

# A Multimodel Approach to the Algonauts Challenge

Noah Syrkis & Sophia De Spiegeleire

June 6, 2023

## Abstract

Understanding how the brain encodes visual information is a key challenge in neuroscience. In this project, we attempt to address this challenge by constructing a multimodal encoding model based on the Algonauts Project 2023 dataset. In addition to the dataset’s image modality, we incorporate a semantic feature vector that describes object categories contained in the image shown to the subject during the functional Magnetic Resonance Imaging (fMRI) data collection. We combine various linear modules to construct two models: one predicting the fMRI data from both the associated image and the image’s associated semantic feature vector; the other predicting both the fMRI data and the semantic vector from the image alone. Bayesian hyperparameter optimization suggests that the latter approach could potentially enhance model performance during inference without increasing the number of parameters. The model’s performance was evaluated using a 5-fold cross-validation strategy and the median Pearson correlation coefficient as the metric. The code for this project is accessible at [github.com/syrkis/neuroscope](https://github.com/syrkis/neuroscope), and training logs are available at [wandb.ai/syrkis/neuroscope](https://wandb.ai/syrkis/neuroscope).

## Introduction

Visual processing is the principal modality through which we interact and decipher our environment. Over the years, substantial progress has been made in understanding how the brain processes visual information, with even surprising parallels observed between artificial and biological vision processing [Cite]. However, as reality can exhibit extraordinarily different visual fingerprints—from simple geometric shapes to complex landscapes and visual noise—any system capable of visual perception is necessarily complicated. Fully capturing this complexity and intricacy remains a challenge. It is this challenge that is the focus of the 2023 Algonauts Project<sup>1</sup>. The Algonauts Project’s 2023 dataset is based on the Natural Scenes Dataset (NSD), which couples images from the Common Objects in Context (COCO) dataset (Lin2014?) with fMRI responses to

---

<sup>1</sup><http://algonauts.csail.mit.edu/>

those images from various participants.

Neuroimaging techniques like fMRI have facilitated valuable insights into the neural correlates of visual perception. However, the potential of these techniques has been somewhat constrained by computational model limitations and the expense and time required to collect large-scale fMRI datasets. Amid these challenges, deep learning has proven to be a powerful tool, that has facilitated a better understanding and emulation of human visual perception. Recent efforts to incorporate multimodality into deep learning models have opened promising avenues to bridge the gap between computational models and the brain’s complexity.

The experiment presented here explores how an additional modality might contribute to developing a model of the brain’s visual encoding system, *without* a large increase in complexity/parameter count. The additional modality used here is a semantic feature vector, derived from the COCO dataset, describing the object categories contained in each image. The two models we developed are tasked with 1) predicting the brain response given the image and knowledge of what is in the image, and 2) predicting the brain response and the semantic contents of the image.

## Literature review

### Visual information processing

Visual information processing, characterized by its hierarchical nature and intricate interconnectivity, plays a vital role in our understanding of the brain and perception. Traditionally, the process is categorized into low, mid, and high-level processing, focusing respectively on elementary visual features, their conjunctions, and abstract representations ((Gonzales-Casillas2018?), (Groen, Silson, and Baker 2017)). However, the hierarchical categorization is insufficient to capture the full complexity of real-world scene perception. It underrepresents the multimodal and interconnected nature of visual perception, particularly when processing complex stimuli such as natural scenes ((Allen et al. 2022), (Groen, Silson, and Baker 2017)). With the development of fMRI, it is now possible to explore and visualize real-time brain activity associated with visual perception ((Allen et al. 2022), (Haxby et al. 2001)). Despite this, the understanding of the intricate interconnectivity in visual perception, particularly for natural scenes, remains limited. Additionally, it is increasingly evident that to unravel the complex network underpinning visual perception, massive amounts of data are required ((Chang et al. 2019), (Allen et al. 2022)).

### Multi modality in visual perception and deep learning

From a human-centered standpoint, multimodality pertains to the multiple sensory systems through which humans perceive and interact with the world ((Parcalabescu, Trost, and Frank 2021)). It is a reflection of the brain’s capacity to integrate and process information from multiple sources. In the context of machine learning, multimodality

refers to utilizing multiple information sources to enhance algorithms' performance. Several studies have pointed out the benefits of multimodal learning in providing richer information about underlying data patterns and creating more complex feature representations ((Ngiam et al. 2011), (Gu et al. 2017)).

## Deep learning and visual encoding models

Deep learning has had a profound impact on neuroscience, particularly in understanding the brain's visual processing mechanisms. The success of deep learning models in neuroscience is attributed to their ability to process high volumes of data, their inherent flexibility, and their structure, which is inspired by the brain's own hierarchical organization ((kriegeskorte2015?), (kell2019?)).

In visual neuroscience, deep learning models have been extensively used to predict brain activity in response to visual stimuli, known as visual encoding models. A number of studies have demonstrated that deep neural networks (DNNs), particularly convolutional neural networks (CNNs), can accurately predict neural responses to various visual stimuli ((khaligh-razavi2014?), (yamins2014?), (cichy2016?)). Notably, these DNNbased models can even outperform traditional hand-engineered models ((guclu2015?), (kell2019?)).

However, most of these models rely on single-modal input data, typically the visual stimuli themselves. One study applied a convolutional recurrent neural network (CRNN) to investigate the computational mechanisms of the retinal circuit involved in interpreting natural scenes. The researchers discovered that recurrent spatiotemporal receptive fields of ganglion cells played a crucial role in encoding dynamic visual scenes. The findings also inciate that the inherent recurrence of the model enhanced the prediction the neural response, but also unveiled corresponding biological counterparts, emphasizing the power and potential of deep learning in visual neuroscience studies (Zheng et al. (2021)).

Zheng et al. (2021), applied a convolutional recurrent neural network (CRNN) to investigate the computational elements of the retinal circuit involved in interpreting the nature of natural scenes. Their findings highlight the instrumental role of the recurrent spatiotemporal receptive fields of ganglion cells in encoding dynamic visual scenes. The findings also inciate that the inherent recurrence of the model enhanced the prediction the neural response, but also unveiled corresponding biological counterparts, emphasizing the power and potential of deep learning in visual neuroscience studies.

Few studies have investigated the potential of multimodal deep learning for predicting fMRI responses. This approach could provide a more comprehensive understanding of visual perception and its underlying neural correlates, particularly when dealing with complex stimuli such as natural scenes.

Han et al. (2019) aimed to investigate the use of

## Methodology

Our methodology is that of a supervised machine learning experiment. We have access to preprocessed fMRI scans showcasing the blood oxygen level-dependent (BOLD) response to a variety of images. Our primary objective is to construct a multimodal model that has as many parameters during inference as its unimodal counterpart, and yet better predicts the brain’s response to a given image. This section outlines the steps and components involved in the execution of our experiment. It should be noted that, in accordance with the Algonauts Project, median Pearson correlation between voxels in the ground truth and the prediction is used as the target metric (though not as a loss function).

## Data

The data underpinning our experiment is provided by the Algonauts Project (Gifford et al. 2023), and is initially derived from the Natural Scenes Dataset (NSD) (Allen et al. 2022). The NSD is currently the largest dataset of its kind, encompassing cortical surface vertices from the left and right hemispheres of eight participants’ brains. These vertices correspond to the neurological responses triggered by 73,000 COCO images used by the NSD, each image depicting natural scenes. In addition to category information for each image, the COCO dataset provides other valuable metadata such as object location boxes and caption lists. Our experiment focused on the COCO object category information. The images in the NSD contains 80 different kinds of objects, with most images containing multiple object kinds (for example a horse and a person). As per the Algonauts guide<sup>2</sup>, we represented each image using the dimensionality reduction method principal component analysis (PCA) of the all the image’s activations in the 2012 image model Alexnet’s second layer (**krizhevsky2012?**). As in the Algonauts guide, PCA was performed reducing each image to a vector of size 100.

Over the course of a year, each participant in the NSD study was exposed to 10,000 unique images, with each image presented three times, resulting in 30,000 image trials per participant. The corresponding fMRI data comprises 19,004 and 20,544 voxels for the left and right hemispheres, respectively. These voxel counts were selected based on preprocessed, high-quality 7T fMRI responses measuring as BOLD response amplitudes. Also included in the dataset are region of interest (ROI) masks for each subject, which aid in extracting specific fMRI data from certain locations in the brain. The fMRI data has been mapped to Harvard’s FsAverage atlas such that the voxels are comparable across individuals. We eliminated subjects 6 and 8 from the experiment due to missing data (voxel counts differed from 19,004 and 20,544 for the left and right hemispheres respectively). We thus trained on six subjects.

---

<sup>2</sup>[https://colab.research.google.com/drive/1bLJGP3bAo\\_hAOwZPHpiSHKlt97X9xsUw](https://colab.research.google.com/drive/1bLJGP3bAo_hAOwZPHpiSHKlt97X9xsUw)

## Models

The purpose of our models is to infer the BOLD response from a given image. The architecture of our primary model involves taking a vector representation of an image  $x$ , and outputting a tuple consisting of the left hemisphere BOLD response  $y_{lh}$ , right hemisphere BOLD response  $y_{rh}$ , and a semantic feature vector  $y_c$  for optimization against the COCO data. This model is partitioned into four submodules, each an MLP processing one of the four variables ( $x$ ,  $y_{lh}$ ,  $y_{rh}$ , and  $y_c$ ). Our baseline will be a unimodal version of this model. We aim to test if including the semantic modality improves performance.

The first module, referred to as the image encoding module, maps the input image vector  $x$  onto a latent space, thereby generating a latent vector  $z$ , which is subsequently fed into the remaining three modules responsible for predicting the outputs. As suggested by the Algonauts challenge baseline, the latent vector  $z$  maintains a dimensionality of 100. Given that each hemisphere’s voxel count is approximately 20k, the linear mapping from the latent space to the voxel space demands around 2 million parameters. Therefore, even with such a compact latent space, the minimum required parameter count is approximately 4 million.

Our second model used  $y_c$  as input, concatenating it with  $x$ . The purpose of this model was to gauge the potential of multimodality on the input side of the network. This model is not our main focus, but rather a test to gauge the usefulness of this particular kind of multimodality.

All hidden layers used the tanh activation function, dropout of 0.1, and weight decay of 0.0001 with the AdamW optimizer from Optax. The models were implemented in Jax with Haiku(CITE). The shared (first) module had two layers, with 100 units each, to create some flexibility as the input to all other modules (the latent vector  $z$ ) flowed through that initial module. The rest of the modules mapped the latent vector input to whatever output dimension their modality had. The learning rate was 0.001 and the batch size was 32. Hyperparameter optimization was not done on the aforementioned hyperparameters due to computational constraints.

The primary model (with the auxiliary task of predicting  $y_c$  during training), had two experiment-specific hyperparameters,  $\alpha$  and  $\beta$ , weighing  $y_c$  and whatever hemisphere was not being optimized for respectively in the loss function. The model used mean squared error for optimizing the fMRI predictions and binary soft f1 loss for  $y_c$  due to a heavy imbalance between categories. Using regular binary cross entropy would yield a low loss by guessing all zeros, as most images contain only a few categories.

## Experiments

#### Incorporating Category Vector Modality and Semantic Vector Representation

To unlock the potential utility of the semantic vector, we designed our experiment with

a multimodal approach. This involved integrating the category vector modality (model 2) by concatenating it with the image vector derived from AlexNet, an auxiliary task to predict the category during training (model 1), and tuning the  $\alpha$  and  $\beta$  parameters weighting the importance of the auxiliary tasks in the loss function. Additional motivation for the inclusion of the auxiliary modalities is the potential avoidance of overfitting; finding inappropriate shortcuts in the data becomes more difficult if the shortcuts also have to make sense of the semantic vector.

**Model Training, Auxiliary Tasks, and Hemisphere Balancing** Two key hyperparameters,  $\alpha$  and  $\beta$ , were used to balance the different aspects of our model’s performance.  $\alpha$  controlled the balance between fMRI loss and category prediction loss, thereby providing weight to the auxiliary task of category prediction. This strategy was based on our hypothesis that having the model solve an auxiliary classification problem could lead to more generalized and versatile representations beneficial for the primary task of predicting fMRI responses.  $\beta$  modulated the balance between the losses of the two hemispheres. By tuning this parameter, we hoped to find out if there is balance that might contribute to a better overall model performance on the subjects.

**Hyperparameter Optimization and Loss Function Design** The cornerstone of our experiment involves hyperparameter optimization, carried out using the Weights & Biases (wandb) sweeps with wandb’s Bayesian optimization techniques. The loss function is expressed as  $(1 - \alpha)((1 - \beta)Loss_{y_{lh}} + \beta Loss_{y_{rh}}) + \alpha Loss_{y_c}$ , when optimizing for  $y_{lh}$  and flipping the  $\beta$  when optimizing for  $y_{rh}$ .  $\alpha$  serves as a weighting factor determining the trade-off between the fMRI prediction task and the category prediction task, while  $\beta$  controls the balance between the losses of the two hemispheres.

**Bayesian Optimization and Cross-Validation** To search for the optimal values of  $\alpha$  and  $\beta$ , we initiated a wandb sweep with Bayesian optimization and optimized with respect to validation left hemisphere correlation in one sweep, and validation right hemisphere correlation in the other sweep. This strategy enables a directed search in hyperparameter space, making it a more efficient and effective approach for hyperparameter tuning than random search or grid search. Additionally, we employed a K-fold cross-validation technique for model evaluation, providing a more robust estimate of the model’s performance and optimal hyperparameters. K was set to 5. Every fold for every subject ran twice to get samples during the Bayesian optimization.

## Results

In **table 1** we see the mean median voxel correlations for versions of model 1 (the primary model with auxiliary task) trained with and without  $\alpha$  and  $\beta$  set to 0. To reiterate: the baseline is model 1 those hyperparameters set to 0, making it unimodal.

Table 1: Mean Median Voxel Correlation (Model 1).

Hemisphere	Train, Alex/COCO	Train, Alex	Test, Alex/COCO	Test, Alex
Left	0.2176	0.2059	0.1959	0.1957
Right	0.2155	0.2046	0.1933	0.1917

In **table 2** we see the mean median voxel correlations across all subjects and folds of model 2 with (Alex + COCO) and without (Alex) the COCO vector concatenated to the Alex vector.

Table 2: Mean Median Voxel Correlation (Model 2).

Hemisphere	Train, Alex/COCO	Train, Alex	Test, Alex/COCO	Test, Alex
Left	<i>0.2176</i>	0.2059	<i>0.1932</i>	0.1927
Right	<i>0.2155</i>	0.2046	<i>0.195</i>	0.1908

In **table 3** we see results of the wandb hyper parameter sweep.

Table 3: Bayesian Hyperparamter Sweep (Model 1).

Hemisphere	$\alpha$ correlation	$\alpha$ importance	$\beta$ correlation	$\beta$ importance
Left	0.063	0.424	- 0.147	0.576
Right	0.076	0.566	- 0.087	0.434

A mean (across all subjects and folds) median voxel correlation projection onto a common cortical atlas is available interactively at [neuroscope.streamlit.app/](https://neuroscope.streamlit.app/).

## Analysis and Discussion

### Future Work

As seen in the Analysis and Discussion, it appears that the semantic vector modality is not particularly useful for the model. A logical next step would be to experiment with extracting the image representations from different, or multiple AlexNet layers, or using an entirely different model for the image representation extraction. We might also explore using more rich COCO modalities such as image captions and object bounding boxes. Lastly, from a neuroscientific perspective, the ROIs of the brain are considered to be different modalities: they function by vastly different rules. Processing the ROIs separately might allow for models tailoring to specific ROI idiosyncrasies.

## Conclusion

## References

- Allen, Emily J., Ghislain St-Yves, Yihan Wu, Jesse L. Breedlove, Jacob S. Prince, Logan T. Dowdle, Matthias Nau, et al. 2022. “A Massive 7T fMRI Dataset to Bridge Cognitive Neuroscience and Artificial Intelligence.” *Nature Neuroscience* 25 (1, 1): 116–26. <https://doi.org/10.1038/s41593-021-00962-x>.
- Gifford, A. T., B. Lahner, S. Saba-Sadiya, M. G. Vilas, A. Lascelles, A. Oliva, K. Kay, G. Roig, and R. M. Cichy. 2023. “The Algonauts Project 2023 Challenge: How the Human Brain Makes Sense of Natural Scenes.” January 10, 2023. <http://arxiv.org/abs/2301.03198>.
- Han, Kuan, Haiguang Wen, Junxing Shi, Kun-Han Lu, Yizhen Zhang, Di Fu, and Zhongming Liu. 2019. “Variational Autoencoder: An Unsupervised Model for Encoding and Decoding fMRI Activity in Visual Cortex.” *NeuroImage* 198 (September): 125–36. <https://doi.org/10.1016/j.neuroimage.2019.05.039>.
- Zheng, Yajing, Shanshan Jia, Zhaofei Yu, Jian K. Liu, and Tiejun Huang. 2021. “Unraveling Neural Coding of Dynamic Natural Visual Scenes via Convolutional Recurrent Neural Networks.” *Patterns* 2 (10): 100350. <https://doi.org/10.1016/j.patter.2021.100350>.

## Appendix