

## **What is a good model for brain encoding in a videogame task ?**

**François Paugam (francois.paugam@umontreal.ca)**

Centre de Recherche de l'Institut Universitaire de Gériatrie de Montréal, Montréal  
Mila - Quebec AI Institute, Montréal  
Université de Montréal, Montréal

**Guillaume Lajoie**

Mila - Quebec AI Institute, Montréal  
Université de Montréal, Montréal

**Pierre Bellec (pierre-louis.bellec@umontreal.ca)**

Centre de Recherche de l'Institut Universitaire de Gériatrie de Montréal, Montréal  
Université de Montréal, Montréal

## Abstract

Videogames represent a promising experimental paradigm for neuroscientists to study active tasks in complex environments. However, interpreting brain dynamics in such complex environments is challenging, though a recent approach is to use brain encoding, i.e. quantify the similarities in activity between the brain and an artificial neural network. A wide range of modelling approaches could potentially be used to encode brain activity in videogames. In this work we compare three machine learning models trained with different objective functions to encode fMRI data collected on 5 subjects playing Super Mario Bros: (1) PPO was trained with reinforcement learning to play the game from video frames; (2) VideoGPT was trained through predictive coding on videos of human gameplay; (3) ResNet was trained for image classification in a diverse set of natural images. All three models produced qualitatively similar brain encoding maps on the levels used for training, though overall ResNet had better brain encoding accuracy and generalised better to new levels. As VideoGPT and PPO were trained from scratch on videogame data, they demonstrate the feasibility of future experiments to explain brain activity during videogames while carefully controlling the nature and size of data used for training.

**Keywords:** brain encoding, artificial neural networks, fMRI, neuroimaging, videogames

## Introduction

Videogames offer a promising experimental framework to study cognitive processes, by strongly engaging participant's attention, emotions, as well as motor and decision making abilities (Anderson et al., 2011; Bavelier, Achtman, Mani, & Föcker, 2012). Brain encoding has emerged as a powerful method to study brain representations evoked by rich environment such as videogames, by quantifying the similarities between the representations of the brain and artificial neural networks. Brain encoding has been widely studied in vision and language tasks (Schrimpf et al., 2020; Oota, Arora, Rowtula, Gupta, & Bapi, 2022; Conwell, Prince, Kay, Alvarez, & Konkle, 2023). By contrast, very few works have used brain encoding techniques in the context of videogames (Cross, Cockburn, Yue, & O'Doherty, 2021; Kemtur et al., 2023). We aim to assess the performance of three different approaches to train a brain encoding model on a fMRI videogame task using the unprecedented CNeuroMod dataset (Boyle et al., 2020). These approaches correspond to markedly different choices of objective functions, i.e. what is optimized by the model during training.

The first approach, reinforcement learning (RL), consists of training a model to play a videogame while maximizing a reward (Cross et al., 2021). The second, "predictive coding" employs an unsupervised generative model trained to generate the continuation of a given input clip of human gameplay. The third, "classification", consists of a vision model trained

to perform classification of a diverse collection of natural still images into many categories. We selected a specific implementation for each objective function, and compared their ability to encode brain activity across different levels of the game "Super Mario Bros" (SMB).

## Methods

### Dataset

We used the Mario dataset of the Courtois NeuroMod databank (Boyle et al., 2020). The dataset (n=5) comprises of about 10h of fMRI data per subject, while they played SMB. fMRI data were acquired on a 3T Siemens Prisma Fit scanner (TR=1.49s, 2mm isotropic), preprocessed using the fMRIprep pipeline (Esteban et al., 2019), and projected onto MIST atlas, at the 1097 parcels scale (Urchs et al., 2019). We left out data from 2 of the 22 SMB levels available in the dataset, to evaluate generalisation to new levels and used the remaining 20 to fit the brain encoding regressions.

### Models

**Reinforcement learning: PPO** The PPO model is a convolutional neural network (CNN) trained to play the game by maximising a reward with the Proximal Policy Optimisation algorithm directly from pixel-level video frames of the game (Schulman, Wolski, Dhariwal, Radford, & Klimov, 2017). The same model is trained to play on all the 20 training levels. It was trained for about 380h of in-game time (~ 20 million images).

**Predictive coding: VideoGPT** The VideoGPT architecture (Yan, Zhang, Abbeel, & Srinivas, 2021) directly processes the video stream of subjects' gameplay, using a short video as input (16 frames at 12 Hz, 1.33 sec). VideoGPT has two modules: a VQ-VAE which compresses the input video in a discrete latent space and a GPT-like transformer, which predicts the VQ-VAE embeddings of the next video, given the VQ-VAE embeddings of the current video. A single VQ-VAE was trained on the pooled data of all subjects and all 22 levels, but we trained a separate transformer for each subject, using one 80%/10% random split of gameplay video across the 20 levels for training/validation per subject (~ 400,000 images, 9h of in-game time).

**Classification: ResNet** ResNet is the pretrained resnet152v2 model (He, Zhang, Ren, & Sun, 2016). It is a CNN with residual connections, trained on a classification task on the ImageNet1K dataset (1.28 million images, 1,000 categories). At the time of writing it is ranked 10<sup>th</sup> on the BrainScore vision leaderboard (Schrimpf et al., 2020).

### Brain encoding

We trained a brain encoding linear readout for each model and subject separately, using scikit-learn (Pedregosa et al., 2011). A 80/10/10 % split of each subjects' data from the 20 levels was used for training, validation and test respectively for all models, with random splits stratified by level identical to the

VideoGPT training. For each model, each layer, and the data of each subject, we extracted the activations for input frames of the 20 levels corresponding to 4.5s, 6s and 7.5s delays before each bold volume. These activations were averaged over time steps, and a PCA was applied to keep 1,000 principle components. The PCA ensured that latent spaces of identical dimensionality were used across models. A ridge regression was fitted and the validation score was used to select the optimal layer and regularisation parameter. For each brain parcel, the test and generalization scores are the  $R^2$  regression scores (averaged across subjects), in the test set and the 2 left-out levels, respectively.

## Results

For all models, the best layer for brain encoding was among the last layers, responsible for decoding latent variables. On average ResNet had the highest brain encoding  $R^2$  test score across brain parcels (0.059), followed by PPO (0.050) and VideoGPT (0.048).  $R^2$  maps of all models had similar topography (Fig.1), with spatial Pearson correlation between maps exceeding 0.98 for all model pairs. The best encoded regions were visual and somatosensory cortices. The ResNet encoded significantly better the visual and sensory regions versus other two models, while PPO and VideoGPT had almost no significant difference in  $R^2$  maps (Fig. 2). Performance of brain encoding was degraded on left-out levels, but followed a similar topography as the original 20 levels for ResNet and PPO, while performance of VideoGPT fell to chance (Fig. 3).

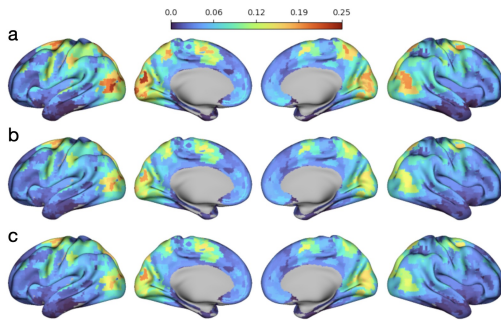


Figure 1: Test  $R^2$  maps of the models, averaged over subjects. **a** ResNet **b** PPO **c** VideoGPT

## Discussion

Our results demonstrate it is *possible* to encode brain data in videogames using models based on different types of objective functions. Surprisingly, the difference in objective functions did not lead to a marked advantage to one of the models for encoding specific brain areas. The ResNet still encoded brain activity better than VideoGPT and PPO. The advantage was moderate for the levels used to train VideoGPT and PPO, and large for left-out levels. This result is likely attributable to its more diverse training set, which made it learn a richer set of features, suggesting that the encoded brain activity is mainly

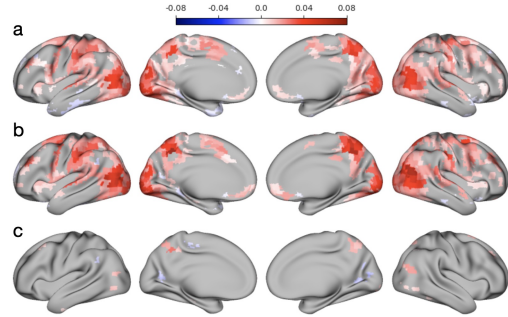


Figure 2: Difference of the test  $R^2$  maps of the models. Shown ROIs are ROIs where the difference is significant ( $p < 0.05$ ) according to a two-sided Wilcoxon signed-rank test, corrected for false discovery rate. **a** ResNet minus PPO, positive (red) values are where the ResNet encodes better **b** ResNet minus VideoGPT **c** PPO minus VideoGPT

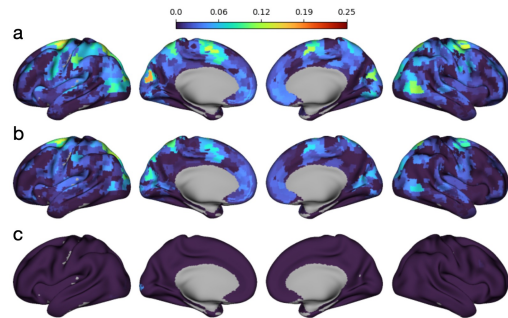


Figure 3: Generalization  $R^2$  maps of the models, averaged over subjects. The  $R^2$  values are clipped from 0. **a** ResNet **b** PPO **c** VideoGPT, almost all values are negative.

driven by high-level visual features. This result is inconsistent with the study by Cross et al. (2021) who concluded to the superiority of RL to vision features. But, this prior work had only investigated a smaller vision model trained on the videogame frames instead of a model like ResNet trained on diverse images. This last point emphasizes that conclusions from this study must be taken carefully, as the three models differed on various aspects other than their objective function, in particular their training dataset and architecture. Given that PPO and VideoGPT were trained on videogame data and still managed to achieve similar brain encoding performance as ResNet, it opens avenues for new experiments to explain brain activity during videogames. For example, they could be re-trained to rigorously establish the potential advantages of different objective function (e.g. different reward function in RL) or architecture (e.g. adding two separate vision branches as in (Bakhtiari, Mineault, Lillicrap, Pack, & Richards, 2021)) in terms explaining brain activity in different brain regions, while carefully controlling for other factors such as dataset size and diversity.

## Acknowledgments

This research was funded by the Courtois foundation awarded to PB. Dr. Lajoie acknowledges support from the Canada CIFAR AI Chair program and the Canada Research Chair in Neural Computations and Interfacing.

## References

- Anderson, J. R., Bothell, D., Fincham, J. M., Anderson, A. R., Poole, B., & Qin, Y. (2011). Brain regions engaged by part-and whole-task performance in a video game: a model-based test of the decomposition hypothesis. *Journal of cognitive neuroscience*, 23(12), 3983–3997.
- Bakhtiari, S., Mineault, P., Lillicrap, T., Pack, C., & Richards, B. (2021). The functional specialization of visual cortex emerges from training parallel pathways with self-supervised predictive learning. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P. S. Liang, & J. W. Vaughan (Eds.), *Advances in Neural Information Processing Systems* (Vol. 34, pp. 25164–25178). Curran Associates, Inc. Retrieved from [https://proceedings.neurips.cc/paper\\_files/paper/2021/file/d384dec9f5f7a64a36b5c8f03b8a6d92-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2021/file/d384dec9f5f7a64a36b5c8f03b8a6d92-Paper.pdf)
- Bavelier, D., Achtman, R. L., Mani, M., & Föcker, J. (2012). Neural bases of selective attention in action video game players. *Vision research*, 61, 132–143.
- Boyle, J., Pinsard, B., Boukhdir, A., Belleville, S., Brambatti, S., Chen, J., ... Bellec, P. (2020). The courtois project on neuronal modelling – first data release. 26th annual meeting of the organization for human brain mapping, 2020.
- Conwell, C., Prince, J. S., Kay, K. N., Alvarez, G. A., & Konkle, T. (2023). What can 1.8 billion regressions tell us about the pressures shaping high-level visual representation in brains and machines? *bioRxiv*. Retrieved from <https://www.biorxiv.org/content/early/2023/07/01/2022.03.28.485868> doi: 10.1101/2022.03.28.485868
- Cross, L., Cockburn, J., Yue, Y., & O’Doherty, J. P. (2021). Using deep reinforcement learning to reveal how the brain encodes abstract state-space representations in high-dimensional environments. *Neuron*, 109(4), 724–738.
- Esteban, O., Markiewicz, C. J., Blair, R. W., Moodie, C. A., Isik, A. I., Erramuzpe, A., ... others (2019). fmriprep: a robust preprocessing pipeline for functional mri. *Nature methods*, 16(1), 111–116.
- He, K., Zhang, X., Ren, S., & Sun, J. (2016, June). Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition (cvpr)*.
- Kemtur, A., Paugam, F., Pinsard, B., Sainath, P., Clei, M. L., Boyle, J., ... Bellec, P. (2023). Behavioral imitation with artificial neural networks leads to personalized models of brain dynamics during videogame play. *bioRxiv*. Retrieved from <https://www.biorxiv.org/content/early/2023/11/01/2023.10.28.564546> doi: 10.1101/2023.10.28.564546
- Oota, S. R., Arora, J., Rowtula, V., Gupta, M., & Bapi, R. S. (2022, October). Visio-linguistic brain encoding. In N. Calzolari et al. (Eds.), *Proceedings of the 29th international conference on computational linguistics* (pp. 116–133). Gyeongju, Republic of Korea: International Committee on Computational Linguistics. Retrieved from <https://aclanthology.org/2022.coling-1.11>
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., ... Duchesnay, E. (2011). Scikitlearn: Machine learning in Python. *Journal of Machine Learning Research*, 12, 2825–2830.
- Schrimpf, M., Kubilius, J., Lee, M. J., Murty, N. A. R., Ajemian, R., & DiCarlo, J. J. (2020). Integrative benchmarking to advance neurally mechanistic models of human intelligence. *Neuron*. Retrieved from [https://www.cell.com/neuron/fulltext/S0896-6273\(20\)30605-X](https://www.cell.com/neuron/fulltext/S0896-6273(20)30605-X)
- Schulman, J., Wolski, F., Dhariwal, P., Radford, A., & Klimov, O. (2017). Proximal policy optimization algorithms. *CoRR*, abs/1707.06347. Retrieved from <http://arxiv.org/abs/1707.06347>
- Urchs, S., Armoza, J., Moreau, C., Benhajali, Y., St-Aubin, J., Orban, P., & Bellec, P. (2019). Mist: A multi-resolution parcellation of functional brain networks. *MNI Open Research*, 1, 3.
- Yan, W., Zhang, Y., Abbeel, P., & Srinivas, A. (2021). Videogpt: Video generation using VQ-VAE and transformers. *CoRR*, abs/2104.10157. Retrieved from <https://arxiv.org/abs/2104.10157>