

1 We thank all reviewers for the valuable and constructive comments.

2 **To Reviewer #1**

3 **What stimulus set was being used to train the GAN, and for augmentation?** The fMRI dataset includes 1200
4 training images (150 categories) and 50 test images with the corresponding fMRI recordings. The GAN was trained
5 with the 1200 training images in the fMRI dataset. In augmentation, we sampled extra 200 images from the ImageNet
6 dataset. The images used for augmentation have no overlap with both training and test images from the fMRI dataset.

7 **Is Figure 2(a) displaying results for test stimuli or training stimuli?** Figure 2(a) displays the results for test stimuli.
8 We will make it clear in the final version. Thanks.

9 **How did the hyperparameter values being selected?** In the training of GAN, 100 fMRI recordings of the training set
10 were kept as a validation set for hyperparameter selection. Then we used the whole training set (containing the samples
11 in validation set) to train the model with the selected hyperparameters. We will make it clear in the final version.

12 **Why use a DNN for the semantic features instead of a linear regression, as for the visual ones? Why Tanh
13 instead of ReLU?** We have evaluated both DNN and linear regression for semantic decoding and did not put it in
14 the paper for space limit. In this experiment, the DNN model obtained better performance (53.9% v.s. 49.1% in
15 classification accuracy). About the activation function, using Tanh function can converge stably in this experiment. We
16 agree that functions like ReLU can also work in this model. We will explore different models for semantic decoding in
17 the future study.

18 **About the evaluation and comparison criterion.** Thanks for your constructive suggestions. We did not conduct other
19 comparison methods because the codes of some papers were not available during the experiment. We are trying to
20 supplement some results of the Friedman test in the final version as suggested.

21 For the other suggestions/issues, we will revise the paper accordingly. Many thanks for your valuable comments.

22 **To Reviewer #2**

23 **About the contributions.** Thanks for confirming our contributions. The previous study [17] provides a good basis and
24 support for our approach. We will clarify the contributions in the introduction, and add extra discussions in the final
25 version. In addition, we will add necessary details to the unclear parts in the article.

26 **To Reviewer #4**

27 **Why using the U-Net structure in the generator?** The advantage of using U-Net lies in that it both processes
28 high-level information (passed through in the bottleneck) and preserves low-level structures (by skip connections). For
29 image reconstruction, we need to combine both high-level (category) and low-level (shape) features, thus U-Net is a
30 suitable choice.

31 **Regarding the necessity of the HVC signals and loss of information.** We agree that theoretically LVC information
32 should contain the information from HVC. However, in the process of recording fMRI signals, the LVC signals cannot
33 be guaranteed to contain all the useful information. Experiment 3.4 compares the performance of decoding semantic
34 information from different cortical areas, which shows that fMRI signals from HVC contain more reliable high-level
35 information (see Figure 3(a)). Therefore, the fMRI signals from HVC should be much helpful for the reconstruction.
36 In reconstruction, our approach did not explicitly model other features (e.g. color). However, part of the information
37 might be implicitly encoded in the weights of the GAN model during the training phase. In the future work, we will try
38 to explicitly decode more critical information for improvement.

39 **Why and how the images' classes were merged?** In the original dataset there are only 8 images within one category
40 (totally 150 categories). Training the DNN model with so few samples within one category is difficult. To facilitate
41 the model's training, the number of images in each category was increased by merging similar classes. We manually
42 grouped similar images from subdivide categories (e.g. gorilla and chimpanzee) into new coarser categories. We will
43 add a supplementary material to explain the details.

44 **Was the testing set withheld for the training of the decoders as well?** Yes. The test set was strictly withheld in the
45 training of the shape decoder, semantic decoder, and the GAN model. Thanks.

46 **About details of implementation.** We will include a supplemental material to describe the preprocessing of fMRI
47 data and network implementation details. And we will explain in detail what is not clear in the paper. For the other
48 suggestions/issues, we will revise the paper accordingly. Many thanks for your valuable comments.