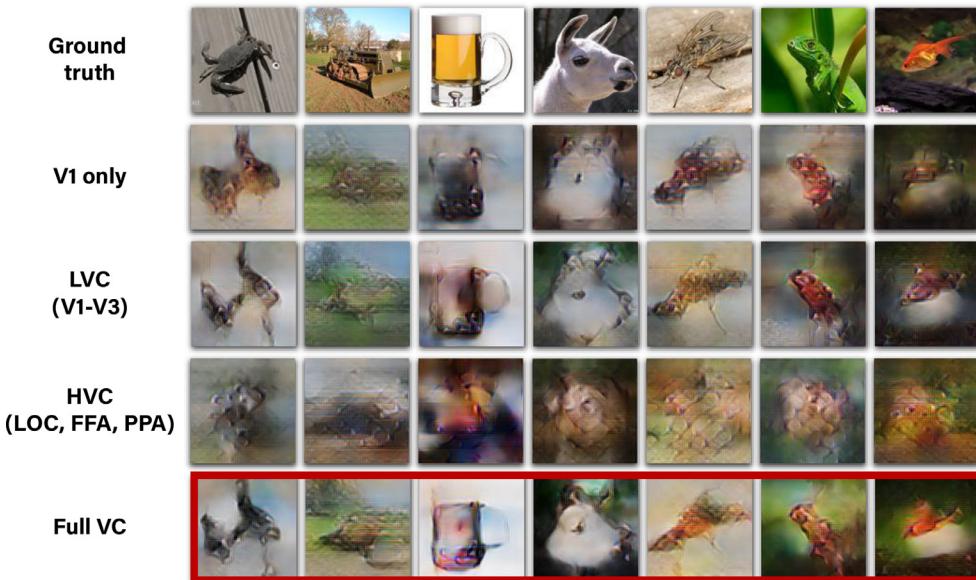
To reconstruct the images of ‘fMRI on ImageNet’, we considered 4600 visual-cortex voxels provided and labeled in Horikawa and Kamitani (2017a). We studied the contribution of different visual areas to our reconstruction performance. To this end, we selected subsets of voxels according to their marked brain areas, and restricted the training of our Encoder/Decoder to those voxels. Fig 9 shows reconstruction results when using voxels only from the following visual areas: (i) V1 ( 870 voxels), (ii) V1-V3, which refer to as Lower Visual Cortex (LVC, 2300 voxels), (iii) Fusiform Face Area (FFA), Parahippocampal Place Area (PPA), and Lateral Occipital Cortex (LOC), which we refer to as Higher Visual Cortex (HVC, 2150 voxels), or (iv) Full Visual Cortex (VC = LVC + V4 + HVC, 4600 voxels). These results show that the early visual areas, particularly V1-V3 (LVC), contain most of the information recoverable by our method. Considering voxels only from HVC leads to substantial degradation in performance despite comprising approximately half of the complete visual cortex voxels. Nevertheless, the higher visual areas clearly add semantic interpretability to the reconstructed images (which is evident when comparing the reconstructions from the Full VC, to those from LVC only).

Importantly, we found that removing any single visual area from our dataset, including V1, does not degrade the results significantly, suggesting the existence of information redundancy across visual areas. The results are strongly affected only when several regions, specifically the entire early visual cortex, are discarded. Furthermore adding V4 to either LVC or HVC did not change the results significantly.

We studied two extensions for the ROI-specific reconstruction results: (i) We controlled for the ROI-voxel count bias by sampling an equal number of voxels from each ROI. This control did not introduce major changes in the cross-ROI comparison or to the conclusion of early-visual-areas predominance; (ii) We analyzed the ROI-specific contribution also in semantic classification. Our results were consistent with those from the ROI-specific reconstruction experiments, showing predominance of early visual areas also in semantic classification. These results can be found in the Supplementary-Material.

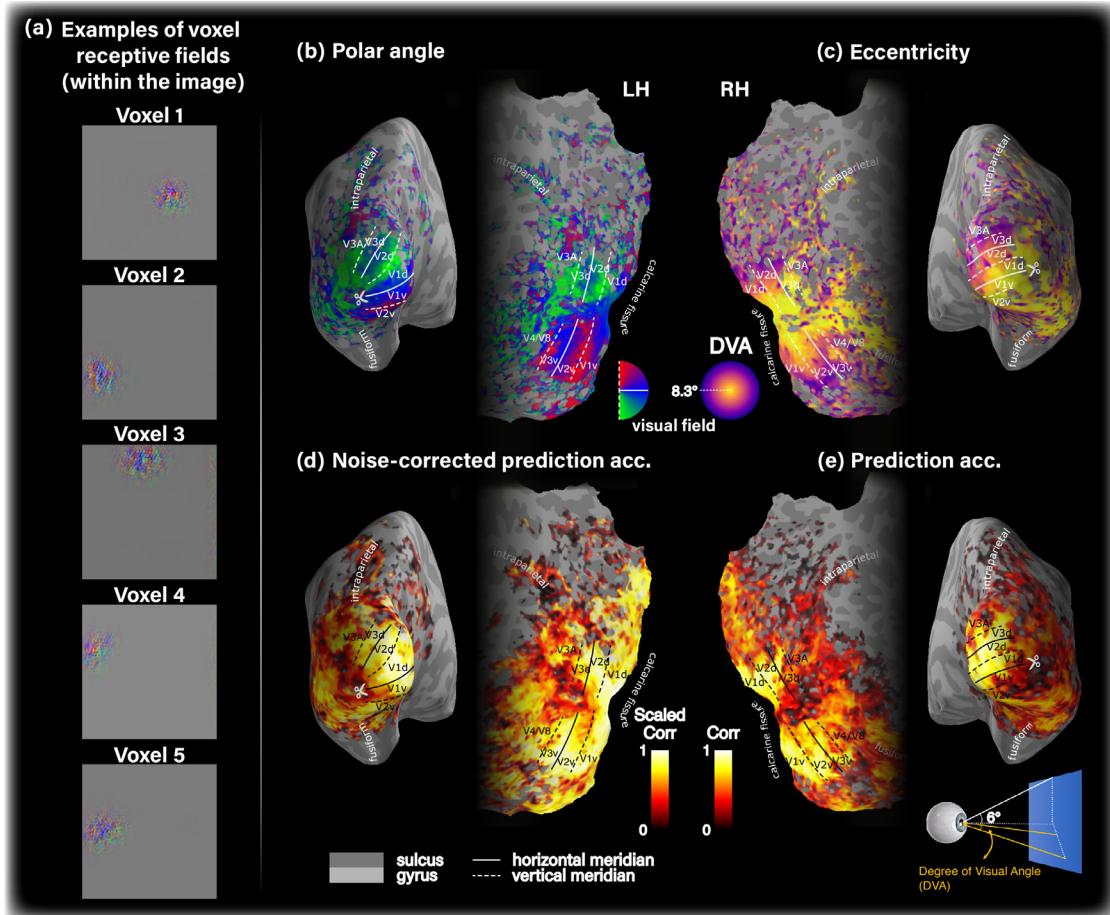
#### 4.4. Emergent biologically consistent population receptive fields

The human visual system is characterized by the well-known primate *retinotopic organization* (Bonhoeffer and Grinvald, 1991; Dow, 2002; Gomez et al., 2018; Mountcastle, 1957; Sereno et al., 1995; St-Yves and Naselaris, 2018; Tootell et al., 1982; Ts'o et al., 2009; Wandell and Winawer, 2015). Retinotopy maps reflect the spatial tuning of cortical hypercolumns or of their aggregation into Population Receptive Field (pRF) as in the case of voxel data. Here, we sought to visualize receptive field of voxels as captured by our Encoder & Decoder. To this end we considered two approaches: (i) Interpreting model weights that associate each voxel with spatial locations as a heatmap, and (ii) A gradient-based approach on the Encoder’s input space. While this approach is applicable to the Encoder only, it is favorable over the former approach because it provides a higher resolution receptive field map (i.e., 112x112, following Encoders *input dimensions* vs. 26x26 or 14x14, following the



**Fig. 9.** Decoding quality is dominated by early visual areas. Columns show reconstructions using our method with fMRI data from various ROIs in the visual cortex including:

- Primary Visual Cortex – V1
- Lower Visual Cortex – V1-V3
- Higher Visual Cortex – Fusiform Face Area (FFA), Parahippocampal Place Area (PPA), Lateral Occipital Cortex (LOC)
- Full Visual Cortex – LVC + V4 + HVC (in red frame).



**Fig. 10.** Our models capture biologically consistent voxel tuning properties. (a) Receptive field of five selected voxels with high SNR from early visual cortex, which indicates their spatial locality in the image. Panels (b)-(e) show single subject data on the corresponding subject-specific cortical surface. (b) Polar angle. (c) Eccentricity tuning, measured by degree of visual angle (DVA). (d) Noise-corrected prediction accuracy. (e) Prediction accuracy (non-scaled Pearson correlation). For simplicity we show the data on either left or right hemisphere. Voxel noise-ceiling is coded by transparency level (alpha channel) in all cortical maps.

layer's architecture in the Encoder/Decoder, see details in Sec 6). Since both methods gave rise to well aligned receptive fields (per voxel, see Supplementary-Material), we henceforth focused our analysis on the receptive fields recovered for the trained Encoder using the gradient-based approach.

[Fig 10a](#) shows receptive fields for several selected voxels, which indicates their spatial locality within the image. Next, we estimated the pRFs eccentricity and polar angle for each voxel. Lastly, we plot these data on the subject-specific cortical surface. [Fig 10b,c](#) shows the resulting tuning maps revealing the expected retinotopic organization. This includes the emergence of horizontal and vertical meridians and their transitions, contra-laterality and up-down inversion, and fovea-periphery gradual transition. We emphasize that this organization automatically emerged on its own from our training method involving natural images and fMRI. No biological atlas or other retinotopic prior was imposed. These results support the biological consistency of our model's predictions.

We sought to analyze the prediction accuracy achieved by our models. [Fig 10d,e](#) show the prediction accuracy distribution (Pearson correlation) of the modeled voxels when normalized by voxel noise-ceiling ([Fig 10d](#)) and when not normalized ([Fig 10e](#)). The prediction noise-ceiling is used to provide an estimate of the best possible prediction accuracy obtainable given infinite data ([Lage-Castellanos et al., 2019](#)). These panels show high prediction accuracy in LVC, and low in HVC. Furthermore, they show that throughout the visual cortex, our model markedly saturates the noise-ceiling of the given data. This indicates the sufficient expressive power of our multi-layer encoding model, enabling it to capture the given data complexity. A comparison of the model's prediction accuracy with a single-layer Encoder architecture from [Belyi et al. \(2019\)](#) can be found in the Supplementary-Material.

## 5. Discussion

We introduced a novel approach for self-supervised training of image decoders on external, image-only datasets. We show that our method gives rise to both state-of-the-art results in image-reconstruction as well as large scale semantic categorization from fMRI data. To date, the performance in the task of natural image reconstruction and semantic categorization from human fMRI recordings is limited by the characteristics of fMRI datasets. In the typical case, the paired training data are scarce, and represent a narrow semantic coverage.

Our self-supervised training on tens of thousands of additional unpaired images from wide coverage, combined with high-level Perceptual Similarity constrains, adapts the decoding model to the statistics of natural images and novel categories. Our framework enables substantial improvement in both image reconstruction quality and classification capabilities compared to methods that rely only on the scarce paired training data. This self-supervised perceptual approach leads to state-of-the-art image reconstructions from fMRI, of unprecedented quality, as supported by image-metric-based evaluations, as well as extensive human-based evaluations. We accomplish this for two substantially different fMRI datasets using a single method (with the same hyperparameters).

Our self-supervised training on tens of thousands of unpaired external images further leads to unprecedented capabilities in the semantic classification of fMRI data (and moreover, of classes never encountered during training). We consider the challenging 1000-way semantic classification task, and demonstrate a striking leap improvement (more than 2x) in classification performance when applying our self-supervised approach over a purely supervised approach. To the best of our knowledge, we are the first to demonstrate such large-scale semantic classification capabilities (1000-way) from fMRI data.

The improved reconstruction quality achieved when minimizing deep image representation error rather than low-level pixel-based one ([Fig 2](#)) supports reports in earlier works that employed similar approaches ([Shen et al., 2019a; 2019b](#)). In [Shen et al. \(2019b\)](#), Perceptual Similarity was used to drive image reconstruction by iterative op-

timization at test time. Here, however, the Decoder itself is trained to minimize the higher-level reconstruction error, allowing a non-iterative feed-forward reconstruction at test-time. In this respect, our approach is more compatible with [Shen et al. \(2019a\)](#), albeit our reconstruction loss is based on the complete Perceptual Similarity metric ([Zhang et al., 2018b](#)), defined via features extracted at multiple intermediate layers of a pretrained object classification network rather than only its almost last layer. Furthermore, this work goes beyond ([Shen et al., 2019a; 2019b](#)) in two aspects: (i) Our Decoder benefited from applying Perceptual Similarity to substantially more data with no fMRI recording. The combined effect of this type of reconstruction criterion with the additional unpaired data is very effective, as evident in [Fig 2](#); (ii) We show that applying Perceptual Similarity not only results in an improved reconstruction quality, but also gives rise to a higher classification accuracy of novel ImageNet classes (classes that were not represented in any of the training data). An important question for future research is whether a natural-image prior, such as the Deep Generator Network employed by [Shen et al. \(2019b\)](#) can be combined with the self-supervised approach in a way that benefits from both the self-supervision on additional unpaired data and the imposed narrower Natural Image search space.

Our ablation studies indicate that reconstruction quality is dominated by data originating from Lower Visual Cortex (V1-V3). The extended architecture of the Encoder, which incorporates high-level deep features, was designed to improve information-harnessing from the Higher Visual Cortex (HVC) as well. Indeed prediction accuracy maps show that the noise-ceiling is saturated throughout the visual cortex, including in higher visual areas. These findings suggest a reasonable representation of the HVC by our model. Nevertheless, the SNR of the data arising from these areas renders them weaker contributors to overall reconstruction quality.

We provide evidence for the retinotopic organization implicitly learned (on its own) by our image-to-fMRI Encoder. This suggests that our models are biologically meaningful, as opposed to tailored and over-fit to a limited dataset. Note that while we show data for the Encoder, we verified in our experiments that model voxels in the Decoder and the Encoder indeed agree (while not explicitly forced to do so).

The proposed method currently focuses on data from individual subjects. A natural extension of the present work is to combine information across multiple subjects (see [Wen et al. \(2018c\)](#)). This is part of our future work.

## 6. Additional technical details

**Hyperparameters.** We trained the Encoder using Adam optimizer for 50 epochs with an initial learning rate of 1e-3, with a 90% learning rate drop using milestones (20, 30, and 35 epochs). During Decoder training with supervised and unsupervised objectives, each training batch contained 16 pairs (supervised training), and 16 unpaired natural images (randomly sampled from the external image database – images without fMRI). We trained the Decoder for 150 epochs using Adam optimizer with an initial learning rate of 1e-3, and 80% learning rate drop after every 30 epochs.

**Runtime.** Our system completes the two-stage training within approximately 1.5 hours using a single Tesla V100 GPU. Once trained, the inference itself (decoding of a new fMRI) is performed in a few milliseconds per image.

**Behavioral experiments.** The participants in the Mechanical Turk behavioral experiments gave their online informed consent to be recorded, and were granted financial incentives for every completed survey. The research protocol was reviewed and approved by the Bioethics and Embryonic Stem Cell Research Oversight (ESCRO) Committee at the Weizmann Institute of Science. In order to assure the validity of the behavioral data (e.g. bot observers, fatigue along the survey), we screened subjects according to their score in interleaved sanity check experiments. The sanity check experiments comprised adding to the actual

experiments also 10% unexpected trivial identification tasks of mildly degraded versions of the ground truth images, instead of the reconstructed images. We further discarded subjects with MTurk success-score (reputation) lower than 97%. Each survey consisted of 50 or 20 trials corresponding to the number of test-images comparison in ‘fMRI on ImageNet’ (Horikawa and Kamitani, 2017a) or ‘vim-1’ (Kay et al., 2008)<sup>4</sup>, all of which were reconstructed using a single particular method. In each trial subjects were presented with a reconstructed image and  $n$  candidate images, the ground-truth image and  $n - 1$  distractor images, and were prompted “Which image at the bottom row is most similar to the image at the top row?”. To assure task difficulty agreement across subjects and reconstruction methods the set of distractor images was randomly selected for each test-image, but remained fixed across surveys; Our results were insensitive to their re-selection.

**Semantic category decoding.** We defined the feature vector underlying the class representatives to be the outputs of block 4 in AlexNet, and used Pearson correlation as the distance metric for ranking class representatives. Our experiments showed that using this intermediate representation level as the embedding of choice yields optimal results for classification.

**Voxel receptive field visualization and estimation.** We analyzed the receptive field of our encoding model in terms of simulated retinotopy as previously proposed in Eickenberg et al. (2017); Wen et al. (2018a). To generate analogous spatial tuning maps for our modeled voxels, we estimated the voxel spatial tuning captured by our models. We start by visualizing each voxel’s receptive field (pRF) using the trained Encoder. Our primary approach to this end was gradient-based (Grossman et al., 2019; Simonyan et al., 2000): Given a random input image, we compute the gradient of a particular voxel with respect to this input. This allows to visualize the image which would drive the maximum change in activity at the target voxel. To produce a heat-map, the values within the resulting gradient-image are squared, averaged across the color-channels, and normalized. Next, we define the pRF center as the center of mass of the preprocessed map. The preprocessing was designed to minimize noise effects. It included map smoothing with a Gaussian kernel,  $\sigma = 3$ , followed by raising the map values to the power of 10 (emphasizing the areas with higher values). About 15% of the voxels had pRF maps which were not confined spatially around a center of mass, and were thus discarded in subsequent analysis. The remaining 85% voxels were considered in the retinotopy maps. In a secondary approach, we recovered voxel receptive field maps for both the Encoder and the Decoder from the learned voxel weights. For the Encoder, we directly visualized the spatial map weights of the space-feature locally-connected layer. For the Decoder, we reshaped the input fully-connected layer to square spatial dimensions and 64 channels and computed the mean square weights across channels to yield a heat map.

**Cortical surface visualization.** We reconstructed and flattened the individual cortical surfaces using FreeSurfer, and plotted the data using a custom extension of PySurfer.

**Noise-Ceiling.** We estimated the fMRI prediction Noise-Ceiling by half-split over the test data repeats following Lage-Castellanos et al. (2019).

**Statistics.** We used Wilcoxon signed-rank (paired) test (two-tailed) for significance testing in the image-metric-based multi-image identification experiments. For the rank-classification experiments, and the (unpaired) behavioral experiments we used Mann-Whitney rank test.

## Data and code availability statements

We used publicly available fMRI datasets as well as image datasets as follows: fMRI on ImageNet, vim-1, and ImageNet. We will make our code publicly available upon publication.

<sup>4</sup> ‘vim-1’ originally contains 120 test-images, however in the behavioral evaluation we considered only the subset of 20 images that were defined in St-Yves and Naselaris (2019) as test-images.

## Credit authorship contribution statement

**Guy Gaziv:** Conceptualization, Writing – original draft, Methodology, Investigation, Formal analysis, Visualization, Writing – review & editing. **Roman Belyi:** Conceptualization, Methodology, Investigation, Software, Writing – review & editing. **Niv Granot:** Conceptualization, Methodology, Investigation, Software, Writing – review & editing. **Assaf Hoogi:** Investigation, Resources, Writing – review & editing. **Francesca Strappini:** Resources, Investigation, Writing – review & editing. **Tal Golan:** Resources, Visualization, Formal analysis, Writing – review & editing. **Michal Irani:** Conceptualization, Methodology, Investigation, Supervision, Funding acquisition, Writing – review & editing.

## Acknowledgments

This project has received funding from the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation programme (grant agreement No 788535).

## References

- Belyi, R., Gaziv, G., Hoogi, A., Strappini, F., Golan, T., Irani, M., 2019. From voxels to pixels and back: Self-supervision in natural-image reconstruction from fMRI. Advances in Neural Information Processing Systems. <https://papers.nips.cc/paper/8879-from-voxels-to-pixels-and-back-self-supervision-in-natural-image-reconstruction-from-fmri.pdf> <http://www.wisdom.weizmann.ac.il/~vision/ssfmri2im/>
- Bonhoeffer, T., Grinvald, A., 1991. Iso-orientation domains in cat visual cortex are arranged in pinwheel-like patterns. Nature 353 (6343), 429–431. doi:10.1038/353429a0. <http://www.nature.com/articles/353429a0>
- Cichy, R.M., Heinze, J., Haynes, J.D., 2012. Imagery and perception share cortical representations of content and location. Cereb. Cortex 22 (2), 372–380. doi:10.1093/cercor/cbr106. <http://www.ncbi.nlm.nih.gov/pubmed/21666128>
- Cox, D.D., Savoy, R.L., 2003. Functional magnetic resonance imaging (fMRI) “brain reading”: detecting and classifying distributed patterns of fMRI activity in human visual cortex. Neuroimage 19 (2 Pt 1), 261–270. <http://www.ncbi.nlm.nih.gov/pubmed/12814577>
- Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., Fei-Fei, L., 2009. ImageNet: A large-scale hierarchical image database. In: 2009 IEEE Conference on Computer Vision and Pattern Recognition. IEEE, pp. 248–255. doi:10.1109/CVPR.2009.5206848. <http://ieeexplore.ieee.org/document/5206848/>
- Dow, B.M., 2002. Orientation and color columns in Monkey visual cortex. Cereb. Cortex 12 (10), 1005–1015. doi:10.1093/cercor/12.10.1005. <https://academicoup.com/cercor/article-lookup/doi/10.1093/cercor/12.10.1005>
- Eickenberg, M., Gramfort, A., Varoquaux, G., Thirion, B., 2017. Seeing it all: convolutional network layers map the function of the human visual system. Neuroimage 152, 184–194. doi:10.1016/j.neuroimage.2016.10.001. <http://www.sciencedirect.com/science/article/pii/S1053811916305481>
- Glasser, M.F., Sotiropoulos, S.N., Wilson, J.A., Coalson, T.S., Fischl, B., Andersson, J.L., Xu, J., Jbabdi, S., Webster, M., Polimeni, J.R., Van Essen, D.C., Jenkinson, M., Consortium, W.-M.H., 2013. The minimal preprocessing pipelines for the human connectome project. Neuroimage 80, 105–124. doi:10.1016/j.neuroimage.2013.04.127. <http://www.ncbi.nlm.nih.gov/pubmed/23668970> <http://www.ncbi.nlm.nih.gov/articlerender.fcgi?artid=PMC3720813> <http://linkinghub.elsevier.com/retrieve/pii/S1053811913005053>
- Glorot, X., Bengio, Y., 2010. Understanding the difficulty of training deep feedforward neural networks. <http://proceedings.mlr.press/v9/glorot10a.html>
- Glover, G. H., 2011. Overview of functional magnetic resonance imaging. /pmc/articles/PMC3073717/ /pmc/articles/PMC3073717/?report=abstract <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3073717/> 10.1016/j.jneu.2010.11.001
- Gomez, J., Natu, V., Jeska, B., Barnett, M., Grill-Spector, K., 2018. Development differentially sculpts receptive fields across early and high-level human visual cortex. Nat. Commun. 9 (1), 788. doi:10.1038/s41467-018-03166-3.
- Grossman, S., Gaziv, G., Yeagle, E.M., Harel, M., Mégevand, P., Groppe, D.M., Khurvis, S., Herrero, J.L., Irani, M., Mehta, A.D., Malach, R., 2019. Convergent evolution of face spaces across human face-selective neuronal groups and deep convolutional networks. Nat. Commun. 10 (1), 1–13. doi:10.1038/s41467-019-12623-6.
- Guclu, U., van Gerven, M.A.J., Güçlü, U., van Gerven, M.A.J., 2015. Deep neural networks reveal a gradient in the complexity of neural representations across the ventral stream. J. Neurosci. 35 (27), 10005–10014. doi:10.1523/JNEUROSCI.5023-14.2015.
- Han, K., Wen, H., Shi, J., Lu, K.-H., Zhang, Y., Fu, D., Liu, Z., 2019. Variational autoencoder: an unsupervised model for encoding and decoding fMRI activity in visual cortex. Neuroimage 198, 125–136. doi:10.1016/J.NEUROIMAGE.2019.05.039.
- Haxby, J.V., Gobbini, M.I., Furey, M.L., Ishai, A., Schouten, J.L., Pietrini, P., 2001. Distributed and overlapping representations of faces and objects in ventral temporal cortex. Science 293 (5539), 2425–2430. doi:10.1126/science.1063736.
- Haynes, J.-D., Rees, G., 2006. Decoding mental states from brain activity in humans. Nat. Rev. Neurosci. 7 (7), 523–534. doi:10.1038/nrn1931.
- Hebart, M.N., Baker, C.I., 2018. Deconstructing multivariate decoding for the study of brain function. Neuroimage 180, 4–18. doi:10.1016/J.NEUROIMAGE.2017.08.005.
- Horikawa, T., Kamitani, Y., 2017. Generic decoding of seen and imagined objects using hierarchical visual features. Nat. Commun. 8 (1), 1–15. doi:10.1038/ncomms15037.

- Horikawa, T., Kamitani, Y., 2017. Hierarchical neural representation of dreamed objects revealed by brain decoding with deep neural network features. *Front. Comput. Neurosci.* 11. doi:10.3389/fncom.2017.00004. <http://journal.frontiersin.org/article/10.3389/fncom.2017.00004/full>.
- Horikawa, T., Tamaki, M., Miyawaki, Y., Kamitani, Y., 2013. Neural decoding of visual imagery during sleep. *Science* 340 (6132). doi:10.1126/science.1234330. <http://www.sciencemag.org/content/340/6132/639.full.html> <http://www.sciencemag.org/content/suppl/2013/04/03/science.1234330.DC1.html> <http://www.sciencemag.org/content/suppl/2013/04/04/science.1234330.DC2.html> <http://www.sciencemag.org/content/340/6132/639>.
- Kamitani, Y., Tong, F., 2005. Decoding the visual and subjective contents of the human brain. *Nat. Neurosci.* 8 (5), 679–685. doi:10.1038/nn1444. <http://www.ncbi.nlm.nih.gov/pubmed/15852014> <http://www.ncbi.nlm.nih.gov/articlerender.fcgi?artid=PMC1808230> <http://www.nature.com/doifinder/10.1038/nn1444>.
- Kamitani, Y., Tong, F., 2006. Decoding seen and attended motion directions from activity in the human visual cortex. *Curr. Biol.* 16 (11), 1096–1102. doi:10.1016/j.cub.2006.04.003. <https://pubmed.ncbi.nlm.nih.gov/16753563/>.
- Kay, K.N., Naselaris, T., Prenger, R.J., Gallant, J.L., 2008. Identifying natural images from human brain activity. *Nature* 452 (7185), 352–355. doi:10.1038/nature06713. <http://www.nature.com/doifinder/10.1038/nature06713> [http://gallantlab.org\\_downloads/2008a.Kay.pdf](http://gallantlab.org_downloads/2008a.Kay.pdf)
- Konkle, T., Caramazza, A., 2013. Tripartite organization of the ventral stream by animacy and object size. *J. Neurosci.* 33 (25), 10235–10242. doi:10.1523/JNEUROSCI.0983-13.2013.
- Lage-Castellanos, A., Valente, G., Formisano, E., De Martino, F., 2019. Methods for computing the maximum performance of computational models of fMRI responses. *PLoS Comput. Biol.* 15 (3), e1006397. doi:10.1371/journal.pcbi.1006397.
- Lin, Y., Li, J., Wang, H., Jiao, S., 2019. DCNN-GAN: Reconstructing Realistic Image from fMRI. Technical Report. <https://arxiv.org/pdf/1901.07368.pdf>
- Milekovic, T., Sarma, A.A., Bacher, D., Simeral, J.D., Saab, J., Pandarinath, C., Sorice, B.L., Blabe, C., Oakley, E.M., Tringale, K.R., Eskandar, E., Cash, S.S., Henderson, J.M., Shenoy, K.V., Donoghue, J.P., Hochberg, L.R., 2018. Stable long-term BCI-enabled communication in ALS and locked-in syndrome using LFP signals. *J. Neurophysiol.* 120 (1), 343–360. doi:10.1152/JN.00493.2017. <https://pubmed.ncbi.nlm.nih.gov/29694279/>.
- Miyawaki, Y., Uchida, H., Yamashita, O., Sato, M.-a., Morito, Y., Tanabe, H.C., Sadato, N., Kamitani, Y., 2008. Visual image reconstruction from human brain activity using a combination of multiscale local image decoders. *Neuron* 60 (5), 915–929. doi:10.1016/j.neuron.2008.11.004. <http://linkinghub.elsevier.com/retrieve/pii/S0896627308009586> <http://www.sciencedirect.com/science/article/pii/S0896627308009586> <http://www.ncbi.nlm.nih.gov/pubmed/19081384>.
- Monti, M.M., Vanhaudenhuyse, A., Coleman, M.R., Boly, M., Pickard, J.D., Tshibanda, L., Owen, A.M., Laureys, S., 2010. Willful modulation of brain activity in disorders of consciousness. *N. Engl. J. Med.* 362 (7), 579–589. doi:10.1056/NEJMoa0905370/SUPPL\_FILE/NEJM\_MONTI\_597SA1.PDF. <https://www.nejm.org/doi/full/10.1056/NEJMoa0905370>.
- Mountcastle, V.B., 1957. Modality and topographic properties of single neurons of cat's somatic sensory cortex. *J. Neurophysiol.* 20 (4), 408–434. doi:10.1152/jn.1957.20.4.408. <http://www.ncbi.nlm.nih.gov/pubmed/13439410>.
- Mozafari, M., Reddy, L., Vanrullen, R., Reconstructing Natural Scenes from fMRI Patterns using BigBiGAN. Technical Report. <https://openneuro.org/datasets/ds001246/>.
- Naci, L., Monti, M.M., Cruse, D., Kübler, A., Sorger, B., Goebel, R., Kotchoubey, B., Owen, A.M., 2012. Braincomputer interfaces for communication with non-responsive patients. *Ann. Neurol.* 72 (3), 312–323. doi:10.1002/ANA.23656. <https://onlinelibrary.wiley.com/doi/full/10.1002/ana.23656> <https://onlinelibrary.wiley.com/doi/10.1002/ana.23656>.
- Naselaris, T., Kay, K.N., Nishimoto, S., Gallant, J.L., 2011. Encoding and decoding in fMRI. *Neuroimage* 56 (2), 400–410. doi:10.1016/j.neuroimage.2010.07.073. <http://www.ncbi.nlm.nih.gov/pubmed/20691790> <http://www.ncbi.nlm.nih.gov/articlerender.fcgi?artid=PMC3037423>.
- Naselaris, T., Olman, C.A., Stanbury, D.E., Ugurbil, K., Gallant, J.L., 2015. A voxel-wise encoding model for early visual areas decodes mental images of remembered scenes. *Neuroimage* 105, 215–228. doi:10.1016/j.neuroimage.2014.10.018. <http://www.ncbi.nlm.nih.gov/pubmed/25451480> <http://www.ncbi.nlm.nih.gov/articlerender.fcgi?artid=PMC4364759> <http://linkinghub.elsevier.com/retrieve/pii/S1053811914008428>.
- Naselaris, T., Prenger, R.J., Kay, K.N., Oliver, M., Gallant, J.L., 2009. Bayesian reconstruction of natural images from human brain activity. *Neuron* 63 (6), 902–915. doi:10.1016/j.neuron.2009.09.006. <http://linkinghub.elsevier.com/retrieve/pii/S0896627309006850> [http://www.cell.com/neuron/pdf/S0896-6273\(09\)00685-0.pdf](http://www.cell.com/neuron/pdf/S0896-6273(09)00685-0.pdf).
- Nishimoto, S., Vu, A.T., Naselaris, T., Benjamini, Y., Yu, B., Gallant, J.L., 2011. Reconstructing visual experiences from brain activity evoked by natural movies. *Curr. Biol.* 21 (19), 1641–1646. doi:10.1016/j.cub.2011.08.031. <http://www.ncbi.nlm.nih.gov/pubmed/21945275> <http://www.ncbi.nlm.nih.gov/articlerender.fcgi?artid=PMC3326357> [http://www.cell.com/current-biology/pdf/S0960-9822\(11\)00937-7.pdf](http://www.cell.com/current-biology/pdf/S0960-9822(11)00937-7.pdf).
- O'Craven, K.M., Kanwisher, N., 2000. Mental imagery of faces and places activates corresponding stimulus-specific brain regions. *J. Cogn. Neurosci.* 12 (6), 1013–1023. doi:10.1162/08989290051137549. <https://pubmed.ncbi.nlm.nih.gov/11177421/>.
- Owen, A.M., Coleman, M.R., Boly, M., Davis, M.H., Laureys, S., Pickard, J.D., 2006. Detecting awareness in the vegetative state. *Science* 313 (5792), 1402.
- Pearson, J., Naselaris, T., Holmes, E.A., Kosslyn, S.M., 2015. Mental imagery: functional mechanisms and clinical applications. *Trends Cogn. Sci.* 19 (10), 590–602. doi:10.1016/j.tics.2015.08.003. <http://www.ncbi.nlm.nih.gov/pubmed/26412097> <http://www.ncbi.nlm.nih.gov/articlerender.fcgi?artid=PMC4595480> <http://linkinghub.elsevier.com/retrieve/pii/S1364661315001801>.
- Qiao, K., Chen, J., Wang, L., Zhang, C., Tong, L., Yan, B., 2020. BigGAN-based bayesian reconstruction of natural images from human brain activity. *Neuroscience* 444, 92–105. doi:10.1016/j.neuroscience.2020.07.040.
- Qiao, K., Chen, J., Wang, L., Zhang, C., Zeng, L., Tong, L., Yan, B., 2019. Category decoding of visual stimuli from human brain activity using a bidirectional recurrent neural network to simulate bidirectional information flows in human visual cortices. *Front. Neurosci.* 13 (JUL). doi:10.3389/fnins.2019.00692.
- Ren, Z., Li, J., Xue, X., Li, X., Yang, F., Jiao, Z., Gao, X., 2021. Reconstructing seen image from brain activity by visually-guided cognitive representation and adversarial learning. *Neuroimage* 228, 117602. doi:10.1016/j.neuroimage.2020.117602.
- Seeliger, K., Güçlü, U., Ambrogioni, L., Güçlütürk, Y., van Gerven, M.A.J., 2018. Generative adversarial networks for reconstructing natural images from brain activity. *Neuroimage* 181, 775–785. doi:10.1016/j.neuroimage.2018.07.043. <https://www.biorxiv.org/content/biorxiv/early/2017/12/08/226688.full.pdf>.
- Sereno, M.I., Dale, A.M., Reppas, J.B., Kwong, K.K., Belliveau, J.W., Brady, T.J., Rosen, B.R., Tootell, R.B.H., 1995. Borders of multiple visual areas in humans revealed by functional magnetic resonance imaging. *Science* 268 (5212), 889–893. doi:10.1126/science.7754376. <https://pubmed.ncbi.nlm.nih.gov/7754376/>
- Shen, G., Dwivedi, K., Majima, K., Horikawa, T., Kamitani, Y., 2019. End-to-end deep image reconstruction from human brain activity. *Front. Comput. Neurosci.* 13, 21. doi:10.3389/fncom.2019.00021. [www.frontiersin.org](http://www.frontiersin.org).
- Shen, G., Horikawa, T., Majima, K., Kamitani, Y., 2019. Deep image reconstruction from human brain activity. *PLoS Comput. Biol.* 15 (1), e1006633. doi:10.1371/journal.pcbi.1006633.
- Simonyan, K., Vedaldi, A., Zisserman, A., Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps. Technical Report. <http://code.google.com/p/cuda-convnet/>.
- Simonyan, K., Zisserman, A., 2014. Very Deep Convolutional Networks for Large-Scale Image Recognition <http://arxiv.org/abs/1409.1556>.
- St-Yves, G., Naselaris, T., 2017. The feature-weighted receptive field: an interpretable encoding model for complex feature spaces. [10.1016/j.neuroimage.2017.06.035](https://doi.org/10.1016/j.neuroimage.2017.06.035)
- St-Yves, G., Naselaris, T., 2018. The feature-weighted receptive field: an interpretable encoding model for complex feature spaces. *Neuroimage* 180 (Pt A), 188–202. doi:10.1016/j.neuroimage.2017.06.035. <http://www.ncbi.nlm.nih.gov/pubmed/28645845> <http://www.ncbi.nlm.nih.gov/articlerender.fcgi?artid=PMC5886832> <http://linkinghub.elsevier.com/retrieve/pii/S1053811917305086>.
- St-Yves, G., Naselaris, T., 2019. Generative Adversarial Networks Conditioned on Brain Activity Reconstruct Seen Images. In: Proceedings - 2018 IEEE International Conference on Systems, Man, and Cybernetics, SMC 2018. Institute of Electrical and Electronics Engineers Inc., pp. 1054–1061. doi:10.1109/SMC.2018.00187. <https://ieeexplore.ieee.org/document/8616183/>.
- Thirion, B., Duchesnay, E., Hubbard, E., Dubois, J., Poline, J.-B., Lebihan, D., Dehaene, S., 2006. Inverse retinotopy: inferring the visual content of images from brain activation patterns. *Neuroimage* 33 (4), 1104–1116. doi:10.1016/j.neuroimage.2006.06.062. <http://linkinghub.elsevier.com/retrieve/pii/S1053811906007373> <http://www.sciencedirect.com/science/article/pii/S1053811906007373>.
- Tootell, R.B.H., Silverman, M.S., Switkes, E., De Valois, R.L., 1982. Deoxyglucose analysis of retinotopic organization in primate striate cortex. *Science* 218 (4575), 900–901. doi:10.1126/science.7134981. <http://science.sciencemag.org/>
- Tripathy, K., Markow, Z.E., Fishell, A.K., Sherafati, A., Burns-Yocom, T.M., Schroeder, M.L., Svoboda, A.M., Eggebrecht, A.T., Anastasio, M.A., Schlaggar, B.L., Culver, J.P., 2021. Decoding visual information from high-density diffuse optical tomography neuroimaging data. *Neuroimage* 226, 117516. doi:10.1016/J.NEUROIMAGE.2020.117516.
- Ts'o, D.Y., Zarella, M., Burkitt, G., 2009. Whither the hypercolumn? *J. Physiol.* 587 (Pt 12), 2791–2805. doi:10.1113/jphysiol.2009.171082. <http://www.ncbi.nlm.nih.gov/pubmed/19525564> <http://www.ncbi.nlm.nih.gov/articlerender.fcgi?artid=PMC2718239>.
- Wandell, B.A., Winawer, J., 2015. Computational neuroimaging and population receptive fields. *Trends Cogn. Sci.* 19 (6), 349–357. doi:10.1016/j.tics.2015.03.009. <http://www.ncbi.nlm.nih.gov/pubmed/25850730> <http://www.ncbi.nlm.nih.gov/articlerender.fcgi?artid=PMC4484758> <http://linkinghub.elsevier.com/retrieve/pii/S1364661315000704>.
- Wen, H., Han, K., Shi, J., Zhang, Y., Culurciello, E., Liu, Z., 2018a. Deep Predictive Coding Network for Object Recognition <https://arxiv.org/pdf/1802.04762.pdf>.
- Wen, H., Shi, J., Chen, W., Liu, Z., 2018. Deep residual network predicts cortical representation and organization of visual features for rapid categorization. *Sci. Rep.* 8 (1), 3752. doi:10.1038/s41598-018-22160-9. <http://www.ncbi.nlm.nih.gov/pubmed/29491405> <http://www.ncbi.nlm.nih.gov/articlerender.fcgi?artid=PMC5830584> <http://www.nature.com/articles/s41598-018-22160-9>.
- Wen, H., Shi, J., Chen, W., Liu, Z., 2018. Transferring and generalizing deep-learning-based neural encoding models across subjects. *Neuroimage* 176, 152–163. doi:10.1016/j.neuroimage.2018.04.053.
- Wen, H., Shi, J., Zhang, Y., Lu, K.-H., Cao, J., Liu, Z., 2018. Neural encoding and decoding with deep learning for dynamic natural vision. *Cereb. Cortex* 28 (12), 4136–4160. doi:10.1093/cercor/bhw268. <https://academic.oup.com/cercor/article/28/12/4136/4560155>.

Zhang, C., Qiao, K., Wang, L., Tong, L., Zeng, Y., Yan, B., 2018. Constraint-Free natural image reconstruction from fMRI signals based on convolutional neural network. *Front. Hum. Neurosci.* 12, 242. doi:[10.3389/fnhum.2018.00242](https://doi.org/10.3389/fnhum.2018.00242). <https://www.frontiersin.org/article/10.3389/fnhum.2018.00242/full>.

Zhang, R., Isola, P., Efros, A.A., Shechtman, E., Wang, O., 2018. The Unreasonable Effectiveness of Deep Features as a Perceptual Metric. In: IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). <https://www.github.com/richzhang/PerceptualSimilarity>